

QoS Driven Per-Request Load-Aware Service Selection in Service Oriented Architectures

Valeria Cardellini, Valerio Di Valerio, Vincenzo Grassi,

Stefano Iannucci, and Francesco Lo Presti

(Dipartimento di Ingegneria Civile e Ingegneria Informatica, Università di Roma “Tor Vergata”
Via del Politecnico 1, 00133 Roma, Italy)

Abstract Service selection has been widely investigated by the SOA research community as an effective adaptation mechanism that allows a service broker, offering a composite service, to bind at runtime each task of the composite service to a corresponding concrete implementation, selecting it from a set of candidates which differ from one another in terms of QoS parameters. In this paper we present a load-aware per-request approach to service selection which aims to combine the relative benefits of the well known per-request and per-flow approaches. Our service selection policy represents the core methodology of the Plan phase of a self-adaptive service oriented system based on the MAPE-K reference loop. Since the service broker operates in a variable and uncertain environment where the QoS levels negotiated with the service providers can fluctuate, it requires some mechanism to enforce the QoS constraints with its users. To this end, we also propose an algorithm for the Analyze phase of MAPE-K which is based on the adaptive Cusum algorithm and allows to determine whether a change in the QoS level requires a service selection replanning. We present experimental results obtained with a prototype implementation of a service broker. Our results show that the proposed load-aware approach is superior to the traditional per-request one and combines the ability of sustaining large volume of service requests, as the per-flow approach, while at the same time offering a finer customizable service selection, as the per-request approach. Furthermore, the results show that the adaptive Cusum algorithm can quickly detect changes in the execution environment and trigger a new optimization plan before the system performance degrades.

Key words: quality of service; service oriented architecture; service selection

Cardellini V, Di Valerio V, Grassi V, Iannucci S, Lo Presti F. QoS driven per-request load-aware service selection in service oriented architectures. *Int J Software Informatics*, Vol.7, No.2 (2013): 195–220. <http://www.ijsi.org/1673-7288/7/i1156.htm>

1 Introduction

Service Oriented Systems (SOSs) are becoming popular thanks to a widely deployed internetworking infrastructure. SOSs are composed by a possibly large number of heterogeneous third party subsystems. As a consequence, as they grow in number and size, they also rapidly increase in complexity, which is further complicated by the highly changing execution environment where they have to

operate. A major trend to tackle this growing complexity is to design SOSs as runtime self-adaptable systems, to make them able to meet both functional and non functional requirements. The former concern the overall logic to be implemented, while the latter concern the Quality of Service (QoS) levels that should be guaranteed to the SOSs users. In both cases, self-adaptation can leverage different mechanisms, including the tuning of control parameters of services used in the composition of a SOS, the selection of the most suitable services within a set of candidates, or even a modification of the overall SOS composition logic. In this paper, we focus on self-adaptation based on QoS driven *service selection*. The service selection goal is to determine the binding of each *abstract task* in the composite service to actual implementations, leaving unchanged the composition logic. The idea at the basis of service selection is to exploit the existence, in the open marketplace as well as in the proprietary service parks, of several services, referred to as *concrete services*, implementing the same functionality with different non functional characteristics and cost^[14, 21, 25].

The service selection problem has been widely investigated in recent years. A first generation of service selection solutions implements a *local* approach^[23, 28], that time by time associates each running task of a SOS with the best available service that implements that task. However, this local approach can guarantee only local QoS constraints, for example the response time of a given task lower than a given threshold. Second generation solutions implement a *global* approach, where the QoS constraints are guaranteed for the whole execution of the composite SOS rather than for its single tasks. They are more suitable in a scenario where a service provider stipulates global Service Level Agreements (SLAs) with users.

Global approaches face the service selection problem at two granularity levels. At the *per-request* grain^[3, 4, 10, 17, 19, 27], the adaptation focuses on each single request submitted to the system and aims at fulfilling the QoS constraints of that specific request. On the contrary, the *per-flow* grain^[2, 7, 9, 16] considers the flow of requests of a user rather than the single request, and the adaptation goal is to fulfill the QoS constraints that concern the global properties of that flow, *e.g.*, the average SOS response time or its availability.

However, the solutions proposed so far for both per-request and per-flow granularities are not satisfactory, either in terms of QoS guarantees or scalability to user requests. The per-request grain exhibits scalability problem under a sustained traffic of requests, because each request is managed independently of all the other concurrent ones. As a consequence, multiple service requests could be assigned to the same concrete service, that could be overloaded. On the other hand, the per-flow grain is not able to ensure QoS guarantees to a single request, and the user perceived QoS could be very different from that stipulated in the SLA.

To overcome these limitations, as a first contribution we propose a new per-request service selection policy, that we call *load-aware per-request*. The proposed policy exploits the multiple available implementations of each abstract task, and realizes a runtime probabilistic binding. In this way, different concurrent requests to the same abstract task are bound to different concrete services, thus realizing a randomized load balancing similarly to the per-flow solutions^[7]. At the same time, however, the QoS constraints are ensured *for each request* that the user

submits.

However, service selection is only one component that is needed to build a self-adaptable system. It allows to reconfigure the system, but the system must also be able to autonomously recognize when a reconfiguration is required, i.e., to understand when the execution environment has changed. The latter is a particularly critical task when a SOA application based on third party services needs to fulfill non-functional requirements, because existing services may disappear or their performance may quickly fluctuate over time. If the reconfiguration trigger is not properly raised, the SOS may no longer be able to fulfill the QoS constraints because the computed service selection becomes almost useless. The MAPE-K loop^[15], a general framework to build self-adaptable systems, identifies four logical essential components, also called phases, that are needed to build a self-adaptable system: Monitor, Analyze, Plan and Execute. This model is based on a feedback-control loop that monitors the execution environment, analyzes the collected data in order to detect changes, plans the necessary actions to reconfigure the system maximizing at the same time some utility function, and executes these actions. In the particular context of SOSs, realizing an MAPE-K loop poses several challenges which range from designing an effective planning mechanism to implementing a methodology able to quickly detect changes in the execution environment, coupled with an efficient methodology for data collection.

As a second contribution, this paper presents a theoretical framework for a self-adaptable SOS that attempts to address all the above challenges. Specifically, we propose:

- a methodology to efficiently collect data from the execution environment;
- an analysis algorithm based on the adaptive Cusum algorithm proposed in Ref. [13];
- a new service selection policy that plans the composite service reconfiguration.

Furthermore, as a third contribution, we provide a fully working implementation of the proposed framework. We plugged it into the MOSES prototype^[6, 9] and performed an extensive experimental evaluation. The results show the effectiveness of our proposal:

- the load-aware per-request policy, due to its inherent scalability, is more suitable to work in a realistic environment than the service selection policy presented in Ref. [3], which is one of the per-request top performing state-of-the-art approaches;
- the adaptive Cusum algorithm is able to quickly detect changes in the execution environment, thus allowing to trigger a new optimization before the SOS performance degrades.

The rest of this paper is organized as follows. In Section 2 we introduce the system model, the QoS model of the SOS, and the basic idea of randomized load balancing. In Section 3 we present the MILP optimization problem that determines the optimal selection policy, while in Section 4 we describe the Monitor

instrumentation and the Analyze algorithm. In Section 5, we first describe the MOSES prototype and then demonstrate the effectiveness of our approach through a set of experiments. An overview of the service selection policies proposed so far and of some self-adaptive frameworks that include the Monitor and Analyze components is discussed in Section 6, while the conclusions are drawn in Section 7.

2 System Model

We consider a broker that offers to the users a composite service P with different QoS levels and monetary prices, exploiting for this purpose a set of existing concrete services. Hence, the broker is an intermediary between users and concrete services, acting as a service provider towards the users and as a service requestor towards the concrete services used to implement the composite service, as shown in Fig. 1. Its main task is to drive the adaptation of the service it manages to fulfill the Service Level Agreements (SLAs) negotiated with its users, given the SLAs it has negotiated with the concrete services and while optimizing a suitable broker utility function, *i.e.*, response time or cost.

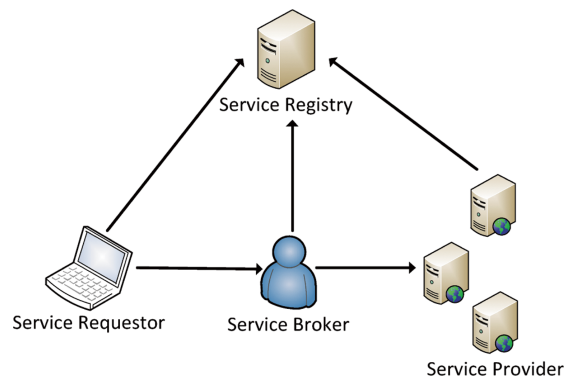


Figure 1. SOA architecture with service broker

Let us denote by $S_i \in \mathcal{S}, i = 1, \dots, m$ the set of abstract tasks that compose the composite service P (for sake of simplicity, in the following we will use S_i and i interchangeably) and by \mathfrak{S}_i the set of concrete services implementing the abstract task S_i . Within this framework, a core task of the service broker is to determine a proper service selection, that is to find for each abstract task S_i a corresponding implementation $cs_{ij} \in \mathfrak{S}_i$, so that the agreed SLAs are fulfilled. The selection criterion corresponds to the optimization of a given utility goal of the broker.

Within this context, the load-aware per-request selection policy realizes a probabilistic binding: for each abstract task, it determines a set of concrete services that can be used to implement that task, with a probability associated with each of these services. Therefore, different concrete services can be bound to the same abstract task, depending on these probabilities. Thus, the actual implementation of a composite service may vary from request to request.

The probabilities used by our selection policy are determined by taking into account both the number of concurrent requests that arrive to the system and the capacity of each concrete service, that we assume known to the service broker (we will

address this issue in Section 5). This probabilistic binding represents a peculiarity of our per-request approach with respect to other service selection policies with the same granularity. It allows us to realize a randomized load balancing which is aware of the capacity of each concrete service as well as of the system load, thus overcoming the limits of the per-request selection policies proposed so far. In fact, in Refs. [3, 4, 10, 28], each abstract task is deterministically bound to only one concrete service that can be overloaded under high load conditions. By exploiting multiple service implementations, our per-request solution is able to share the load.

However, the service broker has to satisfy the QoS levels agreed in the SLAs with its users. Since different concrete services implementing the same task can differ in their QoS attributes values, we may obtain different composite service implementations with distinct overall QoS attributes, depending on the actual service binding. To guarantee the QoS constraints, the services must be chosen so that all the feasible implementations respect them.

The rest of this section is organized as follows. In Section 2.1 we introduce the SLA model and the QoS attributes. In Section 2.2 we explain how the randomized load balancing is realized. Finally, in Section 2.3 we illustrate how the QoS attributes are calculated along the composite service.

2.1 SLA model

The QoS levels offered by the service providers are defined in SLAs. Since our selection policy aims at satisfying every single request, the SLA states conditions that are restricted to the single request. In general, SLA conditions may refer to different kinds of functional and non-functional attributes of the service. We consider in this work the following attributes:

- *response time*: the upper bound on the interval of time elapsed from the service invocation to its completion;
- *availability*: the lower bound on the probability that the service is accessible when invoked;
- *cost*: the upper bound on the price charged for the service invocation.

We model the SLA between the service broker and its users as a tuple $\langle R_{\max}, A_{\min}, C_{\max} \rangle$, where R_{\max} and C_{\max} represent the upper bound on the response time and cost, respectively, and A_{\min} the lower bound on the availability. Since the service broker plays an intermediary role, in its turn it also acts as a service user against the providers of the concrete services it uses to implement the abstract tasks in the service composition. Each of these concrete services is characterized by its own SLA, that we model as an extension of the SLA offered by the broker to its users: response time, availability, and cost maintain the same semantics, but we add a new parameter, L_{ij} , which is a threshold on the maximum request rate that can be submitted to the concrete service cs_{ij} . The cs_{ij} provider satisfies the QoS attributes in the advertised SLA as long as the request rate does not exceed L_{ij} ; in case of violation of the load threshold, there is no guarantee on the QoS levels of cs_{ij} .

2.2 Randomized load balancing

The core of our per-request selection policy is the randomized load balancing of the requests directed to each abstract task S_i , so that they are switched to multiple concrete services implementing it. The load balancing is tuned on the basis of the capacity of each concrete service cs_{ij} , defined by the parameter L_{ij} , and of the rate of requests submitted by the users to S_i . As we already mentioned, an abstract task is bound by the broker to a set of concrete services, each one having an associated probability. So, at abstract task binding time, only one of these services is probabilistically chosen. As a consequence, only a fraction of the incoming requests is switched to a given concrete service cs_{ij} , and this fraction depends on the probability x_{ij} determined by the broker. We define a service selection policy as the set of all these probabilities, that we represent with the vector $\mathbf{x} = [x_1, \dots, x_m]$, where for each entry $\mathbf{x}_i = [x_{ij}]$, $i \in \mathcal{S}$, $j \in \mathfrak{S}_i$, the constraints $x_{ij} \in [0, 1]$ and $\sum_{j \in \mathfrak{S}_i} x_{ij} = 1$ hold. Our idea is to drive the value of the x_{ij} probabilities, forcing on each x_{ij} an upper bound P_{ij} , so that the fraction of requests switched by the broker to the concrete service cs_{ij} does not overload it. The upper bound P_{ij} is calculated through the ratio $P_{ij} = \frac{L_{ij}}{\lambda_i}$, where λ_i is the actual request rate to the abstract task S_i and L_{ij} is the load threshold for cs_{ij} . If P_{ij} is greater than 1, it means that there is no upper bound because cs_{ij} is able to satisfy all the incoming requests to S_i on its own. Vice versa, if P_{ij} is less than 1, cs_{ij} alone cannot satisfy all the requests directed to S_i but it must be backed by other concrete services, so that their overall capacity can sustain the submitted load.

2.3 QoS attributes of the composite service

The QoS attributes of the composite service can be calculated using the service selection policy \mathbf{x} , but a model of the composite service workflow is also needed. In the following subsections we introduce this model and, afterwards, the QoS attributes computation.

2.3.1. Composite service graph

We assume that a composite service workflow managed by the service broker either has a single initial task or starts with a fork-join parallel sequence. Furthermore, we assume that for each conditional branch we know the probability of executing it; similarly, we assume to know the probability of reiterating loops.

The composite service graph is obtained by transforming the workflow of the composite service as in Ref. [3]. In particular, loops are peeled, *i.e.*, they are transformed in a sequence of branch conditions, each of which evaluates if the loop has to continue with the next iteration or it has to exit, according to the branch probability. A pre-requisite for loop peeling is the knowledge of the maximum number of supposed iterations. We calculate this value as the *p*-percentile of the distribution of reiterating the loop. After loop peeling, the composite service can be modeled as a Directed Acyclic Graph (DAG). As in Ref. [3], we define:

- *Execution path.* An execution path ep_n is a multiset of tasks $ep_n = \{S_1, S_2, \dots, S_I\} \subseteq \mathcal{S}$, such that S_1 and S_I are respectively the initial and final tasks of the path and no pair $S_i, S_j \in ep_n$ belongs to alternative

branches. We need a multiset rather than a simple set because a single task may appear several times in the execution path. An execution path may also contain parallel sequences but no loops, which are peeled. A probability of execution p_n is associated with every execution path and can be calculated as the product of the probabilities of executing the branch conditions included in the path. Similarly, the branch conditions that arise from loop peeling produce other execution paths.

- *Subpath.* A subpath of an execution path ep_n is a sequence of tasks $[S_1, \dots, S_I]$, from the initial to the end task, that does not contain any parallel sequence. In other words, each branch b of a parallel sequence identifies a subpath inside the execution path ep_n . We denote a subpath by sp_b^n .

Therefore, the set of all the execution paths identifies all the possible execution scenarios of the composite service. The QoS constraints must hold for every execution path to guarantee the SLAs the service broker stipulated with its users.

2.3.2. QoS attributes computation

Given the service selection policy \mathbf{x} and the execution paths that arise from the composition logic, we can calculate the QoS attributes of each abstract task and then the overall QoS attributes of the composite service. We are interested in the average QoS perceived by the users as well as in its worst case value. As discussed below, we need both these values to maximize the broker utility function and satisfy the QoS constraints.

Let r_{ij} be the response time of the concrete service cs_{ij} , a_{ij} its availability, and c_{ij} its cost. The average QoS values of the abstract task S_i , namely, the average response time R_i , the availability A_i , and the average cost C_i , are given by the following expressions:

$$R_i = \sum_{j \in \mathfrak{S}_i} r_{ij} x_{ij} \quad (1)$$

$$A_i = \sum_{j \in \mathfrak{S}_i} a_{ij} x_{ij} \quad (2)$$

$$C_i = \sum_{j \in \mathfrak{S}_i} c_{ij} x_{ij} \quad (3)$$

The worst case QoS values, denoted by R_i^w , A_i^w , and C_i^w , are given by:

$$R_i^w = \max_{j \in \mathfrak{S}_i} r_{ij} y_{ij} \quad (4)$$

$$A_i^w = \min_{j \in \mathfrak{S}_i} a_{ij} y_{ij} \quad (5)$$

$$C_i^w = \max_{j \in \mathfrak{S}_i} c_{ij} y_{ij} \quad (6)$$

where y_{ij} is a binary variable indicating whether the concrete service cs_{ij} can be ever bound to the abstract task S_i , *i.e.*, $y_{ij} = 1$ if $x_{ij} > 0$ and 0 otherwise.

Using these formulas and the notion of execution paths and subpaths, we can calculate the QoS attributes along each execution path itself, using the aggregation

formulas presented in Ref. [3]. Note that the same formulas apply both for the average case and for the worst case; therefore, for the sake of simplicity, we show only the latter case. We denote by R_n the maximum response time of the execution path ep_n , with A_n its minimum availability, and with C_n its maximum cost; these are, in other words, the QoS attributes values calculated in the worst case scenario. They are:

$$R_n = \max_{sp_b^n \in ep_n} \sum_{S_i \in sp_b^n} R_i^w \quad (7)$$

$$A_n = \prod_{S_i \in ep_n} A_i^w \quad (8)$$

$$C_n = \sum_{S_i \in ep_n} C_i^w \quad (9)$$

While the cost and the availability are simply obtained, respectively, as sum and multiplication of the QoS attributes of each abstract task in the execution path, the matter is slightly different for the response time. Indeed, the response time of an execution path is equal to the response time of the longest subpath inside the execution path itself.

3 Optimization Problem

Given a composite service P , the goal of the service broker is to find a selection policy \mathbf{x} that ensures the QoS constraints for every execution scenario, *i.e.*, for each execution path ep_n that arises from P , while realizing the randomized load balancing. The selection policy \mathbf{x} is calculated by solving a suitable optimization problem. We formulate this optimization problem as a Mixed Integer Linear Problem (MILP), with the following decision variables:

- x_{ij} : it takes value in the range $[0, 1]$ and represents the probability that the concrete service $cs_{ij} \in \mathfrak{S}_i$ is bound to the abstract task S_i ; it is used to drive the randomized load balancing.
- y_{ij} : it is equal to 1 if cs_{ij} is bound to S_i with a given probability defined by x_{ij} , 0 otherwise. We use it to ensure that the QoS constraints are met.

While the QoS constraints are evaluated using the worst case values of the QoS attributes for each abstract task, the objective function is maximized using the average values, because it is the value that is expected along multiple executions of the composite service. In particular, the optimization problem maximizes the aggregated QoS values, which are calculated over all the possible execution paths that arise from the composite service workflow, taking into account the relative probability p_n . We obtain the aggregated values by applying the Simple Additive Weighting (SAW) technique as scalarization method.

In the first phase, each quality dimension along an execution path is normalized according to the following formulas, depending on whether the QoS attribute is a positive (10) or a negative (11) one. A QoS attribute is defined positive (negative) if the greater the value is, the greater (lower) the quality of that attribute. Availability is an example of positive attribute (the higher the availability, the better the quality

is), while response time is an example of negative attribute (the lower the response time, the better the quality is).

$$z_n^h(\mathbf{x}) = \begin{cases} \frac{q_n^h(\mathbf{x}) - \min q_n^h}{\max q_n^h - \min q_n^h}, & \text{if } \max q_n^h \neq \min q_n^h \\ 1, & \text{if } \max q_n^h = \min q_n^h \end{cases} \quad (10)$$

$$z_n^h(\mathbf{x}) = \begin{cases} \frac{\max q_n^h - q_n^h(\mathbf{x})}{\max q_n^h - \min q_n^h}, & \text{if } \max q_n^h \neq \min q_n^h \\ 1, & \text{if } \max q_n^h = \min q_n^h \end{cases} \quad (11)$$

In the above formulas, $q_n^h(\mathbf{x})$ is the h -th quality dimension value calculated over the execution path ep_n using the selection policy \mathbf{x} . $\max q_n^h$ and $\min q_n^h$ are its maximum and minimum values and can be estimated across several composite service executions.

In the second phase a score is obtained using a weighted sum of the normalized quality attributes, as follows:

$$score_n = \sum_h w_h z_n^h(\mathbf{x}) \quad (12)$$

where the weight w_h specifies the relative importance that the broker assigns to a QoS attribute with respect to the others.

Finally, the objective function is obtained using the following weighted formula:

$$F(\mathbf{x}) = \sum_n p_n score_n(\mathbf{x}) \quad (13)$$

The optimal service selection policy \mathbf{x} can be obtained solving the following optimization problem (for sake of simplicity, we use n instead of ep_n):

$$\begin{aligned} & \mathbf{max} \quad F(\mathbf{x}) \\ \mathbf{subject \ to:} & \quad \sum_{j \in \mathfrak{S}_i} x_{ij} = 1 \quad \forall i \end{aligned} \quad (14)$$

$$x_{ij} \leq P_{ij} \quad \forall i, \forall cs_{ij} \in \mathfrak{S}_i \quad (15)$$

$$x_{ij} \leq y_{ij} \quad \forall i, \forall cs_{ij} \in \mathfrak{S}_i \quad (16)$$

$$r_{ij} y_{ij} \leq R_i^w \quad \forall i, \forall cs_{ij} \in \mathfrak{S}_i \quad (17)$$

$$a_{ij} y_{ij} \geq A_i^w \quad \forall i, \forall cs_{ij} \in \mathfrak{S}_i \quad (18)$$

$$c_{ij} y_{ij} \leq C_i^w \quad \forall i, \forall cs_{ij} \in \mathfrak{S}_i \quad (19)$$

$$\sum_{i \in sp_b^n} R_i^w \leq R_n \quad \forall sp_b^n \in ep_n, \forall n \quad (20)$$

$$\sum_{i \in ep_n} \log(A_i^w) = A_n \quad \forall n \quad (21)$$

$$\sum_{i \in ep_n} C_i^w = C_n \quad \forall n \quad (22)$$

$$R_n \leq R_{\max} \quad \forall n \quad (23)$$

$$A_n \geq \log(A_{\min}) \quad \forall n \quad (24)$$

$$C_n \leq C_{\max} \quad \forall n \quad (25)$$

$$\begin{aligned}
x_{ij} &\in \mathbb{R}^+ && \forall i, \forall cs_{ij} \in \mathfrak{S}_i \\
y_{ij} &\in \{0, 1\} && \forall i, \forall cs_{ij} \in \mathfrak{S}_i \\
R_i^w, A_i^w, C_i^w &\in \mathbb{R}^+ && \forall i \\
R_n, C_n &\in \mathbb{R}^+ && \forall n \\
A_n &\in \mathbb{R}^- && \forall n
\end{aligned}$$

Constraints (14) guarantee that the sum of the probabilities of choosing the concrete services is equal to 1 for each task S_i . Constraints (15) define the upper bound to the probability of choosing a concrete service cs_{ij} . These two constraints families implement the randomized load balancing policy. Constraints from (17) to (19) express the response time, availability, and cost of every abstract task S_i in terms of the worst concrete services that are selected to implement that task, as we discussed in Section 2.3. Constraints (20) evaluate the response time of each execution path ep_n as the response time of its longest subpath sp_b^n , while constraints (21) and (22) refer respectively to the availability and cost of each execution path ep_n . Finally, constraints from (23) to (25) are the QoS constraints to be fulfilled. We use the logarithm of the availability instead of the availability in our optimization problem because we need to linearize Equation (8) to put it in our MILP problem.

4 Monitoring and Analysis

In the previous sections, we have presented the load-aware per-request service selection policy that guarantees the QoS constraints to each user request by properly composing the available concrete services. Nevertheless, since these services are offered by third party providers, even in presence of SLAs there is no actual guarantee that the services abide to the negotiated QoS parameters as network overload, service overload, and/or power outages may cause a service to not respect the expected QoS level. As a consequence, to provide QoS guarantees, we have to account for the actual services QoS, rather than the QoS values stated in the SLAs. To this end, we need to monitor the execution of the concrete services and to analyze the collected data, so to be able to determine whether a change in the QoS level has occurred and a new service selection needs to be determined by using the updated QoS parameters.

In the following, we discuss the methodologies used to monitor the concrete services and to analyze the collected data focusing on the response time, while in Section 5 we demonstrate their effectiveness in improving the QoS experienced by the composite service users.

4.1 Monitoring the concrete services

The QoS attributes stated in a SLA are the targets of our monitoring activity. The monitored data can be collected at two different locations: at the service provider side or at the broker side. The former is made possible when the service provider collects data for itself and makes them available to its clients, like the Amazon CloudWatch service. However, the most common solution adopted in the SOA context is to collect data on the service broker that manages the composite service^[1, 6, 11, 20], because the monitoring service is hardly provided by the concrete

services providers.

Generally, the methodology used to collect the data can be either active, if the data are collected sending proper inputs to the monitored entities, or passive, if the data are collected without injecting additional load but rather observing the system behavior. We preferred the latter solution, because in the context of the SOA applications each service invocation has a cost. Another important question regards the frequency at which data are collected, i.e., after how many invocations of a service we measure the QoS. Clearly, a low frequency approach requires less computing power than a high frequency one, but we adopt the latter to react more quickly to changes in the service QoS. Therefore, the monitoring activity of our framework is straight: we passively measure and store the QoS of each single concrete service invocation on a continuous time basis.

4.2 Analyzing the collected data

To analyze the collected data, we propose to use the online adaptive cumulative sum (Cusum) algorithm^[22] for service response time monitoring and abrupt change detection. We chose not to use the standard Cusum algorithm because it has been designed to detect changes in stationary time series with known statistical characteristics. In a non-stationary context, like the SOA context, where the variance can exhibit significant variation over time, standard Cusum performance is quite poor^[13]. Indeed, given a time series, its standard deviation is used to properly tune the Cusum algorithm. Furthermore, although Cusum detects changes in the time series mean, it assumes that the variance is always constant over time.

The online adaptive Cusum detector we implemented^[13] combines an Exponential Weighted Moving Average (EWMA) filter to tracks the slow varying response time series average with a two-sided Cusum test to detect abrupt changes in the series average which cannot be timely accounted by EWMA filter.

We consider the following tracking EMWA filter:

$$\mu_i = \alpha y_i + (1 - \alpha)\mu_{i-1} \quad (26)$$

where μ_i represents the i -th current estimate of the average response time, y_i the i -th collected response time sample, and α is a small constant. To timely detect the occurrence of significant changes, the Cusum algorithm uses two variables, g^+ and g^- , for positive and negative changes, respectively, which measure positive and negative deviation of the time series with respect to its average value. They are initialized to 0 and updated at each step as follows:

$$g_i^+ = \max\{0, g_{i-1}^+ + y_i - (\mu_i + K^+)\} \quad (27)$$

$$g_i^- = \max\{0, g_{i-1}^- + (\mu_i - K^-) - y_i\} \quad (28)$$

where K^+ (K^-) is the smallest shift we want to detect on the leading (trailing) edge. In our experiments, we set it equal to 25% of the response time stated in the SLOs of the monitored services. A change is detected whenever g_i^+ or g_i^- are greater than a suitable threshold H^* , which represents a trade-off between detection delay and probability of false positive. To compute H^* we followed the approach in Ref. [13], which ensures a small probability of false detection (measured in terms of

expected number of samples between false positives, we set to 1000). This requires the numerical inversion of the Siegmund approximation^[22] which typically yields $H^* \approx 5\sigma_y$, where σ_y is the time series standard deviation. Since σ_y is unknown, we resort to a widely adopted approximation, which basically replaces the standard deviation with the estimate of the mean deviation $E[|y_i - E[y]|]$:

$$\sigma_i = \beta|y_i - \mu_i| + (1 - \beta)\sigma_{i-1} \quad (29)$$

where we set β to 0.5.

Whenever an abrupt change is detected, the average response time is updated according to the following equations^[22]:

$$\mu_i = \begin{cases} \mu_{i-1} + K + g_i^+ / N^+ & \text{if } g_i^+ > H^* \\ \mu_{i-1} - K - g_i^- / N^- & \text{if } g_i^- > H^* \end{cases} \quad (30)$$

in place of (26), where N^+ (N^-) is the number of samples since the last time g_i^+ (g_i^-) was equal to zero. Upon a change detection, g_i^+ and g_i^- are reset to 0.

The estimates of the services response time, obtained via (26) (or (30) when an abrupt change is detected), are used by the broker in the actual formulation of the service selection optimization problem. More precisely, since we want to guarantee an upper bound on the composite service response time, in the problem formulation for each concrete service cs_{ij} the broker uses the largest value between r_{ij} , defined in the SLA, and the current estimate computed above, which we denote by r_{ij}^{CUSUM} , *i.e.*, in (17) we replace r_{ij} by $\tilde{r}_{ij} = \max\{r_{ij}, r_{ij}^{CUSUM}\}$.

5 Experimental Analysis

In this section, we present the experimental analysis we have conducted using the MOSES prototype to demonstrate:

- the effectiveness of the proposed load-aware per-request approach with respect to the traditional per-request approach proposed by Ardagna and Pernici in Ref. [3];
- the effectiveness of the adaptive Cusum algorithm, and in turn of the whole theoretical framework for a self-adaptive SOS.

First, we briefly describe the MOSES prototype and the experimental environment, then we analyze the experimental results.

5.1 MOSES prototype

MOSES, which stands for *MOdel-based SElf adaptation of SOA systems*, is a QoS-driven runtime adaptation framework for service-oriented systems^[9]. It is intended to act as a *service broker* which offers to prospective users a composite service with QoS guarantees, exploiting for this purpose the runtime binding between the abstract functionalities of the composite service and a pool of existing concrete services that implement the abstract functionalities. Its main task is to drive the adaptation of the composite service to fulfill the QoS goals stipulated with its users.

The MOSES prototype has a flexible and modular system architecture, where each module performs a specific functionality and its implementation can be changed

without impacting on the rest of the system. In the following, we provide an overview on the MOSES system; a detailed description is in^[6, 9].

We first describe the core MOSES modules (namely, the *BPEL Engine*, the *Adaptation Manager*, and the *Optimization Engine*, together with the *Storage* layer); then, we present the remaining modules that enrich the basic functionalities.

The *Optimization Engine* computes the selection policy α that drives the runtime binding. The modular architecture of MOSES allows us to develop multiple implementations of the service selection optimization policies, possibly using external tools for finding the optimal solution. To this end, the Optimization Engine exposes the same interface to other MOSES modules irrespectively of the specific policy. Examples of already implemented policies are the per-flow optimization policies in Refs. [7, 8] and the per-request optimization policy in Ref. [3], as well as the load-aware selection policy presented in this paper.

The *BPEL Engine* executes the composite service, described in BPEL^[24], that defines the user-relevant business logic.

The *Adaptation Manager* is the actuator of the adaptation actions determined by the Optimization Engine: it is actually a proxy interposed between the BPEL Engine and any external service provider. Its task is to dynamically bind each abstract task's invocation to the real endpoint according to the optimal solution determined by the Optimization Engine.

MOSES is architected as a self-adaptive system based on the MAPE-K (Monitor, Analyze, Plan, Execute, and Knowledge) reference model for autonomic systems^[26]. Figure 2 shows how the MOSES modules implement each MAPE-K macro-component, together with the system inputs (i.e., the composite service and the pool of candidate concrete services). This input is used to build a model (Execute), which in turn is used and is kept up-to-date at runtime (Monitor). The monitored parameters are analyzed (Analyze) in order to know if adaptation actions (i.e., a new service selection policy) have to be taken. A new selection policy α is calculated (Plan) to react to some external significant event, such as a significant

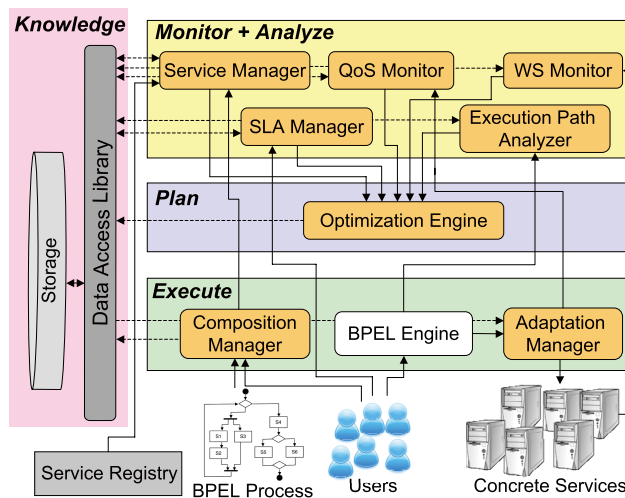


Figure 2. MOSES high-level architecture based on the MAPE-K model

violation of the SLA parameters. The BPEL Engine, together with the Adaptation Manager, belong to the Execute macro-component because their task is to execute the logic of the business processes; on the other hand, the Optimization Engine constitutes the Plan macro-component because it is involved in computing the selection policy.

The basic functionalities implemented in the Execute and Plan macro-components are enriched by the modules belonging to the Monitor+Analyze macro-component: they capture changes in the MOSES environment and, if they are relevant, modify at runtime the system model kept in the Storage layer, also triggering the Optimization Engine to make it calculate a new service selection policy. The new service selection policy α will be then calculated using the system model view as updated by the monitoring modules. Specifically, the *Service Manager* and *WS Monitor* respectively detect the addition or removal of concrete services (the latter due either to graceful failures or crashes). The *QoS Monitor* measures the actual values of the QoS attributes provided by MOSES to its users and offered to MOSES by its service providers and detects violations of the service level objectives stated in the SLAs. It implements the adaptive Cusum algorithm presented in Section 4 and triggers the Optimization Engine whenever a change detection occurs. The *Execution Path Analyzer* tracks variations in the usage profile of the abstract tasks, allowing for example to dynamically update the probability of executing the conditional branches in the composite service workflow. Finally, the *SLA Manager* manages the user registration with the associated SLA, possibly performing an admission control.

We have implemented the MOSES prototype exploiting the Java Business Integration (JBI) implementation provided by OpenESB and MySQL for the storage layer. We have used Sun BPEL Service Engine to orchestrate the service composition. The Optimization Engine relies on IBM ILOG CPLEX Optimizer^[12] as optimization software package to solve the per-request optimization problems.

5.2 Experimental setup

The testing environment consists of 3 Intel Xeon quad-core servers (2 Ghz/core) with 8 GB RAM each (nodes 1, 2, and 3), and 1 KVM virtual machine with 1 CPU and 1 GB RAM (node 4); a Gb Ethernet connects all the machines. The MOSES prototype is deployed as follows: node 1 hosts all the MOSES modules in the Execute macro-component, node 2 the storage layer together with the candidate concrete services, and node 3 the modules in the Monitor+Analyze and Plan macro-components. Finally, node 4 hosts the workload generator.

We consider the composite service defined by the workflow in Fig. 3, composed by 6 stateless tasks, and assume that 4 concrete services (with their respective SLAs) have been identified for each task, except for tasks S_1 and S_3 for which 5 implementations have been identified. The respective SLA parameters, shown in Table 1, differ in terms of cost c_{ij} , availability a_{ij} , and response time r_{ij} (measured in sec). In the experiments, we used this *baseline* set composed of 26 concrete services, as well as an enlarged set of concrete services where we doubled the baseline set (in the following, we refer to the latter as *2x baseline*).

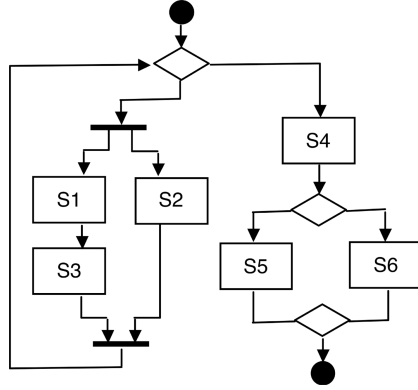


Figure 3. Workflow of the composite service managed by MOSES

Table 1 SLA parameters for concrete services

<i>cs</i>	r_{ij}	a_{ij}	c_{ij}	<i>cs</i>	r_{ij}	a_{ij}	c_{ij}
<i>cs</i> ₁₁	2	0.995	6	<i>cs</i> ₃₁	1	0.995	5
<i>cs</i> ₁₂	1.8	0.99	6	<i>cs</i> ₃₂	1	0.99	4.5
<i>cs</i> ₁₃	2	0.99	5.5	<i>cs</i> ₃₃	2	0.99	4
<i>cs</i> ₁₄	3	0.995	4.5	<i>cs</i> ₃₄	4	0.95	2
<i>cs</i> ₁₅	4	0.99	3	<i>cs</i> ₃₅	5	0.95	1
<i>cs</i> ₂₁	1	0.995	2	<i>cs</i> ₄₁	0.5	0.995	1
<i>cs</i> ₂₂	2	0.995	1.8	<i>cs</i> ₄₂	0.5	0.99	0.8
<i>cs</i> ₂₃	1.8	0.99	1.8	<i>cs</i> ₄₃	1	0.995	0.8
<i>cs</i> ₂₄	3	0.99	1	<i>cs</i> ₄₄	1	0.95	0.6
<i>cs</i> ₅₁	1	0.995	3	<i>cs</i> ₆₁	1.8	0.99	1
<i>cs</i> ₅₂	2	0.99	2	<i>cs</i> ₆₂	2	0.995	0.8
<i>cs</i> ₅₃	3	0.99	1.5	<i>cs</i> ₆₃	3	0.99	0.6
<i>cs</i> ₅₄	4	0.95	1	<i>cs</i> ₆₄	4	0.95	0.4

The concrete services are simple stubs, without internal logic; however, their non-functional behavior conforms to the guaranteed levels expressed in their SLA. The perceived response time is obtained by modeling each concrete service as an $M/D/m/PS$ queue implemented inside the Web service deployed in a Tomcat container. The $M/D/m/PS$ model is parameterized in such a way to have an average CPU usage between 65% and 70% when the request rate is equal to 10 req/sec. Table 2 shows the SLAs offered by MOSES to the composite service users according to different service classes.

Table 2 SLA parameters for service classes

Service class k	R_{\max}^k	A_{\min}^k	C_{\max}^k
1	16	0.88	55
2	18	0.85	50
3	20	0.82	45
4	22	0.79	40

To issue requests to the composite service managed by MOSES we developed a workload generator in C language using the Pthreads library. The workload generator

can be configured to issue requests to MOSES service class by service class at a fixed rate.

We conducted five sets of experiments to evaluate the service selection policies under different scenarios and loads as well as the monitoring and analysis components. The first set of experiments was performed to point out the scalability problems of the traditional per-request approach; the second set was carried out to compare the performance of the traditional per-request policy versus the load-aware one; the third set was performed to analyze the scalability of the load-aware per-request policy. In the first three sets of experiments, each experiment is composed by several runs lasting 15 minutes each, during which the workload generator generates requests corresponding to the service class 2 (see the $k=2$ row in Table 2 for its SLA) at a constant rate. The request rate is then increased run by run until the system keeps stable. The fourth set of experiments was performed in order to prove the effectiveness of our load-aware policy. Specifically, for every service class we generated a constant request rate for the first half of the experiment, then increasing the request rate only for class 2 in the second half on the experiment. Finally, the last set of experiments was performed to demonstrate the effectiveness of the adaptive Cusum algorithm, when some services do not behave according to the stipulated SLAs. As in the first three sets of experiments, each experiment is composed by several runs lasting 15 minutes each. The request rate submitted to MOSES is the same in each experiment, but an external load, increased run by run, is submitted to cs_{11} .

The main performance metric we measured is the response time of the composite service. We also measured the CPU utilization of the concrete services to analyze the different effects of the request load distribution among the concrete services achieved by the traditional and the load-aware per-request policies.

5.3 Experimental results

In the first set of experiments, we ran three load tests on MOSES using the traditional per-request policy in Ref. [3]. In the first test, we used the baseline set of concrete services without instrumenting any of the MOSES modules in the Monitor and Analyze macro-components. In the second test, we used the previous configuration, but we exploited the 2x baseline set of concrete services. Finally, in the last test we used the 2x baseline set of concrete services and we added the support of the QoS Monitor, in order to detect SLA violations of the response time of the concrete services and, in positive case, to determine a new service selection policy that exploits different concrete services implementations.

Figure 4 shows the average response time perceived by the users of the composite service for different request rates submitted to MOSES. We observe that for all the three tests the response time is nearly constant until the request rate reaches 7 req/sec. From this point on, the response time of both the tests without QoS Monitor (*[3], no QoS, baseline* and *[3], no QoS, 2x baseline* curves), regardless of the used concrete services set, rapidly grows because the per-request service selection does not exploit the presence of different service implementations, always using the same service identified as the best one. In the test with the QoS Monitor enabled (*[3], QoS, 2x baseline* curve), the response times of the concrete services are

collected, their average calculated every 2 sec. and analyzed; if the QoS Monitor finds out that the currently used concrete implementations do not have an adequate performance (*i.e.*, they are violating the response time contractualized in the SLA), it triggers the Optimization Engine to compute a new optimal policy \mathbf{x} , using the actual response times of the concrete services instead of those declared into the SLAs. As a result, the currently used overloaded services will not be used in the near future, but they will be candidate for re-usage when the new selected concrete services will in their turn become overloaded. However, the introduction of the QoS Monitor provides only a modest performance improvement; even if the QoS Monitor invocation frequency is relatively high (every 2 sec.), the reaction is not quick enough to address higher request rates. We can conclude that the traditional per-request approach is not able to scale out the available services implementations, and thus it is unable to sustain higher request rates than those sustainable by the bottleneck concrete service.

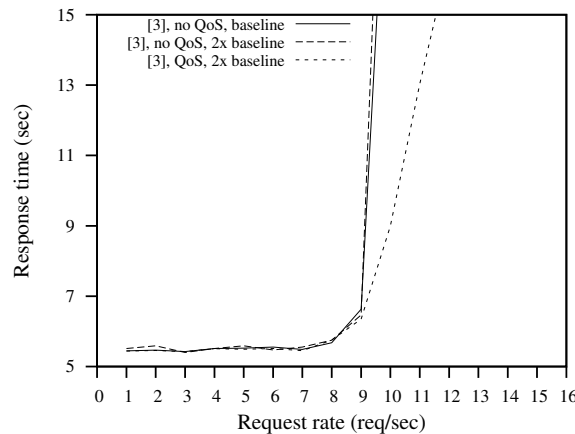


Figure 4. Response time of the traditional per-request service selection policy

The second set of experiments compares the traditional per-request and the load-aware per-request selection policies. These experiments use the baseline set of concrete services and do not involve any Monitor or Analyze MOSES component. Figure 5 compares the average response time according to the request rate submitted to MOSES when using the two different policies. We observe that the response times achieved by the two policies perfectly overlap until the request rate reaches the saturation point of the traditional per-request policy. From this point on, the former (*[3], no QoS, baseline* curve) is not able to exploit the available implementations, while the load-aware policy (*no QoS, baseline* curve) performs better, scaling out the available concrete services. Therefore, the load-aware approach is able to sustain higher request rates than the traditional per-request, given that there are available concrete services to be exploited.

To show the load balancing effectiveness, we monitored the CPU usage of the concrete services during the experiments. Since every concrete service is implemented as an $M/D/m/PS$ queue, the CPU usage has been computed with the formula $\frac{\lambda_{ij}T_{ij}}{nCPU_{ij}}$, where λ_{ij} is the request rate directed to the j -th implementation of

S_i (that is, cs_{ij}), T_{ij} its service time, and $nCPU_{ij}$ the number of CPUs available to that service implementation.

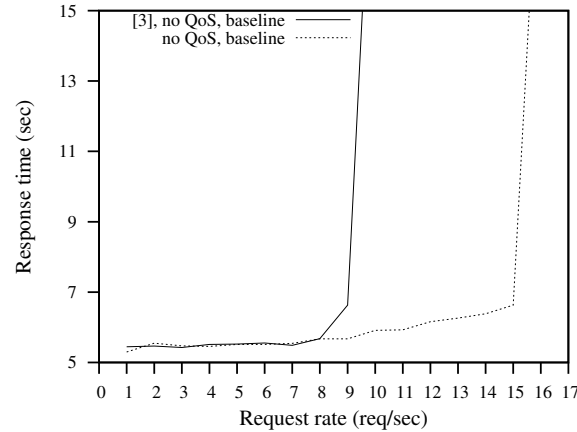


Figure 5. Response time of the traditional versus load-aware per-request service selection policies

Figure 6 shows the CPU usage of cs_{13} , which is the single concrete service used by the traditional per-request optimization approach to implement S_1 . We can see that the load increases almost linearly until it reaches the CPU usage equal to 85%; at that value, the system becomes unstable (see Fig. 5).

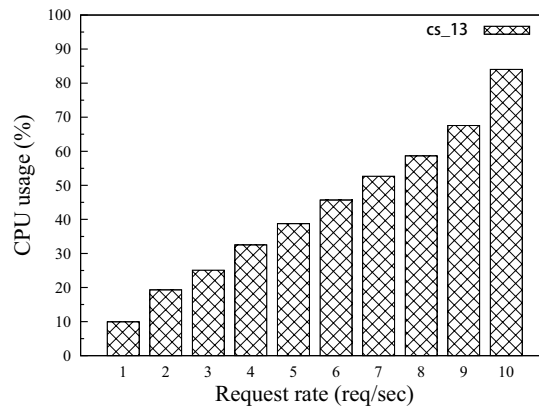


Figure 6. CPU usage of the concrete service selected for S_1 by the traditional per-request policy

Figure 7 shows the CPU usage of the concrete services used by the load-aware policy to implement the abstract task S_1 . Differently from the traditional strategy, with the load-aware policy multiple concrete services can be used to implement the same task. In particular, when the request rate is low (from 1 req/sec to 6 req/sec), there is no need to use multiple concrete services (we recall that each concrete service is modeled so to have an average CPU usage between 65% and 70% when the request rate is equal to 10 req/sec). Therefore, for the low request rate only cs_{13}

is used, like in the traditional per-request policy. When the request rate increases from 7 req/sec to 9 req/sec, the concrete services cs_{13} and cs_{11} are both used. From 10 req/sec on, cs_{15} is also used to implement S_1 , therefore the load is balanced across three concrete services. We observe that the cumulative load does not increase monotonically, because the concrete services model different underlying hardware: cs_{11} and cs_{15} have 29 CPUs each, while cs_{13} has 25 CPUs. Therefore, when more load is directed to a concrete service with a larger capacity, the overall load decreases.

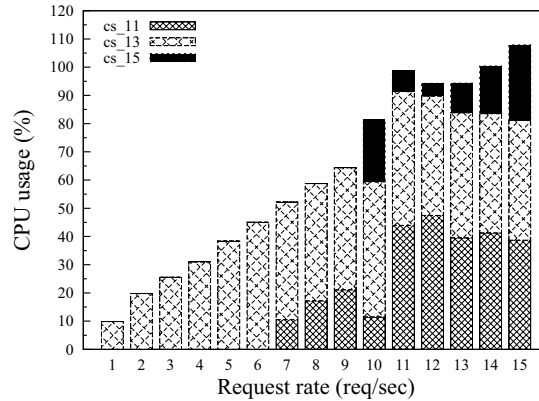


Figure 7. CPU usage of the concrete services selected for S_1 by the load-aware per-request policy

We carried out the third set of experiments to show the scalability capabilities of the load-aware per-request policy. In these experiments we used both the baseline and the 2x baseline sets of concrete services, without deploying the QoS Monitor module. Figure 8 shows the scaling capabilities of the load-aware per-request service selection: until one concrete service at a time can sustain the load (*i.e.*, around 7 req/sec), it does not matter to have a larger number of available implementations; therefore, the results with the baseline set of concrete implementations resemble those with the 2x baseline set. However, at higher request rates, the availability of a larger set of candidate services provides better response times and allows to manage the request rates without incurring in overloading, because the load can be better shared among the available implementations.

We conducted the fourth set of experiments to study the effectiveness of the proposed load-aware per-request approach. We simulated several concurrent users characterized by different service classes. The goal is to prove the effectiveness of the MOSES adaptation under the load-aware per-request policy despite variations in the submitted workload. To this end, each service class submits requests at a constant rate equal to 1 req/sec, except class 2 for which we increased the request rate from 1 to 10 req/sec in the second half of the test. Therefore, in the first half of the experiment, the aggregate workload is equal to 4 req/sec, which can be also easily managed by the traditional per-request policy; on the other hand, in the second half of the experiment we submitted to the SOS an aggregated workload equal to 13 req/sec, which cannot be sustained by the traditional per-request policy (see the concrete services model described in Section 5.2). The overall experiment lasted 1 hour.

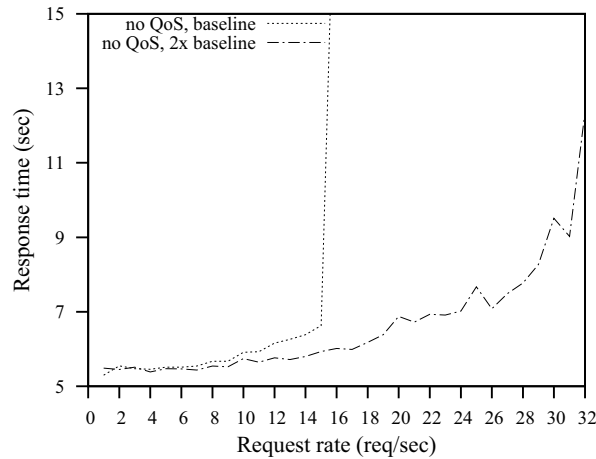


Figure 8. Response time of the load-aware per-request policy under the two sets of concrete services

Figures 9(a)–9(d) show that the perceived response times are far below the response times agreed in the SLAs and represented by the horizontal lines; this can be explained by observing that the average behavior is very different from the worst case considered in the formulation of the optimization policy. This latter issue could be addressed by considering SLAs where the response time constraint is specified in terms of bounds on the percentile.

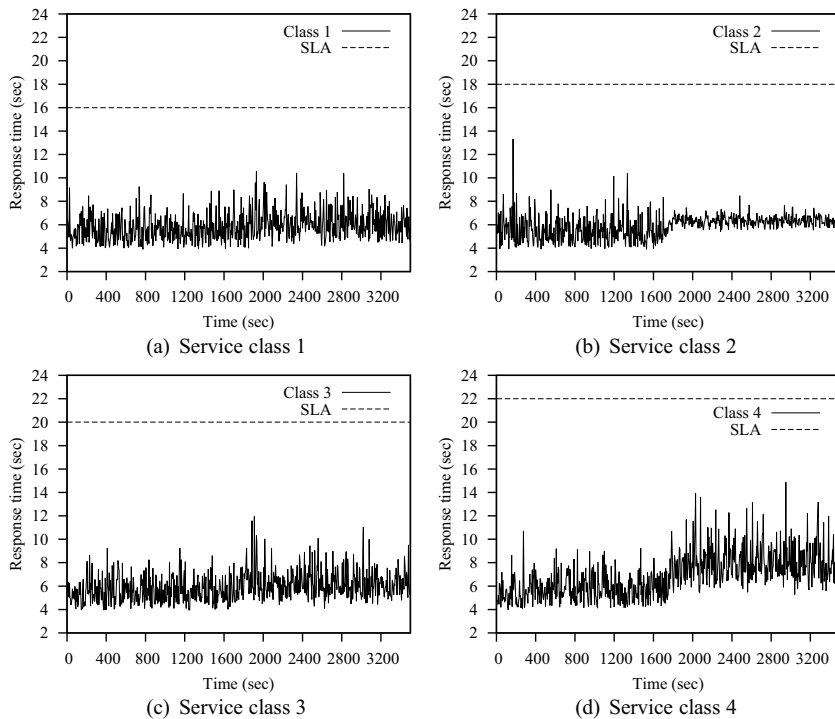


Figure 9. Response time of the load-aware per-request policy for all service classes over time

Table 3 shows the average response times perceived by the users when issuing requests either to a lightly loaded or to an heavy loaded system according to the service class. When the system is subject to a light load, there are not appreciable differences among the service classes. On the other hand, when the load increases, the average response time perceived by class 4 (which has the least stringent SLA) suffers more the load increase. The motivation is that class 4 requests can only exploit a limited number of concrete services, because of the lowest maximum cost in the SLA (see Table 2); therefore, to satisfy the cost constraint they cannot be distributed among all the available concrete services.

Table 3 Average response times of the load-aware per-request policy for all service classes under light and heavy loads

Service class	Light load	Heavy load
1	5.514 sec	6.254 sec
2	5.485 sec	6.350 sec
3	5.509 sec	6.357 sec
4	5.794 sec	8.112 sec

Finally, in the fifth and last set of experiments we demonstrate the effectiveness of the adaptive Cusum algorithm. In the first test, we used the baseline set of concrete services without instrumenting any of the MOSES modules in the Monitor and Analyze macro-components. Conversely, in the second test we added the support of the QoS Monitor, in order to detect SLA violations of the response time of the concrete services and, in positive case, to determine a new service selection policy that exploits different concrete services implementations. Both the test were conducted submitting requests corresponding to the service class 2 at a constant rate equal to 4 req/sec. We also submitted an external load to concrete service cs_{11} in order to overload it. The external load has been incremented by one unit every 15 minutes, starting from 1 req/sec. Figure 10 shows the result of the first test without the QoS Monitor. We can see that the response time of the composite service keeps constant until the external request rate is less than 7 req/sec, . After this threshold, the concrete service cs_{11} begins to be overloaded and therefore the response time of the composite service sharply increases. The SOS performance changes significantly when the QoS Monitor with the adaptive Cusum algorithm is turned on, as shown in Fig. 11. The response time of the composite service is constant, regardless of the external request rate submitted to cs_{11} .

6 Related Work

The service selection problem has been widely investigated in the last years and many solutions can be found in literature. They can be broadly classified into two categories, depending on whether they propose a *local* or a *global* approach. In the local approach, *e.g.*, Refs. [3, 23, 28], only QoS constraints on the execution of single abstract tasks can be predicated: the concrete services are selected one at the time by associating the running abstract task to the best candidate concrete service which supports its execution. On the other hand, the global approach, *e.g.*, Refs. [2-5, 7, 8,

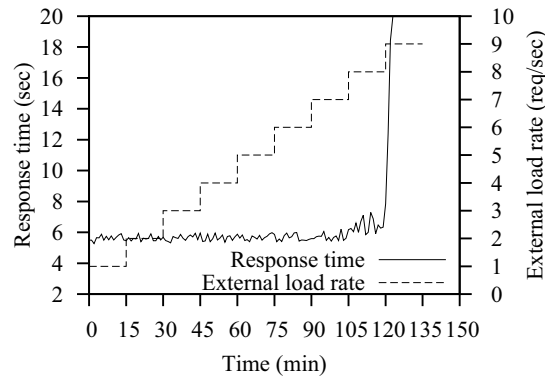


Figure 10. Response time of the load-aware per-request policy under external load without QoS Monitor

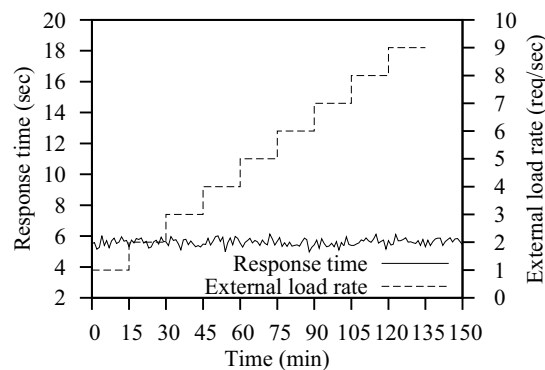


Figure 11. Response time of the load-aware per-request policy under external load with QoS Monitor

10, 16 17, 19, 27, 28], aims to ensure the QoS constraints on the whole execution of the composite service.

Most of the proposed methodologies for service selection focus on the global approach and adopt in particular the *per-request* granularity level, formalizing the service selection as an instance of suitable optimization problems^[3-5, 10, 17, 19, 27, 28]. At the per-request granularity level, the service selection concerns each single request submitted to the service oriented system and has the goal to fulfill the QoS constraints of that specific request, independently of the concurrent requests that may be addressed to the system. Zeng et al.^[28] present a global planning approach to select an optimal execution plan by means of integer programming. Their QoS model consider price, availability, reliability, and reputation as parameters. Ardagna and Pernici^[3] model the service composition as an MILP problem and their technique is particularly efficient for long running process instances. Their approach is formulated as an optimization problem handling the whole application instead of each execution path separately. The proposal by Alrifai and Risse^[4] is slightly different, as the global approach is combined with local selection techniques to reduce the optimization complexity. The global constraints are reduced to local constraints using integer programming in such a way that satisfying the latter also

ensure the former. Another approach to address the complexity of the optimization problem is proposed by Alrifai *et al.* in Ref. [5], where the set of the available concrete services is pruned on the basis of the skyline notion before resolving the optimization problem itself. Canfora *et al.*^[10] follow a quite different strategy for optimal selection, relying on genetic algorithms. They define an iterative procedure to search for the best solution of a given problem among a constant size population without the need for linearization required by integer programming. Since the per-request service selection problem is NP-hard, other heuristic policies have been also proposed (e.g., see Refs. [17, 19, 27]).

The common factor to all the solutions discussed so far is that each abstract task is only bound, from time to time, to a single concrete service. It seems reasonable to suppose that, for a given class of requests, the same optimal binding between an abstract task and a corresponding concrete service holds until a significant change detected in the execution environment triggers the calculation of a new binding. Hence, the per-request policies have the drawback of possibly overloading the selected services during the time interval that interlapses between two subsequent changes, because each request is handled independently of all the others.

This drawback is partially solved by those proposals that adopt the *per-flow* granularity level^[2, 7, 8, 16], where the focus is on the overall flow of requests of a user, rather than on a single request. Under the per-flow granularity, the service selection goal is to fulfill the QoS constraints that concern the global properties of that flow, *e.g.*, the average response time of the composite service. In Ref. [16] the service selection problem is addressed with an LP formulation: an abstract task is probabilistically bound at runtime to several concrete services thus realizing a request load balancing. However, the actual load submitted to each concrete service is not taken into account. Also in Refs. [7-9] the service selection problem is formulated as an LP problem that probabilistically binds each abstract task to multiple corresponding concrete services, but, differently from the work in Ref. [16], in these proposals the load submitted to each concrete service is accounted. The system incoming workload is also taken into account by Ardagna and Mirandola in Ref. [2], but the service selection is formulated as a constrained non-linear optimization problem.

The solutions to the service selection problem presented in Refs. [2, 7-9] take into account the load balancing issue and can scale better than the per-request approaches because of the corresponding formulation of the optimization problem. However, they work on a per-flow basis; therefore, the QoS constraints are ensured on average and in the long term, but no QoS guarantee is given to each single submitted request. In this paper we have proposed a new approach to service selection, that combines the different advantages of the per-request and per-flow approaches proposed so far: it scales similarly to the per-flow ones with respect to the submitted request load, but allows to ensure the QoS constraints on a per-request basis.

With regard to full frameworks for the self-adaptation of SOS, we have few works that extensively describe also the Monitor and Analyze phases of the MAPE-k loop^[1, 11, 18].

With respect to the Monitor phase, all the frameworks focus on monitoring the QoS attributes stated in the SLAs stipulated with the concrete service providers, on

a continuous basis, and using a passive methodology. The only exception is Ref. [11], because also the workload submitted to each service in the service composition and the in-house hardware resources are monitored. Anyway, using a passive approach on a continuous time basis appears to be the facto standard for Web service monitoring.

For what regards the Analyze phase, each framework has its own approach. The work in Ref. [11] uses the data collected in the Monitor phase to parametrize a Markov model of the composite service, and then solves this model to discover if the QoS constraints will be violated. The proposal in Ref. [1] uses a rule-based approach. Before and after each concrete service invocation a rule is evaluated in order to discover if the concrete service behavior differs from what expected. For example, the rule can be constructed using boolean, relational, and mathematical operators or other useful constructs, like the maximum, minimum, and average values of collected data. Finally, the framework presented in Ref. [18] performs the analysis using statistical methods. It does not suggest a single method, but rather it provides as example of suitable methodologies the application of either the Bayesian inference or the Student t-test statistical significance test.

7 Conclusions

In this paper we have presented a theoretical framework for a self-adaptable service oriented system. The core of the framework is a new per-request load-aware policy that addresses the service selection issue for a service broker which offers a composite service with QoS constraints. The proposed policy realizes a randomized load balancing of the requests submitted to each abstract task, exploiting the multiple concrete implementations available in the open service market-place. To avoid overloading the chosen concrete services, the load balancing is tuned by taking into account the capacity of each concrete service and the load submitted to the abstract task. In particular, a probability is assigned to each concrete service proportionally to its capacity and, for each single request, the concrete services to be invoked are selected according to these probabilities. The other component of the proposed framework is the adaptive Cusum algorithm to detect changes in the execution environment. By exploiting it, the SOS can trigger a reconfiguration, thus changing the subset of used concrete services.

Using a prototype implementation, we have first compared our approach with one of the top performing per-request service selection policies, presented in Ref. [3]. Our experimental results show the scalability of the load-aware per-request policy: it can sustain higher request rates than the per-request policy in Ref. [3], because it allows to concurrently exploit multiple concrete services using the load balancing mechanism, while the approach in Ref. [3] uses only one concrete service at a time, *i.e.*, it directs all the concurrent requests in the same QoS class to the same concrete service. The experimental results also show that for a service broker using our policy the maximum sustainable load grows with the number of available concrete services, while this does not happen with the policy in Ref. [3]. Indeed, the latter is completely insensitive to the number of available implementations for each abstract task, while the load-aware service selection has proved suitable to work in a real scenario.

The second experimental contribution has demonstrated the effectiveness of the whole theoretical framework, that includes the adaptive Cusum algorithm for the

Analyze phase. Our experimental results show that with adaptive Cusum the SOS is able to reconfigure itself quickly. Furthermore, the changes in the execution environment do not affect the response time of the composite service and so its users do not perceive any performance degradation.

As future work we will consider QoS constraints which specify bounds on the percentile of the response time, because the user perceived QoS can be better expressed in terms of percentiles rather than mean values. Furthermore, due to the computational complexity of the MILP formulation we used for the per-request load-aware policy, we will consider to apply techniques for pruning the set of available concrete services in such a way to speed-up the problem resolution in case of large problem instances.

References

- [1] Ardagna D, Baresi L, Comai S, Comuzzi M, Pernici B. A service-based framework for flexible business processes. *IEEE Software*, March 2011, 28(2): 61–67.
- [2] Ardagna D, Mirandola R. Per-flow optimal service selection for web services based processes. *J. Syst. Softw.*, 2010, 83(8): 1512–1523.
- [3] Ardagna D, Pernici B. Adaptive service composition in flexible processes. *IEEE Trans. Softw. Eng.*, 2007, 33(6): 369–384.
- [4] Alrifai M, Risse T. Combining global optimization with local selection for efficient qos-aware service composition. *Proc. WWW '09*. ACM. 2009. 881–890.
- [5] Alrifai M, Skoutas D, Risse T. Selecting skyline services for QoS-based web service composition. *Proc. WWW '10*. ACM. 2010. 11–20.
- [6] Bellucci A, Cardellini V, Di Valerio V, Iannucci S. A scalable and highly available brokering service for SLA-based composite services. *Proc. ICSOC '10*. LNCS 6470. Springer. December 2010. 527–541.
- [7] Cardellini V, Casalicchio E, Grassi V, Lo Presti F, Mirandola R. Flow-based service selection for web service composition supporting multiple qos classes. *Proc. IEEE ICWS '07*. 2007. 743–750.
- [8] Cardellini V, Casalicchio E, Grassi V, Lo Presti F, Mirandola R. QoS-driven runtime adaptation of service oriented architectures. *Proc. ACM ESEC/SIGSOFT FSE*. 2009. 131–140.
- [9] Cardellini V, Casalicchio E, Grassi V, Iannucci S, Lo Presti F, Mirandola R. MOSES: a framework for QoS driven runtime adaptation of service-oriented systems. *IEEE Trans. Softw. Eng.*, September 2012, 38(5): 1138–1159.
- [10] Canfora G, Di Penta M, Esposito R, Villani ML. A framework for QoS-aware binding and re-binding of composite web services. *J. Syst. Softw.*, 2008, 81(10): 1754–1769.
- [11] R. Calinescu, L. Grunske, M. Kwiatkowska, R. Mirandola, and G. Tamburrelli. Dynamic QoS management and optimization in service-based systems. *IEEE Trans. Softw. Eng.*, May 2011, 37(3): 387–409.
- [12] IBM ILOG CPLEX Optimizer. <http://www.ibm.com/software/integration/optimization/cplex-optimizer/>.
- [13] Casolari S, Tosi S, Lo Presti F. An adaptive model for online detection of relevant state changes in internet-based systems. *Perform. Eval.*, May 2012, 69(5): 206–226.
- [14] Di Nitto E, Ghezzi C, Metzger A, Papazoglou MP, Pohl K. A journey to highly dynamic, self-adaptive service-based applications. *Autom. Softw. Eng.*, 2008, 15(3–4): 313–341.
- [15] Kephart JO, Chess DM. The vision of autonomic computing. *IEEE Computer*, 2003, 36(1): 41–50.
- [16] Klein A, Ishikawa F, Honiden S. Efficient QoS-aware service composition with a probabilistic service selection policy. *Proc. ICSOC '10*. LNCS 6470. Springer. December 2010. 182–196.
- [17] Liang Q, Wu X, Lau HC. Optimizing service systems based on application-level QoS. *IEEE Trans. Serv. Comput.*, April 2009, 2(2): 108–121.
- [18] Mosincat A, Binder W, Jazayeri M. Achieving runtime adaptability through automated model

- evolution and variant selection. *Enterprise Information Systems*. 2013. to appear.
- [19] Menascé DA, Casalicchio E, Dubey V. On optimal service selection in service oriented architectures. *Perform. Eval.*, August 2010, 67(8): 659–675.
 - [20] Menascé DA, Gomaa H, Malek S, Sousa J. SASSY: a framework for self-architecting service-oriented systems. *IEEE Software*, November 2011, 28(6): 78–85.
 - [21] Motahari-Nezhad HR, Li J, Stephenson B, Graupner S, Singhal S. Solution reuse for service composition and integration. *Proc. WSCA '09*. 2009.
 - [22] Montgomery DC. *Introduction to Statistical Quality Control*. Wiley, 2008.
 - [23] Maamar Z, Sheng QZ, Benatallah B. Interleaving web services composition and execution using software agents and delegation. *Proc. WSABE '03*. 2003.
 - [24] OASIS. *Web Services Business Process Execution Language Version 2.0*. January 2007.
 - [25] Papazoglou MP. Service-oriented computing: Concepts, characteristics and directions. *Proc. IEEE WISE '03*. 2003.
 - [26] Salehie M, Tahvildari L. Self-adaptive software: Landscape and research challenges. *ACM Trans. Auton. Adapt. Syst.*, 2009, 4(2): 1–42.
 - [27] Yu T, Zhang Y, Lin KJ. Efficient algorithms for Web services selection with end-to-end QoS constraints. *ACM Trans. Web*, 2007, 1(1): 1–26.
 - [28] Zeng L, Benatallah B, Dumas M, Kalagnamam J, Chang H. QoS-aware middleware for web services composition. *IEEE Trans. Softw. Eng.*, 2004, 30(5).