

Article

Multiformalism models for performance engineering

Enrico Barbierato¹, Marco Gribaudo² and Giuseppe Serazzi³

¹ Dip. di Matematica e Fisica, Università Cattolica del Sacro Cuore, Via dei Musei 41, 25121 Brescia, Italy; enrico.barbierato@unicatt.it

² Dip. di Elettronica, Informazione e Bioingegneria, Politecnico di Milano, via Ponzio 345, 20133 Milano, Italy; marco.gribaudo@polimi.it

³ Dip. di Elettronica, Informazione e Bioingegneria, Politecnico di Milano, via Ponzio 345, 20133 Milano, Italy; giuseppe.serazzi@polimi.it

* Correspondence: enrico.barbierato@unicatt.it

Version March 5, 2020 submitted to Future Internet

Abstract: Nowadays, the necessity to predict the performance of cloud and edge computing-based architectures has become paramount, in order to respond to the pressure of data growth and more aggressive Service Level of Agreements. In this respect, the problem can be analyzed by creating a model of a given system and studying the performance indices values generated by the model's simulation. This process requires to take in account a set of paradigms, carefully balancing the benefits and the disadvantages of each one. While Queuing Networks are particularly suited to modeling cloud and edge computing architectures, however, particular occurrences - such as autoscaling - require different techniques to be analyzed. This work presents a review of paradigms designed to model specific events in different scenarios, such as Timeout with Quorum-based Join, Approximate Computing with Finite Capacity region, MapReduce with Class Switch, Dynamic Provisioning in Hybrid Clouds and Batching of requests in e-Health applications. The case studies are investigated by implementing models based on the above mentioned paradigms and analyzed with discrete event simulation techniques.

Keywords: Quorum-based Join; Multiformalism; Finite Capacity region; Class Switch

1. Introduction

In 2017, the Economist designated the growing diffusion of data as the "the new oil" emergency [1], quoting the profit reached in the first quarter of the year generated by the main market actors equal to 25 billion dollars. Though the remark was not exactly original (the same perspective had been already shared by the mathematician Clive Humby in 2006 [2]) and was far from being an obvious analogy [3], it highlighted how applications and services could be characterized in order to maximize benefit from this new commodity: i) network access (the capability to grant access to the internet to a broad range of devices), ii) service metering (a user is billed accordingly to how much service it is consumed), iii) shared architecture (the capability of providing shared resources and infrastructures to all sort of users) and iv) provisioning according to demand (it is necessary to provide dynamically the requested services). All these requirements converged in a model called *cloud computing*. In Gartner's 2018 Hype Cycle for Cloud Computing¹, it is reported that "Cloud computing has reached the Slope of Enlightenment" and Forbes announced it is "the new kid in the block"², 5G giving a boost to cloud

¹ <https://www.gartner.com/en/documents/3884671>

² <https://www.forbes.com/sites/forbestechcouncil/2019/11/15/the-next-evolutionary-step-for-cloud-computing/#3e04c1646dd7>

28 technology in different market areas, such as autonomous cars, virtual reality, health care and smart
29 homes.

30 It is more and more evident that the capacity to offer services such as servers, databases and
31 storage has become a powerful attractor for the main tech players in the world market, including
32 Amazon, Microsoft, Google, IBM and Oracle to name a few. An additional aspect concerns security,
33 chiefly the ability to shield data in the cloud from any kind of malicious intrusion (ranging from
34 hacking to theft) by properly configuring firewalls and VPNs initially and establishing specific policies.

35 Moreover, to guarantee an efficient usage of the cloud architecture, it is necessary to define a
36 Quality of Service (QoS) policy. In [4], quality is defined as consideration of a few key indicators
37 such as i) flexibility (managing a functionality without affecting the system), ii) maintainability, iii)
38 performance, and iv) scalability, among others. The authors stress the level of difficulty hidden in
39 the choice of the process able to provision a valid QoS agreement for Cloud Computing architectures,
40 identifying scheduling, admission control and dynamic resource provisioning the main keys to solving
41 to this problem.

42 User expectations drive the design of the cloud model, which is articulated through four pillars: i)
43 public (a kind of cloud available to a large audience, which deploys services at a low cost), ii) private
44 (the access is restricted to the members of a particular organization), iii) community (access is shared
45 by organizations sharing analog motivations) and finally iv) hybrid.

46 Such pillars present both advantages and disadvantages. Among the former, the most evident is
47 the cost cutting in IT companies infrastructures regarding implementation, scalability, maintenance
48 and reliability. Having fluid access to an ubiquitous cloud represents another significant landmark.
49 The most interesting challenge consists of security and privacy matters, followed by some ambiguities
50 in defining the services (this problem has been partially mitigated by the Open Cloud Consortium).
51 If cloud computing relies on a centralized architecture (usually, a data center), in edge computing,
52 processing occurs at the edge of the network. This choice has been proved to be a viable solution to
53 overcome data latency typical of cloud computing, at the cost of limited processing power.

54 The two proposed architectures are not mutually exclusive: instead, they complement each other. In
55 this sense, *fog computing* extend the concept of centralized network by taking into account localized
56 data centers or *fog nodes*, deployed in order to store and elaborate data at a shorter distance from the
57 source. A fog node can filter which data need to be referred to the central server of the cloud structure
58 from that which can be processed locally.

59 As the success of a cloud network architecture depends on its performance, it is crucial to identify
60 those factors influencing its implementation in order to build a valid model, which can be either
61 simulated or studied analytically. The scientific literature provides numerous approaches to this task,
62 ranging from stochastic models to machine learning-inspired methods.

63 The contribution of this work consists of a survey of specific extensions of Queuing Networks
64 formalism and solving methods to efficiently analyze the aforementioned scenarios. To be more specific,
65 we start from the observation that Queuing Networks are particularly suited to model cloud and edge
66 computing architectures, as this formalism denotes some limitations when taking in account different
67 scenarios. For instance, one can imagine a simple auto-scaling scenario, where virtual machines are
68 shut down and restated only when needed³. Restart requires a non-negligible amount of time, and
69 when the system is fed with a very low workload, this becomes the main component of the average
70 response time. As the workload increases, the chance of finding the system already in *on-state* increases
71 too, resulting in a reduced response time. However, then requests increase, resources start to saturate,
72 extending again the average response time. This behavior cannot be modeled easily with queuing
73 modeling primitives even if some specific computing techniques (such as *Fork/Join*, *Finite Capacity*

³ See, for example, *Overview of autoscale in Microsoft Azure Virtual Machines, Cloud Services, and Web Apps* at <https://docs.microsoft.com/en-us/azure/azure-monitor/platform/autoscale-overview>

74 *Regions* and *Class Switch*) and approaches (such as *Multiformalism*) allow the modeler to take in account
75 more complex architectures.

76 Table 1 briefly recaps the main characteristics of the case studies presented here, including specific
77 characteristics for which special techniques are required and describes the minimal extension that can
78 capture such behavior.

Table 1. Case studies summary

Scenario	Solution	Section
Timeout with Quorum based Join	Fork/Join paradigm	3.1
Approximate Computing with Finite Capacity Region	Finite capacity regions	3.2
MapReduce with Class Switch	Class Switch	3.3
Dynamic provisioning in Hybrid Clouds	Multiformalism	4.1
Batching of requests in e-Health applications	Multiformalism	4.2

79 All the models presented in this work have been solved using JMT⁴ [5]. All the results have
80 been computed with 99% confidence intervals that, for simplicity of presentation, are not graphically
81 represented.

82 This work is organized as follows. Section 2 takes in account the related work in the scientific
83 literature and Section 3 focuses on computing problems that can be addressed with extensions to
84 classical Queuing Networks modelling features, while Section 4 describes cases where this formalism
85 alone cannot be applied, but where multi-formalism modelling techniques that combine Petri Nets with
86 Queuing Networks, can provide meaningful and clear solutions. Both sections include a description of
87 some of the known results, explaining how the contribution was able to provide better insights in the
88 considered scenarios. Finally, Section 5 draws the conclusions from the paper.

89 2. Related Work

90 The complex scenarios introduced by cloud computing needs to be analyzed in order to be able
91 to predict the future of the current technological landscape. To this respect, Varghese and others
92 discuss in [6] some of the current paradigms, including i) the cloud pitfalls, ii) hybrid paradigms, iii)
93 micro-clouds and cloudlets, iv) ad-hoc clouds (such as SETI), v) heterogeneous clouds (at high level,
94 considering multiple providers and at low level, from the perspective of different processors), vi) fog
95 and mobile edge computing and finally vii) serverless computing. The remaining part of the article
96 examines the impact on society of next generation cloud paradigms.

97 In reviewing the existing work on modeling and simulating cloud and edge computing
98 technologies, various aspects need to be included. For instance, an important guideline to driving
99 a system performance analysis consists of determine which metrics characterize the considered
100 architecture.

101 In [7], the authors examine a specific class of services (Infrastructure as a Service or IaaS) discussing
102 Application Response Time (ART, literally "time taken by the application to respond to other users'
103 requests.") as a metric. Specifically, the accounted tasks are described according from two perspectives,
104 computation intensive and communication intensive tasks: the former is decomposed in the analysis of
105 CPU and memory consumption, the latter being monitored by network tools (or even by using SMNP
106 agents). In [8], Qiang Duan considers an extended set of parameters determining the performance
107 of a cloud architecture service, such as i) response time, ii) throughput, iii) availability, iv) utilization,
108 v) resilience, vi) scalability and vii) elasticity. Hybrid architecture performance is the focus of [9].
109 Maheshwari and others review the paramount design parameters such as the proportion of resources

⁴ <http://jmt.sourceforge.net/Download.html>

110 on edge side vs cloud side and the latency of edge clouds in order to determine measure indices in the
111 shape of the average response time and service goodput.

112 With respect to the study of computing infrastructures workload in general (including the
113 identification of patterns), Calzarossa et al. present [10], an interesting survey on the subject. Workload
114 remains a key indicator to correctly match the Quality of Service (QoS) and Quality of Experience
115 (QoE) and to adequately respond to energy saving policies and resource provisioning requests. The
116 capacity of performing passive or active data collection (by storing data with logging capacities or
117 by deploying ad hoc tools) from an application generates a volume of data that can be exploited for
118 monitoring purposes. However, such tasks imply a risk, since collecting big volumes of data can add
119 a significant workload to the system, which can be mitigated by choosing appropriated sampling
120 techniques, assuming that the sample is correctly selected (this and other issues are discussed in [11]).
121 The work in [10] reviews also the main techniques used to analyze the system performance, ranging
122 from statistical analysis to the usage of graphs and finally stochastic processes.

123 Multiformalism is used to model different components of a system whereby the modeler's
124 aim is to compute its performance, through different formalisms. The literature presents several
125 methodologies (see [12] for an overview of historical evolution of the field and especially for
126 what concerns performance modeling techniques through QN and PNs), each one supported by
127 a corresponding tool. For example, SMART [13] is a software package used to design complex
128 discrete-state systems, providing both numerical solution algorithms and discrete-event simulation
129 techniques. PEPA (Performance Evaluation Process Algebra, [14]) is one of the many extensions of
130 Process Algebra (a set of abstract languages capable of describing concurrent systems consisting of a
131 set of agents performing one or more actions specifying concurrent behaviors and the synchronization
132 between them. Typical examples of process algebras are Communicating Sequential Processes (CSP,
133 [15]) and Calculus of Communicating Systems (CCS, [16]). Its novelty consists in the deployment
134 of the concept of duration (an exponentially distributed random variable) of an action that makes
135 explicit the relationship between the Process Algebra model and a Continuous Time Markov Chain.
136 Different components of a system work together by using a kind of cooperation technique. More recent
137 approaches include Möbius [17], OsMoSys [18] and SIMTHESys [19]. Users who need to design new
138 heterogeneous formalisms on Möbius are requested to refer to a meta-model interface called Abstract
139 Functional Interface (AFI). Möbius supports Stochastic Activity Networks (SANs), Petri nets and
140 Markov chains. OsMoSys can create multiformalism models and workflow management to achieve
141 multi-solution. The key idea of OsMoSys relies on meta-modeling and the concept of object-oriented
142 paradigms.

143 From the point of view of model structure, OsMoSys represents a main metaformalism
144 (metametamodel), supports formalism inheritance (at formalism and element level), and can extend
145 formalisms by adding new elements. It allows model composition by the inclusion of submodels,
146 supporting generic submodels and hidden information and multiformalism models with bridge
147 formalisms. Möbius presents a more complex model architecture, where several different model types
148 (organized in a logic tree, parameterize and solve a model. OsMoSys supports the development of
149 multiformalism models by composition of submodels written in different formalisms by exploiting the
150 benefits of metamodeling.

151 Multisolution is dealt with Möbius and OsMoSys in a different way. In the former, the solver is
152 obtained in the form of an optimized executable model, based on the description given by the user. The
153 latter solves models by (semiautomatically) generating a business process, executed by its workflow
154 engine, which describes the solution in terms of external solvers activations.

155 SIMTHESys (Structured Infrastructure for Multiformalism modeling and Testing of Heterogeneous
156 formalisms and Extensions for SYStems, [20]) is a framework for defining new formalisms and
157 generating other related solvers, that allows the combination of more formalisms in the same models.
158 The solution architecture of SIMTHESys is designed to automatic generating solvers based on several

159 formalism families, Exponential Events, Exponential and Immediate Events, Labeled Exponential
160 Events Formalisms and so forth).

161 In regard to surveys on cloud and edge computing, the literature offers several examples: [21]
162 is an exhaustive analysis of the main paradigms, while in [22], the authors classify different fog
163 computing-based systems, studying their principal requirements and features. In [23] Mao and others
164 review a selection of papers on mobile edge computing by using Task model, Design Objective and
165 Proposed Solution as guidance.

166 With respect to other paradigms than those reviewed in this paper, Machine Learning algorithms
167 have achieved a considerable popularity to predict the performance of cloud architectures. In [24],
168 the authors discuss issues related to scaling of VMs resources in cloud computing implementing
169 proactive strategies based on Neural Networks, Linear Regression and Support Vector Regression,
170 the latter providing the best accuracy.

171 In [25], Ardagna and others evaluate a queuing-based analytical model and a novel fast ad-hoc
172 simulator in various scenarios and infrastructure setups. Such approaches are able to predict
173 average application execution times with an error of 12% on average. Machine Learning and
174 analytical modeling can be combined as discussed in [26], where different hybrid applications, such as
175 Transactional Auto Scaler, IRON Model and Chorus, are based respectively on divide and conquer,
176 bootstrapping and ensemble techniques.

177 3. Queuing Networks techniques

178 The focus of this section concerns cloud related features that need advanced Queuing Network
179 techniques to be properly modeled. Firstly, this section presents an approach for timeout modelling,
180 which is followed by a study of its impact on cloud applications (Section 3.1). Secondly, we propose an
181 example of Approximate Computing, exploited to provide complex services in an environment with
182 highly variable workload. Note that the reduction of the solutions quality is used to provide service
183 level agreements guarantees while posing a limit to the automatic scalability of a cloud deployment
184 (Section 3.2). Finally, Big Data techniques, exploiting the parallelization of resources deployed on a
185 large number of virtual machines and using paradigms such as MapReduce (Section 3.3), conclude
186 this section.

187 3.1. Timeout with Quorum based Join

188 The implementation of a timeout value can be useful in several situations, for example to
189 understand when a particular computing process should be arrested before exhausting the system
190 resources or because its execution is taking so long that is preventing the execution of other components.
191 In cloud computing it can be used to model either user behavior, as detailed in Section 3.1.1, or specific
192 cloud features. In the latter case, two notable examples are Spot Instances on Amazon Web Services,
193 and maximum execution times in Function as a Service (FaaS) Serverless deployments.

194 Spot instances are very cheap virtual machines that are provided with a bidding process. The user
195 places a bid, defining the maximum price she is willing to pay for running the VM. The provider -
196 based on its current workload - and the offered price decide whether to assign the VM or not. If the
197 user receives a VM, this might be later shut down by the provider, in case of changes in the workload
198 that would increase the price of the VMs to an amount larger than the user bid. This event, can be
199 modelled as a timeout that might occur after a random amount of time.

200 Serverless computing, is currently a very popular cloud software deployment paradigm: users writes
201 their applications as a set of functions that exploits other cloud services (such as authentication and
202 storage), and that is written in a specific programming language. The cloud provider decides where
203 such functions will be run, based on the users' requests. Due to small overhead, functions can be
204 started and parallelized in a relatively quick way, reducing the users costs and increasing the providers
205 resource utilization. However, to keep a balance in the overall architecture, all functions executions are
206 characterized by a maximum running time. Should they take longer, they will be interrupted, and all

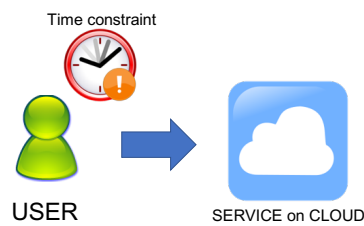


Figure 1. An application scenario including user-generated timeouts

207 the partial computation will be lost. This type of behavior can be easily modelled with a timeout event
 208 of constant duration.

209 3.1.1. Description of the problem

210 Let us consider the case study presented in Figure 1. In this scenario, a web application is deployed
 211 over the cloud, and parallelized on K different virtual machines. If the cloud is not able to provide
 212 an answer within a reasonable time, the user might decide to leave the request before the response is
 213 delivered. In this scenario, a genetic algorithm is used to solve an optimization problem to propose the
 214 user with a set of recommendations to download from a media service. The service demands of the
 215 algorithms, executed by the genetic program, are highly variable: for this reason, the execution time
 216 might become sometimes very large, forcing the user to leave the application before an actual response
 217 is provided.

218 3.1.2. Fork/Join paradigm

219 The Fork/Join paradigm (see [27] for a formal discussion) is one of the most common extensions
 220 added to conventional Queuing Networks. For instance, one of the most interesting scenarios from
 221 the point of view of performance analysis consists of a load-balancer requiring to split a request into
 222 parallel processing units. In this context, a Fork/Join paradigm stems from the need of efficiently
 223 parallelizing *divide-and-conquer* algorithms (see [28]), which are usually decomposed into a base case
 224 (that is immediately solved) and a recursive case, where the problem is decomposed into smaller
 225 subproblems assumed to be disjunct: in the end, the responses deriving from the solution of each
 226 subproblem are merged. As a result, the main idea is to fork the initial problem into more subproblems,
 227 execute them in parallel and finally synchronize them joining the solutions.

228 Modern extensions of Fork and Join allow for *Quorum based joins*. Let us assume a request is split
 229 into M tasks by a Fork node. A conventional join node would wait for all M tasks to complete the
 230 executions to reconstruct the request and let it continue through the model. Quorum based joins
 231 are characterized by an extra parameter $Q \leq M$. In this case, the Fork/Join part of the request is
 232 considered to be finished when the first Q out of M tasks reach the Join node. The late $M - Q$ tasks
 233 will be simply discarded when they will reach the Join node.

234 3.1.3. Model description

235 The Queuing Network depicted in Fig. 2 models the application of Fig. 1. It is composed of a
 236 Fork/Join that splits a request into two tasks: one models the real request execution (upper branch),
 237 and the other represents the occurrence of timeout due to user abandonment (lower branch). The
 238 execution of the service demand is modeled by K servers queuing node, while the timeout is modeled
 239 by a Delay node. The Join node waits until the first task is executed (in this case, the Quorum is
 240 equal to one). The probability distribution taken into account to describe the service demand is
 241 hyperexponential with an average duration of 1 second and a coefficient variation of 5. The timeout is
 242 considered to be deterministic, and arrivals are assumed to be generated according a Poisson process:
 243 both the duration of the former and the rate of the latter are varied during the study of the model. The
 244 two loggers $L1$ and $L2$ in Fig. 2 track the instant of time the request in execution flows through them

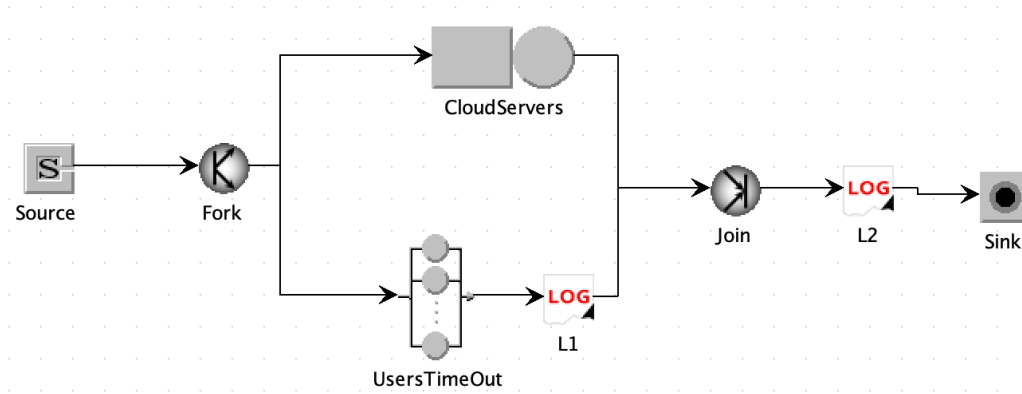


Figure 2. Model of a cloud application with user's timeout

245 and allow to compute the probability $P(Timeout)$ that a request ended due to timeout. In particular,
 246 since the Join has zero duration, a request that ends due to timeout will have the same timestamps
 247 recorded by both the loggers.

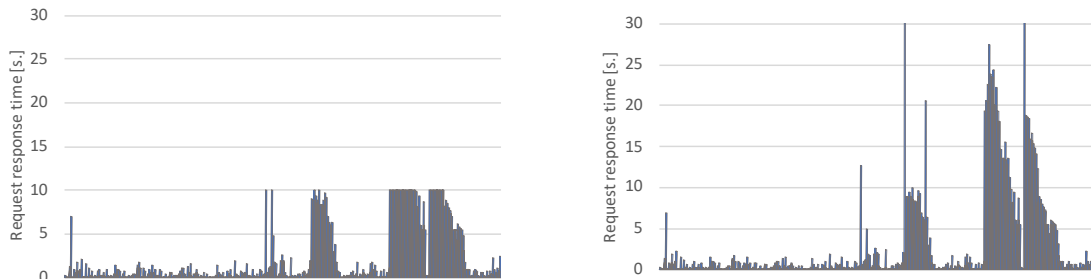


Figure 3. Response times with (left) and without (right) timeout of 10s, for $\lambda = 1.8$ r/s, and $K = 2$ VMs

248 3.1.4. Model results

249 We start considering requests arriving at rate $\lambda = 1.8$ r/s, and a timeout $T = 10$ s. It is interesting
 250 to compare the behavior denoted on the right of Fig. 3, where the spikes generated by the services are
 251 not depending on a timeout variable and the scenario shown on the left, where the timeout is cutting
 252 the service times of the requests exceeding it.

253 This effect is further investigated in Fig. 4, where the response time distribution for different number
 254 K of VMs, ranging from 1 to 3 and an arrival rate of $\lambda = 0.9$ r/s is presented. Although for $K = 1$ the

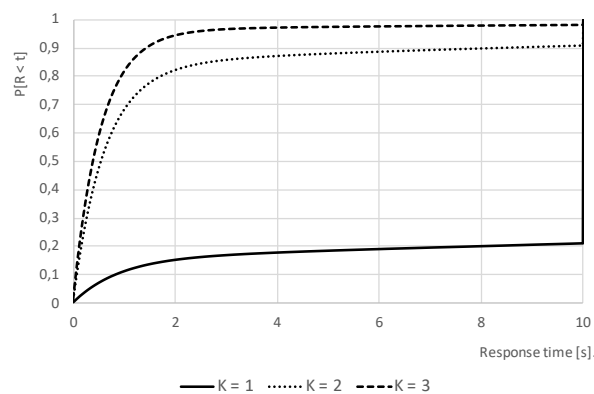


Figure 4. Distributions of Response Time with Timeout=10 sec and $\lambda = 0.9$ r/s for a variable number K of VMs

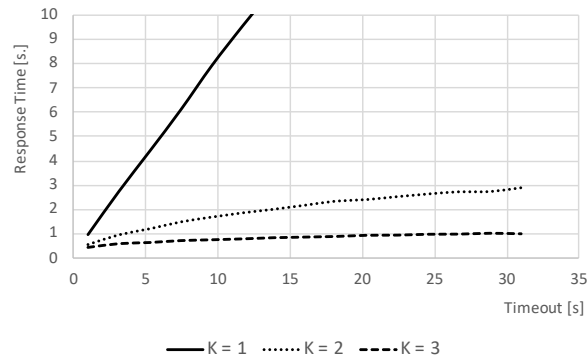


Figure 5. Behavior of Response Time for timeout ranging from 1s to 30s, $\lambda = 0.9$ r/s for a variable number K of VMs

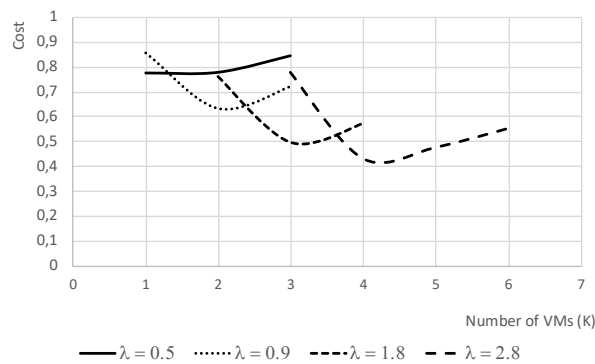


Figure 6. Trade-off between deployment cost (number K of VMs) and user satisfaction for different arrival rates λ .

255 system is stable and characterized by an average utilization $U = 0.9$, most of the requests will undergo
 256 a timeout (the steep rises at time $t = 10$ s): users are then either be served very quickly, or they quit the
 257 system otherwise. With $K = 3$, timeout almost never occurs, while $K = 2$ seems to be a reasonable
 258 tradeoff between deployment cost and user experience, where most of the requests are served on time,
 259 and only less than 10% are subject to timeouts.

260 The effect of the deployment on a different number of VMs is then investigated in Fig. 5, for increasing
 261 timeout durations. With only one VM ($K = 1$), the users are almost always subject to timeouts, and
 262 its value basically defines the time requests will take before being discarded due to abandonment.
 263 Parallel deployments ($K = 2$ and $K = 3$), are instead only marginally influenced by timeouts, since in
 264 any case, most of the users will be served before the event occurs.

265 Finally, Fig. 6, studies the trade-off user satisfaction and deployment cost, by defining a metric:
 266 $cost = 1 - U + P(Timeout)$ where U is the average utilization of the VMs and $P(Timeout)$ is the
 267 probability that a timeout event occurs, as identified by the two loggers. When the system is deployed
 268 on a small number of VMs, the utilization is very high (hence, $1 - U$ is close to zero), but also the
 269 timeout probability is close to $P(Timeout) \approx 1$. Conversely, when the system is deployed on a large
 270 number K of VMs, the time probability is close to $P(Timeout) \approx 0$, but also the utilization is very low,
 271 making $1 - U \approx 1$. The best configuration of the system, is thus the one that minimizes this cost. For
 272 an arrival rate of $\lambda = 0.9$ r/s this occurs at $K = 2$, while for $\lambda = 1.8$ r/s and $\lambda = 2.8$ r/s respectively at
 273 $K = 3$ and $K = 4$. Note that in these cases, respectively $K = 1$ and $K = 2$ could have not been used
 274 since in such circumstances the system would have not been stable. For $\lambda = 0.5$ r/s, the best solution
 275 seems to be $K = 1$, since having a second VM does not sufficiently decreases the timeout probability.

276 3.2. Approximate Computing with Finite Capacity Region

277 The evolution of hardware and software architectures is making approximate computing an
278 emerging technology with great potential. This computing paradigm is based on a very simple
279 concept: an approximate result obtained in a short amount of time is very often sufficient to achieve
280 the objectives of the application with the required accuracy. Essentially, this method trades off accuracy
281 of results with the requested time for their computation.

282 Examples of applications range from parallel search engines to data analytics, video streaming
283 compression, statistical analysis of large data sets, MapReduce transactions, and special hardware
284 implemented for this purpose, e.g., Neural Network Accelerators.

285 The wide range of applications requires the implementation of many approximation techniques.
286 Among them all, the one based on the parallel execution of several tasks of an application and a
287 stopping condition based on the accuracy level required is very popular. Deploying Approximate
288 Computing in a cloud computing scenario is particularly relevant, as this approach allows to contain
289 the costs of autoscaling algorithms by imposing a threshold to the acquired VM numbers (without
290 affecting the performance and maintaining the system responsiveness to the requests' bursts. The VMs
291 parallelism degree is not increased, at the cost of lowering down the results quality).

292 In the following, we show a simple model of Approximate Computing based on parallel
293 computations, implemented with a *Fork* primitive, and a control condition based on the accuracy level
294 required, implemented with a *Join* primitive and advanced synchronization [29].

295 Thus, Approximate Computing has the potential to benefit a wide range of application frameworks
296 e.g. data analytics, scientific computing, multimedia and signal processing, machine learning and
297 MapReduce, and so forth.

298 3.2.1. Description of the problem

299 A car navigation system designed to meet the needs of the users of a smart city must be adaptive
300 to the highly variable conditions set by traffic, municipality constraints and other unexpected events.
301 Typically, a pipeline of various stages is used. The first one is the *Alternative Route Planning* that consists
302 of the identification of k alternative paths from source to destination. This step is referred to as the
303 *K-Shortest Paths with Limited Overlap* [30] problem, since the alternative paths should overlap less than
304 a given threshold. The shortest path in terms of distance or average traveling time might not be the
305 optimal one, as several other parameters may have a significant influence on the computation of the
306 best solution. Examples of conditions that must be considered are i) traffic status, ii) layout of the
307 road map, iii) municipality constraints, iv) transient work in progress, v) car accidents, vi) severe
308 climate conditions, and other unexpected events. Consequently, the computation times required by the
309 algorithms are highly variable non only because of the different numerical algorithms implemented
310 but also for those other variables considered.

311 To be effective, the *mean response time* R (i.e., the time spent by a request in execution in the *Fork/Join*
312 *region* plus the time spent at the *Fork* waiting to enter the region) of this step for the computation of the
313 best route must be ≤ 3 sec. To achieve this performance goal, sub-optimal solutions are computed
314 considering only the results of the first k algorithms of N that completed their executions. The model is
315 used to investigate the impact on the best route computation time of the variability of the execution
316 times of the N algorithms and of the subset size k .

317 3.2.2. Finite capacity regions

318 In this case, Queueing Network models contain nodes that include, in turn, regions applying
319 specific policies (*finite capacity constraints or FCC*). A policy denotes upper bounds on the number
320 of requests that can reside in its service nodes at the same time. Ready customers (contraposed to
321 customers that are waiting at terminals) wait to enter into a FCC and, once they are inside, compete to
322 acquire the available resources. An interesting case occurs when the population belongs to multiple

323 classes. In this case, the constraints undertaken are the following: i) one bounding the number of
 324 customers in a region for a specific class and ii) a shared one limiting the number of customers without
 325 class distinctions. FCC allows to capture the occurrence of performance saturation effects determined,
 326 for example, by memory constraints.

327 3.2.3. Model description

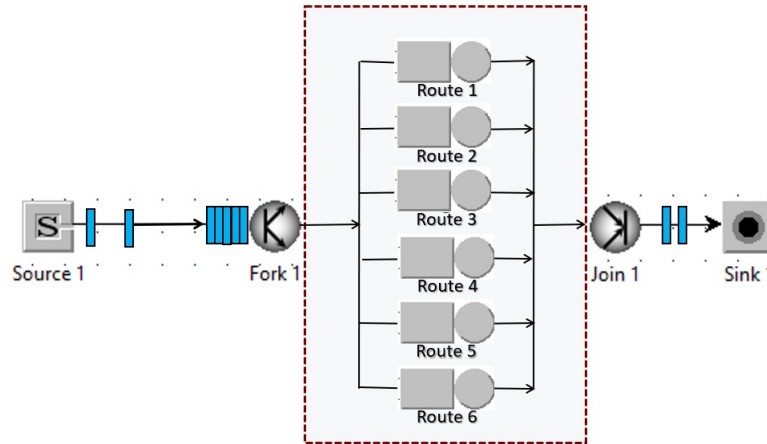


Figure 7. Layout of the model for the computation of the best route

328 The layout of the model is shown in Fig. 7. It consists of a *Fork/Join* region with a capacity limited
 329 of one request. Each arriving request is split at the *Fork* into *six tasks* that are executed in parallel. Each
 330 task represent a different algorithm for the route computation and is executed by a dedicated processor
 331 represented by a *Delay* node.

332 In this specific case, it is assumed that the storage does not origin a bottleneck, since all the nodes
 333 exploit a local content and their execution follow a fairness policy.

334 The parameters of the service demands of the six algorithms are reported in Table 2. Their mean
 values are highly variable (from 1 to 5 sec.) and their coefficients of variation vary from 0.5 to 5.

Table 2. Service demands [sec], coefficients of variation, and distributions of the computation times of the six algorithms

Component	Parameters		
	mean	Coeff. of Var.	distribution
Algorithm 1	1	cv=5	hyperexp
Algorithm 2	3	cv=3	hyperexp
Algorithm 3	1	cv=1	exp
Algorithm 4	2	cv=1	exp
Algorithm 5	4	cv=0.7	Erlang
Algorithm 6	5	cv=0.5	Erlang

335 The data of Table 2 have been assumed to consider a good variety of computation times. For
 336 the sake of simplicity, and without affecting the accuracy of the model, we limit the capacity of the
 337 *Fork/Join region* to one request, i.e., only the tasks generated by one request can be executed in the
 338 region. Requests that arrive when another one is in execution wait at *Fork* node.

339 A typical characteristic of car navigation systems is the presence of fluctuations in the arrival flow of
 340 requests. The arrival rate considered in the study is $\lambda = 2 \text{ req/sec}$ with the exponential distribution of
 341 the interarrival times.

342 The request exits the *Fork/Join region* when at least k algorithms have completed their executions. The
 343 *Join* implements a synchronization rule referred to as *Quorum k*.

3.2.4. Model results

Fig. 8 shows the behavior of the mean response times and of the 95th percentiles for the computation of the best routes with respect to the number of algorithms considered from $k = 1 \div 6$. The mean values are obtained with 99% confidence intervals. As can be seen from the figure, it is

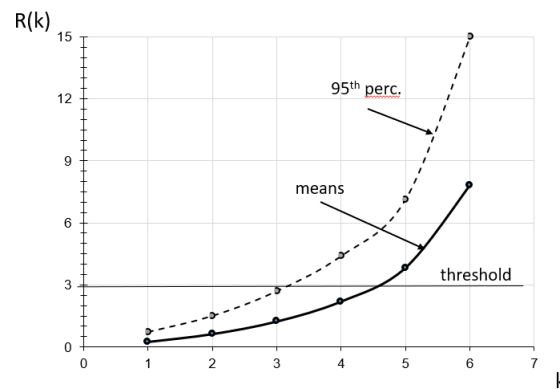


Figure 8. Mean response times and 95th percentiles for the computation of the best route waiting at the *Join* the first k algorithms that completed their executions.

sufficient to decrease the number k of algorithms waiting for synchronization at the *Join* from six to five to obtain a 51% reduction in the mean response time (from 7.81 to 3.82 sec) and of 52% of the 95th percentile (from 15 to 7.15 sec). With $k = 4$ the mean response time is 2.19 sec. and the 95th percentile is 4.41 sec. This value meets the 3 sec threshold. The accuracy of the results provided was positively assessed by the data collected in various periods of the day with representative traffic requests.

3.3. MapReduce with Class Switch

MapReduce is one of the first techniques used to support Big Data applications. From the original proposal by Apache foundation, which was supported by the Hadoop project, several different extensions have been proposed. All the techniques however are based on very similar principles. In short, (see fig. 9), all the data that need to be processed (which might consist of Exabytes of data), are split into chunks that are distributed over a large set of participating nodes. Each node, beside storing part of the data, can also perform operations on it. Using specific software patterns, complex operations can be performed to obtain insight information from the considered Big Data collection. MapReduce is the simplest of these patterns: first, the Map operation applies a function to each entry of the database, generally performing searching and sorting tasks. The Reduce phase, collects the intermediate results to produce the final answer. The number of chunks does not necessarily needs to correspond to the number of nodes: in many best practices, the number of chunks is much higher than the number of nodes to increase the parallelism, so that faster nodes can start working on new chunks while the slower ones are still finishing their job. MapReduce is closely related to Cloud computing, since VMs represents the most natural way of acquiring nodes to support the parallel execution process.

3.3.1. Description of the problem

Sizing a MapReduce application is a non-trivial task. Indeed, increasing the number of nodes decreases the running time of the application. However, due to the hyperbolic decrease and the increasing synchronization complexity, after a given point (depending on the application and on its demands), performance improvements becomes negligible at first, up to another point, where increasing the nodes makes only the system's behavior worse. Performance model are thus of paramount importance to correctly decide the optimal number of resources, to obtain the best results

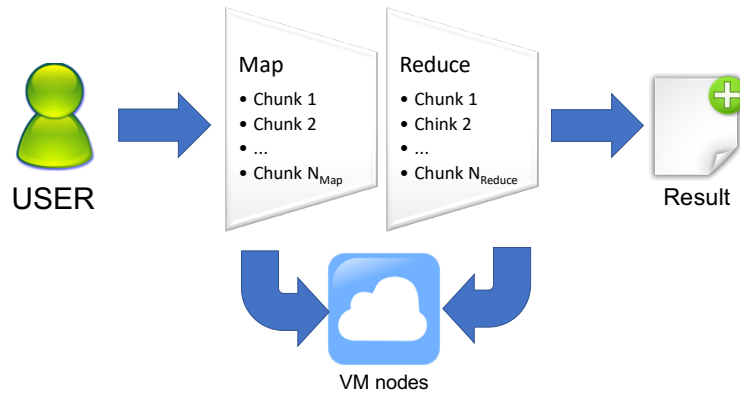


Figure 9. A MapReduce application deployed on Cloud

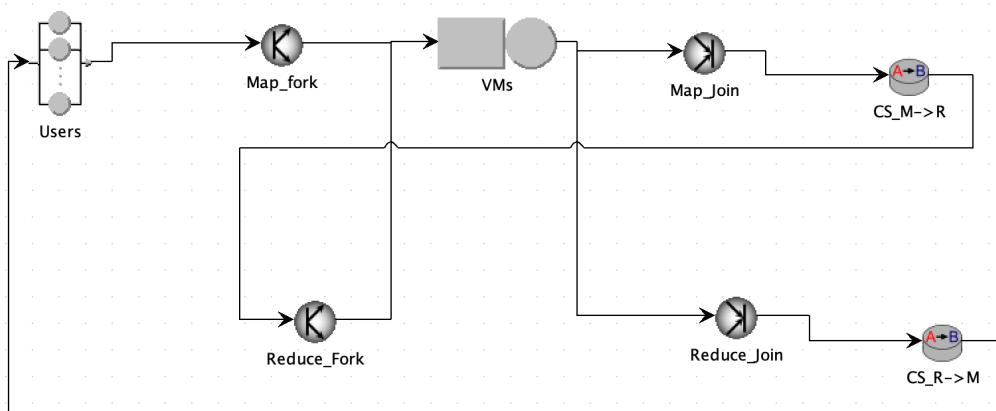


Figure 10. Model of a MapReduce application

377 with the least expense.

378 While fork/join paradigms are the basis for modeling MapReduce such as Big Data processing
 379 environments, they cannot capture the fact that the same processing nodes are used in different ways
 380 for different stages of the algorithm (i.e. Map and Reduce). The *Class Switch* feature, introduced in the
 381 next section, becomes the key tool to model such behavior.

382 3.3.2. Class Switch

383 One of the main assumptions adopted when working with simple QNs, concerns the type of
 384 customers, which are assumed identical from a statistical point of view. However, when regarding a
 385 real system, this assumption doesn't necessarily hold as more parameters, such as the service time and
 386 the routing probabilities, must be taken in account. Therefore, it is necessary to postulate the existence
 387 of different types of customers by introducing the terms of *chain* and *class*. The former identifies a
 388 situation where a customer belongs to the same type during the entire execution. The latter denotes
 389 instead a temporary classification: from this perspective, a customer is able to change from a class to a
 390 different one while executing within the system according to a probability. The type of class plays an
 391 important role in characterizing the customer service time (in each node) and the routing probability.
 392 Classes can be partitioned into chains, which prevents the case where a job switches from classes
 393 that are part of different chains. A *class switch* operation, allows a customer in chain to change its
 394 appearance to the server, by presenting itself with another class.

3.3.3. Model description

A simple MapReduce application is presented in fig. 10. In particular, we consider N users, that every $Z = 100$ s (modelled by delay node *Users*) submit a map-reduce job to the system. There are two classes in the system, representing respectively the Map and Reduce stages of the application. All requests starts in the Map phase, and are characterized by the corresponding class. The *Map_fork* node splits the job into M_{Map} tasks, which are then rejoined in node *Map_join*. The class switch node *CS_M->R* changes the class of the job to the Reduce stage, which is immediately split into M_{Red} reduce tasks in node *Reduce_Fork*. The job is finished when all the Reduce tasks terminate: this is modelled by the join operation performed by node *Reduce_join*, immediately followed by class switch node *CS_R->M* that restores requests to the starting Map class. All tasks are served by the K server queueing node VMs modeling the cloud environment running the application. The service time of the all the VMs is assumed to be exponential, with a different duration D_{Map} and D_{Red} respectively for the Map and Reduce stages. In the next experiments, we assume $M_{Map} = 64$, $M_{Red} = 32$, $D_{Map} = 1$ s and $D_{Red} = 2$ s.

3.3.4. Model results

The system is studied for a different number of users N , starting with $N = 10$ and up to $N = 50$, for different parallelization levels K (corresponding to the number of acquired VMs). Fig. 11 shows

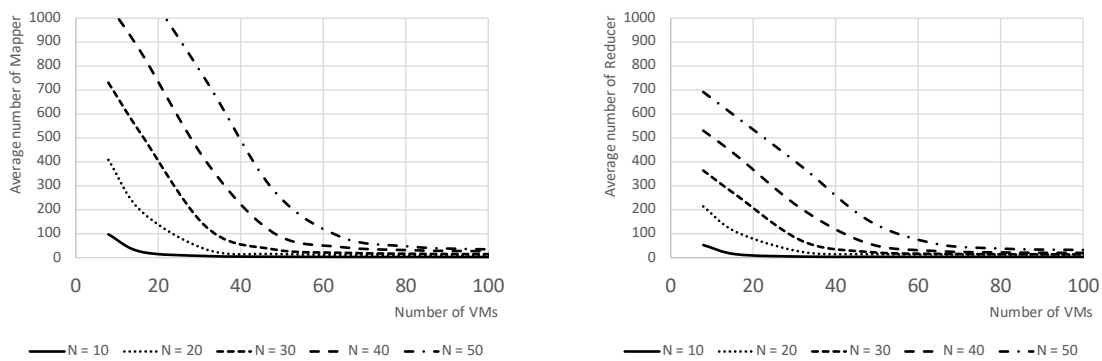


Figure 11. Average number of mappers (left) and reducer (right), for different number of users N , and K VMs.

the average number of tasks in execution in the VMs for the Map and the Reduce stage. As expected, it grows with the number of users N , and decrease with the number of K VMs. It is interesting however to note how the decay is not linear, but not even clearly hyperbolic: this is due to the different configurations of the Map and Reduce stages.

Fig. 12, shows an exploitation of the model to determine the best number of VMs to support a given number of users N . In particular, it explains how the model can be used to size the system to achieve a target average response time $\tau = 200$ s. When the number of users is very low ($N = 10$), even a small amount of VMs ($K = 8$) is able to provide average response time much lower than the threshold. However, with a population of $N = 50$, at least $K = 24$ VMs are necessary to achieve acceptable performances. It is also interesting to see that for $K \geq 64$, all the considered configurations have basically the same lowest performances: this is due to the fact that each job is split into $M_{Map} = 64$ tasks. If no other job is currently running, with $K \geq 64$ all tasks can be executed in parallel in one shot.

4. Multiformalism QN/PN techniques

According to multiformalism, different part of a system can be modeled by using different formalisms flavors (the choice depends on the modeler's familiarity with one or more languages). In this way, it is possible to lower the learning curve and match the user abstraction as a result. The model so derived can be solved by defining the proper combination of formalisms by mapping the model

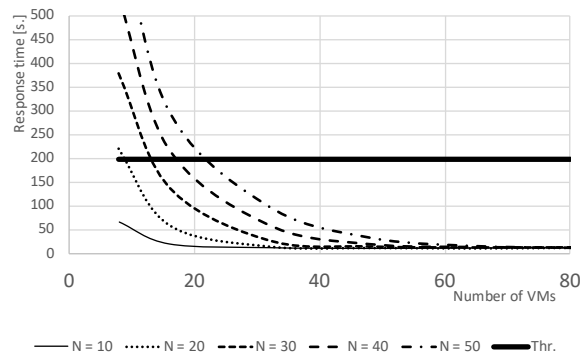


Figure 12. Determining the optimal number of VMs to obtain an average response time R below a given threshold ($\tau = 200$ s. in this example).

4.29 concepts into solvers primitives. Multiformalism has proven to be successful in different areas such as
 4.30 biology, fault-tolerant computing and disaster recovery. As a result, this interdisciplinary aspect has
 4.31 created interesting links between different communities of modelers. A variety of different software
 4.32 tools have been implemented to date. With regard to JMT, the considered multiformalism approach
 4.33 allows the integration of Queueing Networks (QN) and Generalized Stochastic Petri Nets (GSPN, see
 4.34 [31] for the theoretical background of this formalism).

4.35 4.1. Hybrid Cloud

4.36 A hybrid cloud consists of an on-premises infrastructure, based on a private cloud, and resources
 4.37 acquired as-a-service from a public cloud provider. Among the various factors motivating a boost
 4.38 in the diffusion of these architectures are high security, controlled performance, large scalability, fast
 4.39 adoption of new technologies, and cost savings.

4.40 The problem of scalability in hybrid clouds is typically addressed through the *dynamic*
 4.41 *provisioning* of resources from the public cloud. The model presented in the following case study
 4.42 addresses this problem by implementing an algorithm that dynamically routes the requests to the
 4.43 public cloud when the load on the private component exceeds a threshold value.

4.44 4.1.1. Description of the problem

4.45 This case study concerns a model of an IT infrastructure based on a hybrid cloud. More precisely,
 4.46 it focuses on modeling the process for dynamic resource provisioning. This method is able to acquire
 4.47 virtual machines (VMs) on-demand from the public cloud when requests exceed the capacity set as a
 4.48 threshold on the private cloud. The multiformalism model implemented consists of elements of Petri
 4.49 Nets and Queueing Networks. The hybrid cloud scenario considered is represented in Fig. 13 .

4.50 The incoming requests are processed by the user interface of the application and, after some formal
 4.51 and security check, are sent to the *Load Controller* module. To satisfy the performance requirements,
 4.52 the software architecture of the app has been designed assuming that a dedicated VM of the local cloud
 4.53 is assigned to each request in execution. Since several highly fluctuating workloads share the resources
 4.54 of the private cloud, a limit is set on the maximum number of VMs dedicated to this application. When
 4.55 this threshold is reached, the new VMs are acquired on-demand from a public cloud.
 4.56 The impact of this threshold on global response time and throughputs of both the two clouds needs
 4.57 to be investigated. The objective of the study is the identification of the computational capacity, in
 4.58 terms of number of cores and power, of the servers of the private cloud that are required to satisfy the
 4.59 performance target with cost savings. In fact, VMs with the same computational power as a private
 4.60 provided by the public cloud are much more expensive. Therefore, to save costs, the VMs provided by
 4.61 the public cloud are much less powerful than the private ones. Thus, there is a tradeoff between the
 4.62 VMs provided by the private and public clouds, performance and the infrastructure costs.

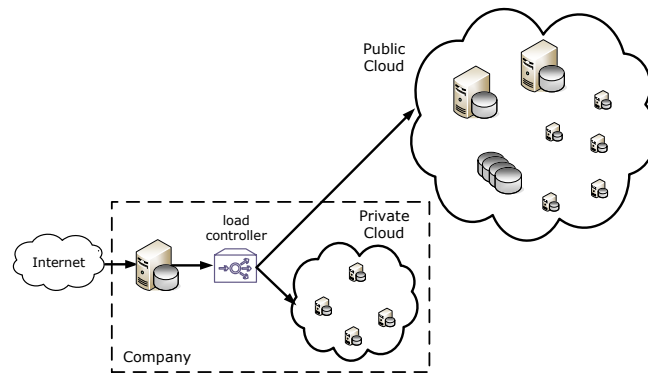


Figure 13. The hybrid cloud scenario considered.

4.1.2. Model description

The layout of the model is shown in Fig. 14. The workload consists of two classes of customers: the incoming requests, representing the user demands of computation time, and the VMs (the tokens), representing the number of VMs available in the private cloud.

The *JoinPrivate transition* is enabled when a request arrives in *Arriving place* and there is at

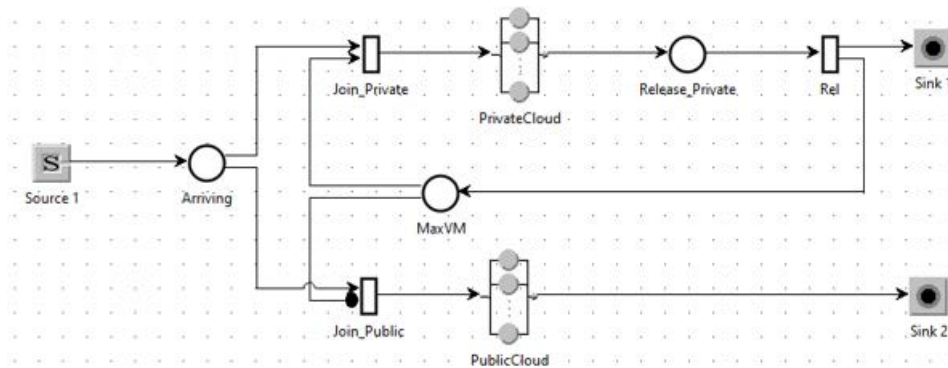


Figure 14. Layout of the model for the dynamic provisioning of VM in a Hybrid cloud.

least one token available in *MaxVM place*. Each time the *transition* is activated, a VM of the private cloud is assigned to the request and the value of *MaxVM* is decreased by one. When this value is zero, the *Inhibitor arc* from *MaxVM place* and *JoinPublic transition* activate the latter and the request is addressed to the public cloud. When a request has been completely executed in the private cloud, the *Rel transition* routes it to *Sink1* and a token is sent to the *MaxVM place* incrementing the number of VMs available.

The two clouds are represented by two delay stations since there is no competition for an available VM in both the clouds. The arrival rate of the requests to be considered in the study is $\lambda = 50 \text{ req/sec}$. This value has been assigned by the application designers since it is representative of a medium/high load that according to the business plan should be achieved in a year. The fluctuations of arrivals has been modeled with a $cv = 4$ of the hyper-exponential distribution of interarrival times. The mean service demands of the private cloud is 2.5 sec while those of the public cloud is 7.5 sec. The high variability of service times was considered assuming their distributions as hyper-exponential with $cv = 6$. The time required by the processing of a request in the other infrastructure components, such as the user interface and the load controller, are negligible compared to the service demands of the VMs, therefore they have been considered as small increases in their values.

4.1.3. Model results

The behavior of the algorithm for dynamic provisioning of virtual machines is highlighted in Fig. 15 that shows the trend of the number of VMs in execution in the two clouds private (a) and public (b) in the interval 0÷50 sec. The arrival rate is $\lambda = 50$ req/sec and the threshold of the VMs in the private

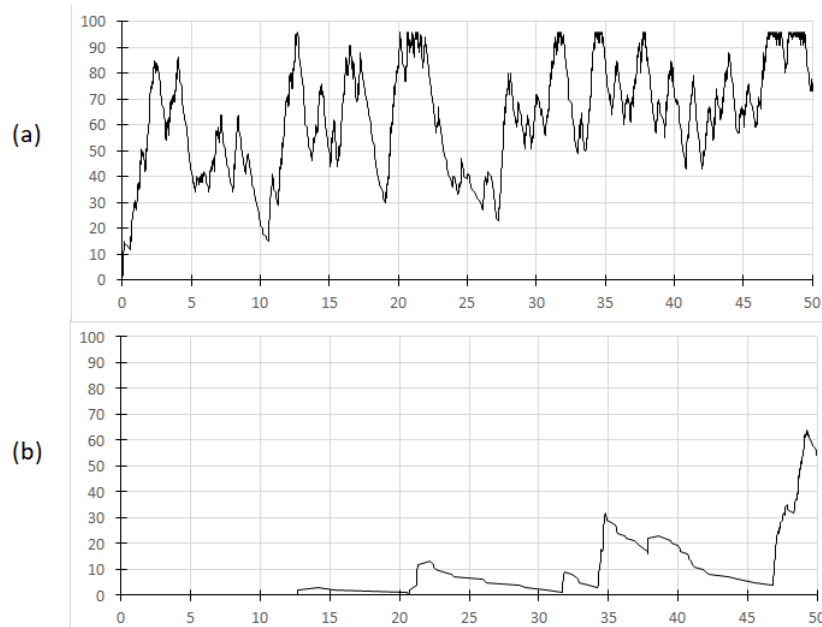


Figure 15. Behavior of the number of VMs in execution in the private (a) and public (b) clouds in the interval 0÷50 sec. The threshold for the VMs in the private cloud is 96 and the arrival rate of requests is $\lambda = 50$ req/sec.

cloud is 96. As can be seen in Fig. 15a, when the number of requests in execution in the private cloud is greater than 96, the provisioning of the new VMs is dynamically routed to the public cloud (see e.g., the interval 45-50 sec in Fig. 15b).

The number of requests in execution in the two clouds as a function of the maximum number of VMs

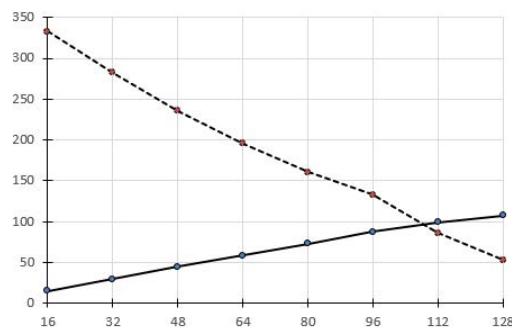


Figure 16. Number of requests in execution in the private (solid line) and public (dashed line) clouds vs max number of VMs in the private cloud.

MaxVM that can be provisioned in the private cloud is shown in Fig. 16. The range of MaxVM evaluated is 16÷128. The mean response time R of the system as a function of the MaxVM is depicted in Fig. 17.

With 96 VMs the average response time is close to 4 sec, a value that is considered acceptable as a performance target with the associated costs. As can be seen in Fig. 16, 87 VMs are in execution in the private cloud and 133 in the public cloud with MaxVM = 96.

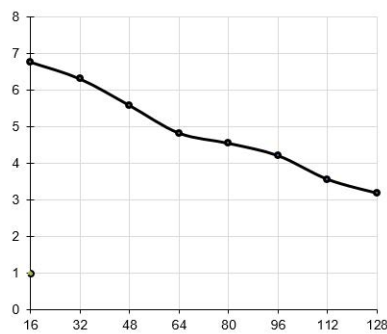


Figure 17. Global Response time vs maximum number of VMs in the private cloud.

497 The costs of the infrastructure can be evaluated considering the throughput of the two clouds as a
 498 function of MaxVM (see Fig. 18).

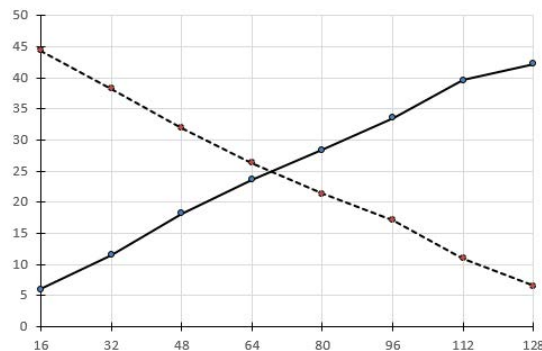


Figure 18. Throughput of private (solid line) and public (dashed line) cloud vs maximum number of VMs in the private cloud.

499 4.2. *Batching in IoT-based healthcare*

500 The proliferation of IoT in the health care scenario has introduced new problems, which have been
 501 faced for effective use of their potential. Important benefits can be obtained in all areas of e-Health, in
 502 particular in those that use IoT integrated in information infrastructures enabling the use of ubiquitous
 503 computing technologies. Patients can be monitored anytime anywhere, either in special hospital wards
 504 or remotely, through the use of wearable sensors and smart medical devices.

505 Sensors may detect a variety of patient physiological signals, such as, temperature, pulse, oxygen
 506 saturation, blood pressure and glucose, ECG, EEG, as well as other body motion related variables that
 507 can help accurately monitoring patient movements. Among the potential benefits that can be achieved
 508 by body sensors, and more generally by IoT smart devices, in e-Health monitoring are the high rate of
 509 data transmission and the minimization of end-to-end data delivery time. The interconnections among
 510 the various components of the networks, e.g., IoT devices, intelligent medical devices, edge and fog
 511 systems, hospital and cloud servers, patients and medical staff, are implemented through cabled or
 512 wireless networks with low-power communication protocols.

513 The following case study focuses on body sensor networks, and more specifically on the study of
 514 the trade-off that exists between performance of the network (data delivery time) and the energy
 515 consumed by the data exchange (the cost of transmission).

516 The implemented model is derived from a more complex version of [32] that considers a
 517 completely different scenario: the smart monitoring of fog computing infrastructures. The key feature
 518 of these models is the dynamic management of the buffer of requests based on the intensity of arrivals

519 and the expiration of a periodic trigger. With the multiformalism models it is possible to implement
 520 algorithms with dynamic behavior as a function of the workload characteristics.

521 4.2.1. Description of the problem

522 Fig. 19 shows the target e-Health scenario considered. The data collected from body sensors are
 523 transmitted through wireless or wired connections to the Hedge nodes located as close as possible,
 524 where they are pre-processed and then sent to the fog nodes (if any) or to the hospital servers for their
 525 complete processing.

The data arriving at the hospital servers are subject to fluctuations generated by the different

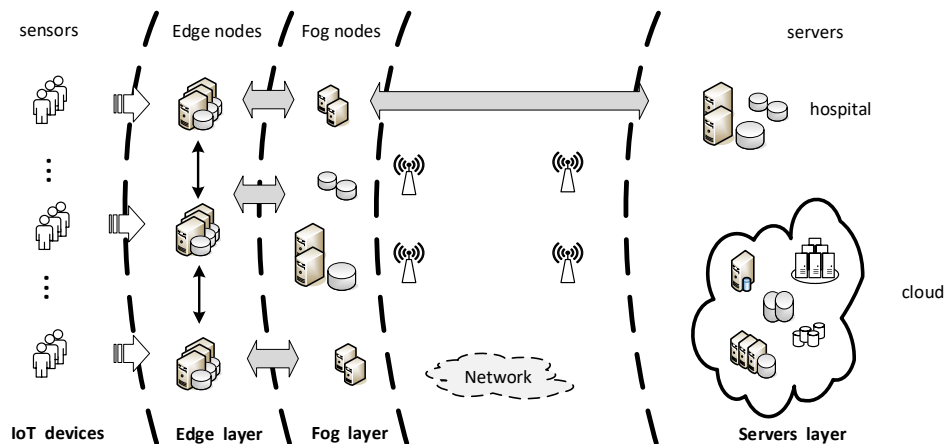


Figure 19. The considered e-Health scenario.

526 type of physiological signals detected and the health conditions of the patients. Indeed, different
 527 type of measured variables require a different frequency at which a detected signal is available for
 528 transmission. For example, in body temperature the sampling rate can be one per minute, in a pulse
 529 oxygen monitoring the rate can be one per second, while in other variables, like ECG or EEG, it can be
 530 of the order of several hundreds per second. Also, when a patient's health condition is assessed as
 531 critical, new sensors are activated and the detection rate of other monitored variables can be increased
 532 under the control of edge or fog nodes. Among the many problems that need to be addressed, this
 533 case study concerns the following:

- 535 • identification of the amount of data that must be considered in each transmission to hospital
 536 servers in order to satisfy the performance requirement in term of end-to-end data delivery time
 537 and minimize the energy consumption of the operations;
- 538 • identification of potential critical health conditions of patients that need urgent investigation, i.e.,
 539 fast response time.
 540

541 The former problem requires studying the trade-off between the time required to deliver the detected
 542 signal to the servers in the upper-layer of the medical infrastructure, and the cost associated with the
 543 transmission operation. The immediate transmission of a detected signal minimizes its end-to-end
 544 response time from either the hospital servers or the cloud. However, the set up costs of the connection
 545 cannot be shared with other signals. The technique of *batching* the data of several signals to be

546 transmitted in a single operation is used to approach this problem. The impact on end-to-end delivery
 547 time of different batch sizes, and thus on the number of operations required to transmit the signals
 548 detected by a set of sensors, must be studied. Knowing the number of sensors connected to an edge
 549 system and the type of signals detected, it is possible to derive the arrival rate of the requests to the
 550 hospital servers. Then, once pre-processing is complete, the data are stored in a buffer where are ready
 551 to be transmitted. The management of this buffer is crucial to achieve the two objectives described
 552 above.

553 The implemented algorithm considers the number and types of signals detected by the sensors
 554 connected to the edge nodes, the fluctuations of arriving traffic considered *regular* and the arrival
 555 patterns that must be transmitted with *priority* as they can be associated with a patient in critical
 556 conditions. The most important elements of the implemented model simulating this algorithm are
 557 described in the next section.

558 4.2.2. The Model

559 Multiformalism models allow the exploitation of Queueing Networks and Petri Net primitives to
 560 represent each concept with the most appropriate technique. In order to describe the dynamic behavior
 561 of the batching algorithm we used the PN primitives while the QN primitives were used to represent
 562 the other components of the e-Health infrastructure.

The layout of the model is shown in Fig. 20. The workload consists of two classes: *signals* detected by

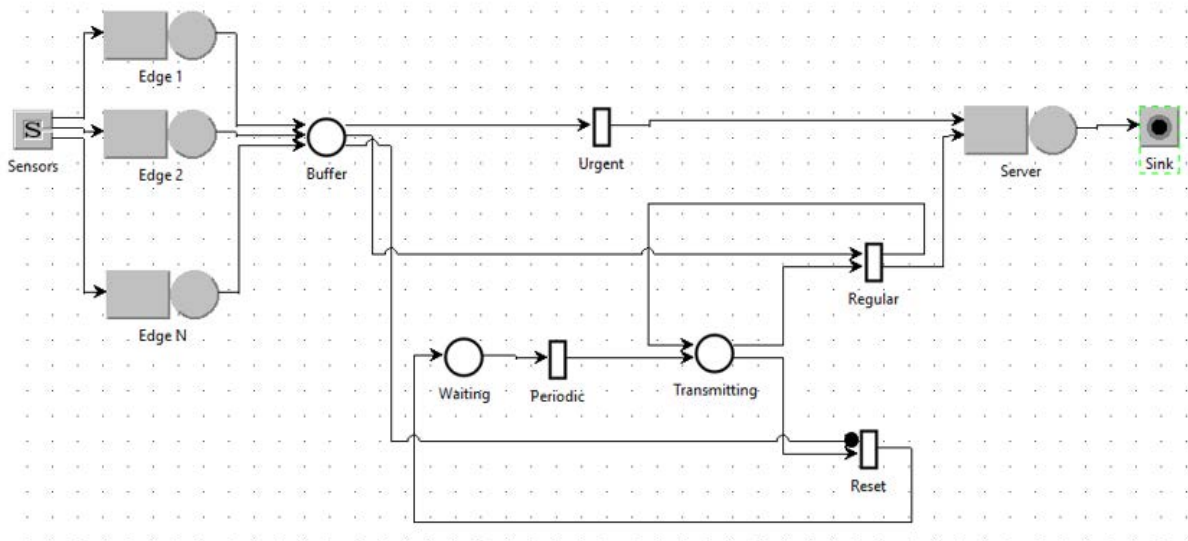


Figure 20. Layout of the model for the smart batching of signals in the e-Health scenario considered.

563 sensors (referred to as requests) and a *token*, needed to model the periodic/triggered management of
 564 the requests.
 565

566 The key feature of the model is the algorithm that manages the transmission requests batch. To ensure
 567 that the model and the presentation of the algorithm remain as simple as possible, we adopted several
 568 assumptions that have a minimal impact on the performance but that greatly simplify the description
 569 of other parts of the model.

570 The set of physiological signals detected by each patient is the same, and each edge system monitors
 571 several patients. The computational power of the hospital servers has been oversized compared to
 572 the processing time required by signals. In the model we do not explicitly represent the fog systems
 573 since their processing time per request is negligible compared to the service time required by the Edge
 574 nodes. Moreover, they were considered as small increases in the service times of the Edge nodes.

575 The global arrival rate λ of data generated by sensors is modeled by all the source Sensors as a single

576 aggregated Poisson process of rate λ . This flow is evenly distributed among the Edge systems modeled
 577 with Edge queuing nodes. The times required to process the data of a signal, i.e., the service time
 578 of a visit by an Edge node, are exponentially distributed with mean $D_{Edge} = 200\text{ ms}$. At the end
 579 of this processing phase, the requests are buffered, i.e., routed to the *place Buffer*, and ready to be
 580 transmitted. They are transmitted to the hospital servers (or the cloud servers) that must perform
 581 their complete analysis requiring a service time exponentially distributed with mean $D_{server} = 500\text{ ms}$.
 582 Requests follow two paths: one for *Regular* requests, and one for *Urgent* requests. The requests in the
 583 buffer are managed according to two different policies:

- 584 • The buffer is emptied (i.e., the requests that are in the buffer are transmitted) *periodically* with a
 585 period defined according to the number and type of signals detected by all sensors. Requests are
 586 assumed to belong to patients under *Regular* conditions and are sent at the end of the period;
- 587 • The buffer is emptied when the number of requests in the buffer reaches a *threshold* value β ,
 588 i.e. the maximum batch size. In this case, requests with such a high arrival rate are assumed to
 589 indicate the presence of a *critical* condition for one or more patients. Therefore, requests in the
 590 buffer are considered *Urgent* and must be sent immediately without waiting for the end of the
 591 emptying period.

592 The *periodic* transmission of *Regular* requests is modeled by the loop between *places* and *transitions*
 593 *Waiting*, *Periodic*, *Transmitting* and *Reset*. The deterministic firing time of *transition Periodic*
 594 represents the duration of the clock for the transmissions of the requests arrived in *Buffer*. This value
 595 is computed by analyzing the detection rate of the sensors in normal operating conditions. According
 596 to the configuration analyzed, we considered *15 sec* as *constant firing time* of the *transition Periodic*
 597 (i.e., as *periodic empty cycle time*). As soon as the empty cycle expires, a token is transferred to *place*
 598 *Transmitting* where two alternatives are possible. If there are requests in the buffer, the immediate
 599 *transition Regular* will be enabled and will transfer them to the transmission channel. When the
 600 buffer is empty, either because all requests have been transferred or because no requests arrived in the
 601 periodic time frame, the immediate *transition Reset* fires, due to an *Inhibitor arc* that connects it to the
 602 *Buffer*, and restarts the timer.

603 The *Urgent* requests are managed by the immediate *transition Urgent* that is connected to *place*
 604 *Buffer* with an input *arc* of weight β . When the threshold value β is reached, the batch of requests in
 605 the *Buffer* are immediately transmitted to the hospital server. Note that also the *arc* that exits the
 606 *transition Urgent* has weight β , since the entire batch of requests is sent to the server.

607 4.2.3. Model Results

608 We considered several configurations of the system by modifying the global arrival rate λ , the
 609 threshold β of *Urgent* requests, and the empty buffer cycle time. In this section, we limit the description
 610 of the results to those that emphasize the impact of the dynamic management of the buffer of requests
 611 on the performance of the system and related costs.

612 The arrival rates of the requests considered in the study are $\lambda = 0.1 \div 1.9\text{ req/sec}$. The interarrival
 613 times are exponentially distributed. These values were assumed to be representative of the potential
 614 number of patients in an emergency ward of small or medium sized hospitals.

615 Fig. 21 shows the behavior of the *System Response Time R*, i.e., the end-to-end time required by a
 616 signal from its detection to the completion of its processing by the hospital server, for the complete
 617 range of arrival rates. The family of curves refers to different values of β , from 2 to 20, i.e., the threshold
 618 for the identification of *Urgent* batches of requests will be transmitted immediately.

619 With a λ increase we model an increase in the number of sensors or in the detection rate and the
 620 workload managed by Edge nodes is also greater. For small values of β (2 and 5) there is an initial
 621 decrease of *R*. This is because with such small sizes nearly all batches are considered *Urgent* and
 622 therefore most requests are transmitted almost immediately when they join the *Buffer*. With larger
 623 values of β (10, 15 and 20) this initial decrease of *R* as λ increases is not present. When λ becomes

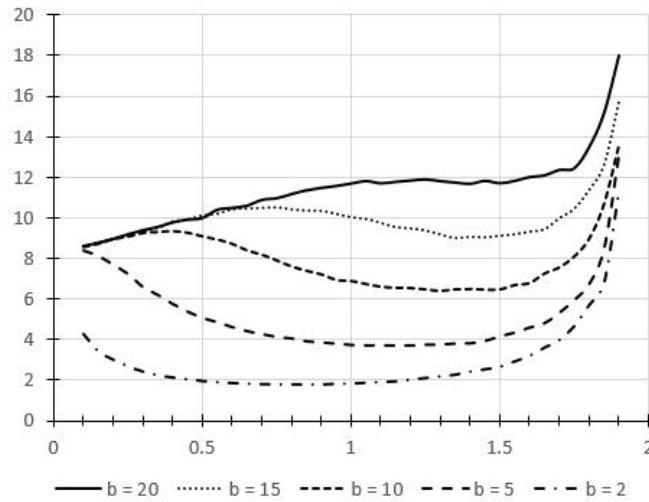


Figure 21. System Response Time vs arrival rate of requests for various values β of the batch sizes for *Urgent* requests.

624 greater than 0.1 the value of R starts to increase from the beginning. With these batches sizes, R begins
 625 to decrease for larger λ values than those obtained with smaller β . The motivation for this behavior
 626 is that as λ increases, the number of requests in the buffer increases as well, and if β has higher
 627 values the threshold is reached with less probability and the number of batches transmitted when the
 628 period expire is greater. However, with further increases of λ the threshold is reached more easily and
 629 therefore there is a shorter waiting time for the requests in the buffer, i.e., the *Urgent* requests increase.
 630 When λ is greater than 1.5 req/sec, R increases for all β since the response time of the highly utilized
 631 hospital server becomes the dominant part of its value.

632 According to the objectives of the study, the R values should be analyzed together with the cost (energy
 633 consumption) for transmissions. It is assumed that the cost is directly proportional to the number of
 634 times a batch is transmitted, i.e., a buffer is emptied. Indeed, the greater the number of times and the
 635 sizes of the transmitted requests at the same time, the better the energy efficiency of the system. In
 Fig. 22 the transmissions per second are shown for all the values of the considered λ . As expected,

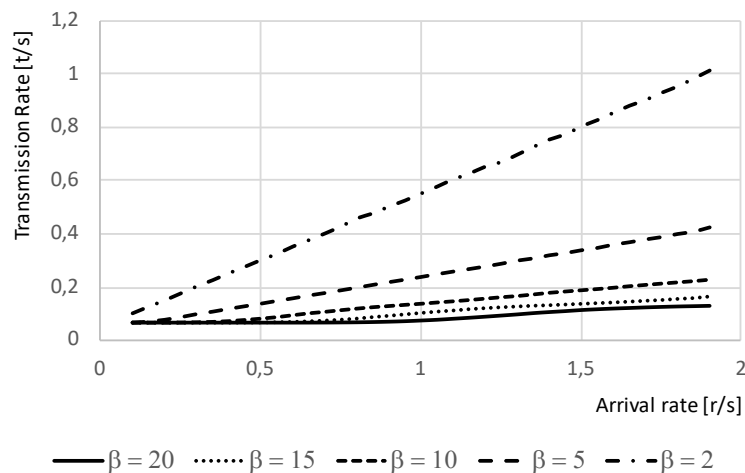


Figure 22. Batch transmissions per second vs arrival rate of requests for various values β of the batch sizes of *Urgent* requests.

636 with batch size $\beta = 2$ the number of transmissions per second are the highest, while the minimum for
 637

638 $\beta = 20$. However, to have a meaningful result, these values should be taken in account together with
 639 the System Response Times. To this end, the metric System Power, introduced in [33] and combining
 640 the throughput X of a system with its response time R , is considered. This metric is the ratio X/R
 641 of throughput and response time, and captures the level of efficiency in executing a workload. The
 642 maximum Power corresponds to the optimal operating point for the system, i.e., the point in which
 643 the throughput is maximized with the minimum response time. In our system, we have considered
 644 the ratio of transmissions/sec and System Response Time. Fig. 23 shows the ratio of the two metrics of
 Fig.s 22 and 21.

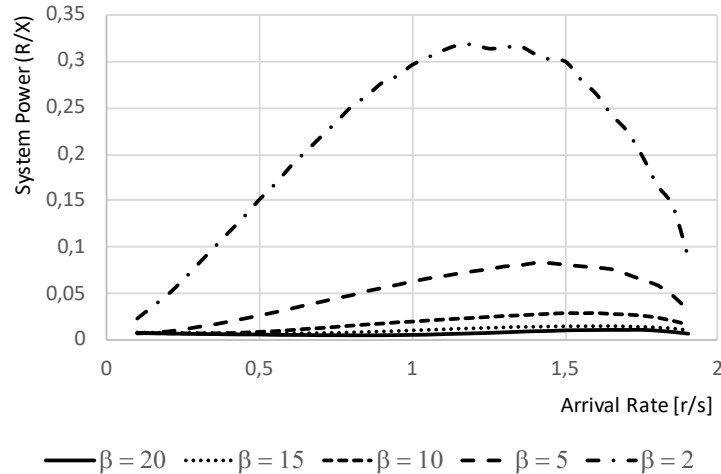


Figure 23. System Power vs arrival rate of requests for various values β of the batch sizes of Urgent requests.

645
 646 The optimal operating points of the system are clearly identified as a function of the batch size β
 647 for the *Urgent* requests and the global arrival rate of signals.

648 5. Conclusions and future work

649 In recent decades, Queuing Networks and Petri Nets techniques have rapidly developed in order
 650 to respond to the needs of more and more complex applications. Furthermore, the growing availability
 651 of computational resources and large capacity memories have greatly enhanced the power of these
 652 formalisms within the boundaries of their semantics. However, taking into account these factors, it is
 653 clear that the implicit limitations of QNs and PNs make the analysis of many of the current problems
 654 unsuitable.

655 Indeed, in such cases, the idea that as a function of the characteristic of the problem, there exists a
 656 technique that is more suitable than others to solve the problem considered is no longer valid. New
 657 methods and tools must be considered in order to broaden the spectrum of scenarios and issues that
 658 can be studied and resolved.

659 To this end, the contribution proposed by this work consists of the review of specific paradigms
 660 applied to complex scenarios, providing more insight into the interpretations of the experimental
 661 results, realized with the support of JMT, a dedicated suite of tools for the performance evaluation and
 662 modeling of systems.

663 From a critical perspective, it is interesting to open a debate comparing the results achieved in the
 664 previously discussed case studies and work from existing literature. With respect to timeout modeling,
 665 Holvoet et al.[34] discuss the requirements that a software modeling approach should meet. In this
 666 work, the authors introduce a new formalism called Object Petri Nets (OPN) based on Colored Petri
 667 Nets (CPNs, see [35]), proposing a new kind of transition called *timeout transition*, which is in turn
 668 based on *non-deterministic input arcs*, and *timeout output arcs*. *Enabledness* and *transition ability* to

669 occur are two distinct concepts: in the former, the transition is enabled if a proper number of tokens
670 are available from all but the “non-deterministic input places”, while the latter is bound to happen if it
671 has been enabled for timeout time without having taken place. In this case, tokens are withdrawn from
672 all but the “non-deterministic input places”; finally, the tokens are shifted only through the timeout
673 output arcs.

674 With the exception of Wyse’s work in [36] (where the author presents REACT, a tool for the
675 evaluation of approximate computing techniques), literature about approximated computing modeling
676 is relatively limited. To the best of our knowledge, we believe that the case study discussed in 3.2 is a
677 novel contribution in this area of research.

678 Modeling of the Map Reduce paradigm, on the other hand, reflects the business requirements
679 related to the need to develop a strategy capable of allocating the optimal number of resource by
680 minimizing, in parallel, the involved costs in deploying complex architectures. In this sense, Hadoop’s
681 performance analysis represents a crucial indicator that can be evaluated only through a valid model.
682 In contrast to the approach used in 3.3, in [37], the authors make use of JMT to develop simulation
683 models based on finite capacity region (FCR) deploying QN and Stochastic Well formed Nets (SWNs).
684 It is worth noting that some behaviors, such as the dynamic container allocation, are rather complex to
685 abstract, though the results achieve 9% of accuracy.

686 Multiformalism is a solid approach meant to augment the modeling capability of complex systems.
687 Both the case studies presented in section 4 can be compared to different techniques used to model
688 dynamic provisioning and e-healthcare scenarios discussed, though it has been noted that both consider
689 one single formalisms, reducing in this way, the degree of freedom of the modeler in choosing the
690 most suitable formalism to model a specific aspect of a problem. In [38], the authors deployed a
691 hybrid model consisting of an M/M/c model and multiple M/M/1 models to provision computing
692 resources within a virtualized application, while [39] presents a modeling architecture composed of an
693 M/M/1/B queue for each fog node with identical service time, a M/M/C queue with infinite waiting
694 buffer and a M/M/c/K queue characterizing each private node in the private cloud datacenter.

695 Future work is likely to follow two directions. Firstly, by providing users with more complex
696 metrics not directly obtainable from single formalisms. This would allow a better understanding of
697 the results produced by the models. For example, these metrics could be oriented to the study of a
698 system’s energy consumption or concerning the evaluation of the efficiency of complex algorithms
699 implementing approximate computing or genetic programming techniques.

700 Secondly, by comparing the single formalisms-based models with those using multiformalism
701 presented in this work, it is possible to note that the latter approach allows a greater elasticity on
702 modeling the dynamic aspects that take into consideration complex algorithms not based on system
703 variables (for example, this is not feasible within the QN paradigm: though the load depending routing
704 factor can be used, it is rather limited).

705 References

- 706 1. Economist. "The world’s most valuable resource is no longer oil, but data". [www.economist.com/
707 leaders/2017/05/06/the-worlds-most-valuable-resource-is-no-longer-oil-but-data](http://www.economist.com/leaders/2017/05/06/the-worlds-most-valuable-resource-is-no-longer-oil-but-data). [Online; accessed
708 15-January-2020].
- 709 2. Guardian. "Tech giants may be huge, but nothing matches big data". [https://www.theguardian.com/
710 technology/2013/aug/23/tech-giants-data](https://www.theguardian.com/technology/2013/aug/23/tech-giants-data), 2013. [Online; accessed 19-december-2019].
- 711 3. Flender, S. "Data is not the new oil". [https://towardsdatascience.com/data-is-not-the-new-oil-
712 bdb31f61bc2d](https://towardsdatascience.com/data-is-not-the-new-oil-bdb31f61bc2d), 2019. [Online; accessed 19-december-2019].
- 713 4. Ramadan, H.; Kashyap, D. "Quality of Service (QoS) in Cloud Computing. *International Journal of
714 Computer Science and Information Technologies*, 2017, pp. 318–320.
- 715 5. Bertoli, M.; Casale, G.; Serazzi, G. JMT: performance engineering tools for system modeling. *SIGMETRICS
716 Perform. Eval. Rev.* **2009**, *36*, 10–15. doi:http://doi.acm.org/10.1145/1530873.1530877.

- 717 6. Varghese, B.; Buyya, R. Next Generation Cloud Computing: New Trends and Research Directions. *Future*
718 *Generation Computer Systems* **2017**. doi:10.1016/j.future.2017.09.020.
- 719 7. Sajjad, M.; Ali, A.; Khan, A.S. Performance Evaluation of Cloud Computing Resources. *Performance*
720 *Evaluation* **2018**, *9*.
- 721 8. Duan, Q. Cloud service performance evaluation: status, challenges, and opportunities—a survey from the
722 system modeling perspective. *Digital Communications and Networks* **2017**, *3*, 101–111.
- 723 9. Maheshwari, S.; Raychaudhuri, D.; Seskar, I.; Bronzino, F. Scalability and performance evaluation of edge
724 cloud systems for latency constrained applications. 2018 IEEE/ACM Symposium on Edge Computing
725 (SEC). IEEE, 2018, pp. 286–299.
- 726 10. Calzarossa, M.C.; Massari, L.; Tessera, D. Workload Characterization: A Survey Revisited. *ACM Computing*
727 *Surveys* **2016**, *48*. doi:10.1145/2856127.
- 728 11. Megyesi, P.; Molnár, S. Analysis of Elephant Users in Broadband Network Traffic. 2013, Vol. 8115.
729 doi:10.1007/978-3-642-40552-5_4.
- 730 12. Casale, G.; Gribaudo, M.; Serazzi, G., Tools for Performance Evaluation of Computer Systems: Historical
731 Evolution and Perspectives. In *Performance Evaluation of Computer and Communication Systems. Milestones*
732 *and Future Challenges: IFIP WG 6.3/7.3 International Workshop, PERFORM 2010, Vienna, Austria, October 14-16,*
733 *2010, Revised Selected Papers*; Hummel, K.A.; Hlavacs, H.; Gansterer, W., Eds.; Springer: Berlin, Heidelberg,
734 2011; pp. 24–37. doi:10.1007/978-3-642-25575-5_3.
- 735 13. Ciardo, G.; Jones, III, R.L.; Miner, A.S.; Siminiceanu, R.I. Logic and stochastic modeling with SMART.
736 *Perform. Eval.* **2006**, *63*, 578–608. doi:10.1016/j.peva.2005.06.001.
- 737 14. Hillston, J. Tuning Systems: From Composition to Performance. *The Computer Journal* **2005**, *48*, 385–400.
738 doi:10.1093/comjnl/bxh097.
- 739 15. Hoare, C.A.R. *Communicating sequential processes*; Prentice-Hall, Inc.: Upper Saddle River, NJ, USA, 1985.
- 740 16. Milner, R. *Communication and concurrency*; Prentice-Hall, Inc.: Upper Saddle River, NJ, USA, 1989.
- 741 17. Sanders, W.H.; Courtney, T.; Deavours, D.; Daly, D.; Derisavi, S.; Lam, V. Multi-Formalism and
742 Multi-Solution-Method Modeling Frameworks: The Mobius Approach, 2007.
- 743 18. Vittorini, V.; Iacono, M.; Mazzocca, N.; Franceschinis, G. The OsMoSys approach to multi-formalism
744 modeling of systems. *Software and System Modeling* **2004**, *3*, 68–81.
- 745 19. Barbierato, E.; Gribaudo, M.; Iacono, M., Multiformalism and Multisolution Strategies for Systems
746 Performance Evaluation; 2015; pp. 201–222. doi:10.1002/9781119131151.ch8.
- 747 20. Barbierato, E.; Iacono, M.; Gribaudo, M. *Multiformalism and multisolution strategies for system performances*
748 *evaluation*; Prentice-Hall, Inc.: USA, 2015.
- 749 21. Khan, W.; Ahmed, E.; Hakak, S.; Yaqoob, I.; Ahmed, A. Edge computing: A survey. *Future Generation*
750 *Computer Systems* **2019**, *97*. doi:10.1016/j.future.2019.02.050.
- 751 22. Mouradian, C.; Naboulsi, D.; Yangui, S.; Glitho, R.; Morrow, M.; Polakos, P. A Comprehensive Survey on
752 Fog Computing: State-of-the-art and Research Challenges. *IEEE Communications Surveys & Tutorials* **2017**,
753 *PP*. doi:10.1109/COMST.2017.2771153.
- 754 23. Mao, Y.; You, C.; Zhang, J.; Huang, K.; Letaief, K.B. Mobile Edge Computing: Survey and Research Outlook.
755 *ArXiv* **2017**, *abs/1701.01090*.
- 756 24. Ajila, S.; Bankole, A. Using Machine Learning Algorithms for Cloud Client Prediction Models in a Web
757 VM Resource Provisioning Environment. *Transactions on Machine Learning and Artificial Intelligence* **2016**, *4*,
758 doi:10.14738/tmlai.41.1690.
- 759 25. Ardagna, D.; Barbierato, E.; Evangelinou, A.; Gianniti, E.; Gribaudo, M.; Pinto, T.B.M.; Guimarães, A.;
760 Couto da Silva, A.P.; Almeida, J.M. Performance Prediction of Cloud-Based Big Data Applications.
761 Proceedings of the 2018 ACM/SPEC International Conference on Performance Engineering; Association
762 for Computing Machinery: New York, NY, USA, 2018; ICPE '18, p. 192–199. doi:10.1145/3184407.3184420.
- 763 26. Didona, D.; Romano, P. Hybrid Machine Learning/Analytical Models for Performance Prediction: A
764 Tutorial. Proceedings of the 6th ACM/SPEC International Conference on Performance Engineering;
765 Association for Computing Machinery: New York, NY, USA, 2015; ICPE '15, p. 341–344.
766 doi:10.1145/2668930.2688823.
- 767 27. Conway, M.E. A Multiprocessor System Design. Proceedings of the November 12-14, 1963, Fall Joint
768 Computer Conference; Association for Computing Machinery: New York, NY, USA, 1963; AFIPS '63 (Fall),
769 p. 139–146. doi:10.1145/1463822.1463838.

- 770 28. Blumofe, R.D.; Leiserson, C.E. Scheduling Multithreaded Computations by Work Stealing. Proceedings of
771 the 35th Annual Symposium on Foundations of Computer Science; IEEE Computer Society: USA, 1994;
772 SFCS '94, p. 356–368. doi:10.1109/SFCS.1994.365680.
- 773 29. Arcari, L.; Gribaudo, M.; Palermo, G.; Serazzi, G. Performance-Driven Analysis for an Adaptive
774 Car-Navigation Service on HPC Systems. *SN Computer Science* **2020**, *1*. doi:10.1007/s42979-019-0035-7.
- 775 30. Chondrogiannis, T.; Bouros, P.; Gamper, J.; Leser, U. Alternative Routing: K-Shortest Paths with Limited
776 Overlap. Proceedings of the 23rd SIGSPATIAL International Conference on Advances in Geographic
777 Information Systems; Association for Computing Machinery: New York, NY, USA, 2015; SIGSPATIAL '15.
778 doi:10.1145/2820783.2820858.
- 779 31. Balbo, G., Introduction to Generalized Stochastic Petri Nets. In *Formal Methods for Performance Evaluation: 7th International School on Formal Methods for the Design of Computer, Communication, and Software Systems, SFM 2007, Bertinoro, Italy, May 28-June 2, 2007, Advanced Lectures*; Bernardo, M.; Hillston, J., Eds.; Springer
780 Berlin Heidelberg: Berlin, Heidelberg, 2007; pp. 83–131. doi:10.1007/978-3-540-72522-0_3.
- 781 32. Pinciroli, R.; Gribaudo, M.; Roveri, M.; Serazzi, G., Capacity Planning of Fog Computing Infrastructures
782 for Smart Monitoring; 2018; pp. 72–81. doi:10.1007/978-3-319-91632-3_6.
- 783 33. Kleinrock, L. Power and Deterministic Rules of Thumb for Probabilistic Problems in Computer
784 Communications. Conference Record, International Conference on Communications; , 1979; pp.
785 43.1.1–43.1.10.
- 786 34. Holvoet, T.; Verbaeten, P., Using Petri Nets for Specifying Active Objects and Generative Communication;
787 2001; pp. 38–72. doi:10.1007/3-540-45397-0_2.
- 788 35. Jensen, K. Coloured Petri Nets. Petri Nets: Central Models and Their Properties; Brauer, W.; Reisig, W.;
789 Rozenberg, G., Eds.; Springer Berlin Heidelberg: Berlin, Heidelberg, 1987; pp. 248–299.
- 790 36. Wyse, M. Modeling Approximate Computing Techniques. 2015.
- 791 37. Bernardi, S.; Gianniti, E.; Aliabadi, S.; Perez-Palacin, D.; Requeno, J. Modeling Performance of Hadoop
792 Applications: A Journey from Queueing Networks to Stochastic Well Formed Nets. 2016, Vol. 10048, pp.
793 599–613. doi:10.1007/978-3-319-49583-5_47.
- 794 38. Bi, J.; Zhu, Z.; Tian, R.; Wang, Q. Dynamic Provisioning Modeling for Virtualized Multi-tier Applications
795 in Cloud Data Center. 2010, pp. 370–377. doi:10.1109/CLOUD.2010.53.
- 796 39. El Kafhali, S.; Salah, K. Performance Modeling and Analysis of Internet of Things enabled Healthcare
797 Monitoring Systems **2019**. *8*, 48–58. doi:10.1049/iet-net.2018.5067.

800 © 2020 by the authors. Submitted to *Future Internet* for possible open access publication
801 under the terms and conditions of the Creative Commons Attribution (CC BY) license
802 (<http://creativecommons.org/licenses/by/4.0/>).