

# Hybrid Optimization Algorithm for Large-Scale QoS-Aware Service Composition

Pablo Rodríguez-Mier, Manuel Mucientes, and Manuel Lama

**Abstract**—In this paper we present a hybrid approach for automatic composition of web services that generates semantic input-output based compositions with optimal end-to-end QoS, minimizing the number of services of the resulting composition. The proposed approach has four main steps: 1) generation of the composition graph for a request; 2) computation of the optimal composition that minimizes a single objective QoS function; 3) multi-step optimizations to reduce the search space by identifying equivalent and dominated services; and 4) hybrid local-global search to extract the optimal QoS with the minimum number of services. An extensive validation with the datasets of the Web Service Challenge 2009-2010 and randomly generated datasets shows that: 1) the combination of local and global optimization is a general and powerful technique to extract optimal compositions in diverse scenarios; and 2) the hybrid strategy performs better than the state-of-the-art, obtaining solutions with less services and optimal QoS.

**Index Terms**—Service composition, service optimization, hybrid algorithm, QoS-aware, Semantic Web services

## 1 INTRODUCTION

WEB services are self-describing software applications that can be published, discovered and invoked across the web using standard technologies [1]. The functionality of a Web service is mainly determined by the functional properties that describe their behaviour in terms of its inputs, outputs, and also possibly additional descriptions that the services may have, such as preconditions and effects. These four characteristics, commonly abbreviated IOPEs, allow the composition and aggregation of web services into composite web services that achieve more complex functionalities and, therefore, solve complex user needs that cannot be satisfied with atomic web services. However, compositions should go beyond achieving a concrete functionality and take into account other requirements such as quality-of-service (QoS) to generate also compositions that fit the needs of different contexts. The QoS determines the value of different quality properties of services such as response time (total time a service takes to respond to a request) or throughput (number of invocations supported in a given time interval), among others characteristics. These properties apply both to single services and to composite services, where each individual service in the composition contributes to the global QoS. For composite services this implies that having many different services with similar or identical functionality, but different QoS, may lead to a large amount of possible compositions that satisfy the same functionality with different QoS but also with a different number of services.

However, the problem of generating automatic compositions that satisfy a given request with an optimal QoS is a

very complex task, specially in large-scale environments, where many service providers offer services with similar functionality but with different QoS. This has motivated researchers to explore efficient strategies to generate QoS-aware web service compositions from different perspectives [2], [3]. But despite the large number of strategies proposed so far, the problem of finding automatic compositions that minimize the number of services while guaranteeing the optimal end-to-end QoS is rarely considered. Instead, most of the work has focused on optimizing the global QoS of a composition or improving the execution time of the composition engines. An analysis of the literature shows that only a few works take into consideration the number of services of the resulting optimal QoS compositions. Some notable examples are [4], [5], [6], [7]. Although most of these composition engines are quite efficient in terms of computation time, none of them are able to effectively minimize the total number of services of the solution while keeping the optimal QoS.

The ability to provide not only optimal QoS but also an optimal number of services is specially important in large-scale scenarios, where the large number of services and the possible interactions among them may lead to a vast amount of possible solutions with different number of services but also with the same optimal QoS for a given problem. Moreover, there can be situations where certain QoS values are missing or cannot be measured. Although the prediction of QoS can partially alleviate this problem [8], it is not always possible to have historical data in order to build statistical models to accurately predict missing QoS. In this context, optimizing not only the available QoS but also the number of services of the composition may indirectly improve other missing properties. This has important benefits for brokers, customers and service providers. From the broker point of view, the generation of smaller compositions is interesting to achieve manageable compositions that are easier to execute, monitor, debug, deploy and scale. On the other hand, customers can also benefit from smaller

- The authors are with the Centro de Investigación en Tecnologías de la Información (CiTIUS), Universidade de Santiago de Compostela, Spain. E-mail: {pablo.rodriquez.mier, manuel.mucientes, manuel.lama}@usc.es.

Manuscript received 9 May 2015; revised 8 Sept. 2015; accepted 16 Sept. 2015. Date of publication 22 Sept. 2015; date of current version 9 Aug. 2017. For information on obtaining reprints of this article, please send e-mail to: reprints@ieee.org, and reference the Digital Object Identifier below. Digital Object Identifier no. 10.1109/TSC.2015.2480396

compositions, specially when there are multiple solutions with the same optimal end-to-end QoS but different number of services. This is even more important when service providers do not offer fine-grained QoS metrics, since decreasing the number of services involved in the composition may indirectly improve other quality parameters such as communication overhead, risk of failure, connection latency, etc. This is also interesting from the perspective of service providers. For example, if the customer wants the cheapest composition, the solution with fewer services from the same provider may also require less resources for the same task.

However, one of the main difficulties when looking for optimal solutions is that it usually requires to explore the complete search space among all possible combinations of services, which is a hard combinatorial problem. In fact, finding the optimal composition with the minimum number of services is NP-hard (see Appendix A, which can be found on the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TSC.2015.2480396>). Thus, achieving a reasonable trade-off between solution quality and execution time in large-scale environments is far from trivial, and hardly achievable without adequate optimizations.

In this paper we focus on the automatic generation of semantic input-output compositions, minimizing both a single QoS criterion and the total number of services subject to the optimal QoS. The main contributions are:

- A multi-step optimization pipeline based on the analysis of non-relevant, equivalent and dominated services in terms of interface functionality and QoS.
- A fast local search strategy that guarantees to obtain a near-optimal number of services while satisfying the optimal end-to-end QoS for an input-output based composition request.
- An optimal combinatorial search that can improve the solution obtained with the local search strategy by performing an exhaustive combinatorial search to select the composition with the minimum number of services for the optimal QoS.

We tested our proposal using the Web Service Challenge (WSC) 2009-2010 datasets and, also, a different randomly generated dataset with a variable number of services. The rest of the paper is organized as follows: Section 4 introduces the composition problem, Section 5 describes the proposed approach, Section 6 presents the results obtained, and Section 7 gives some final remarks.

## 2 RELATED WORK

Automatic composition of services is a fundamental and complex problem in the field of service oriented computing, which has been approached from many different perspectives depending on what kinds of assumptions are made [2], [3], [9], [10]. AI Planning techniques have been traditionally used in service composition to generate valid composition plans by mapping services to actions in the planning domain [11], [12], [13], [14], [15], [16]. These techniques work under the assumption that services are complex operators that are well defined in terms of IOPEs, so the problem can be translated to a planning problem and solved using classical

planning algorithms. Most of these approaches have been mainly focused on exploiting semantic techniques [13], [16], [17] and developing heuristics [15], [16], [18] to improve the performance of the planners. As a result, and partly given by the complexity of generating satisfiable plans in the planning domain, these approaches do not generate neither optimal plans (minimizing the number of actions) nor optimal QoS-aware compositions.

Other approaches have studied the QoS-aware composition problem from the perspective of Operation Research, providing interesting strategies for optimal selection of services and optimizing the global QoS of the composition subject to multiple QoS constraints. A common strategy is to reduce the composition problem to a combinatorial Knapsack-based problem, which is generally solved using constraint satisfaction algorithms (such as Integer Programming) [19], [20], [21], [22], [23] or Evolutionary Algorithms [24], [25]. Some relevant approaches are [19], [22]. In [19] the authors present AgFlow, a QoS middleware for service composition. They analyze two different methods for QoS optimization, a local selection and a global selection strategy. The second strategy is able to optimize the global end-to-end QoS of the composition using a Integer Linear Programming method, which performs better than the suboptimal local selection strategy. Similarly, in [22] the authors propose a hybrid QoS selection approach that combines a global optimization strategy with local selection for large-scale QoS composition. The assumption made by all these approaches is that there is only one composition workflow with a fixed set of abstract tasks, where each abstract task can be implemented by a concrete service. Both the composition workflow and the service candidates for each abstract task are assumed to be predefined beforehand, so these techniques are not able to produce compositions with variable size.

A different category of techniques are graph-based approaches that 1) generate the entire composition by selecting and combining relevant services and 2) optimize the global QoS of the composition. These techniques usually combine variants or new ideas inspired by different fields, such as AI Planning, Operations Research or Heuristic Search, in order to resolve more efficiently the automatic QoS composition, usually for a single QoS criterion. Some relevant approaches in this category are the top-3 winners of the Web Service Challenge 2009-2010 [4], [5], [6]. Concretely, the winners of the WSC challenge [4], presented an approach that automatically discovers and composes services, optimizing the global QoS. This approach also includes an optimization phase to reduce the number of services of the solution. Although the proposed algorithm has in general good performance, as demonstrated in the WSC, it cannot guarantee to obtain optimal solutions in terms of number of services. The other participants of the WSC have also the same limitation.

A recent and interesting approach in this category has been recently presented by Jiang et al. [26]. In this paper, the authors analyze the problem of generating top  $K$  query compositions by relaxing the optimality of the QoS in order to introduce service variability. However, the compositions are generated at the expense of worsening the optimal QoS, instead of looking first for all possible composition alternatives with the minimum number of services that guarantee the optimal QoS.

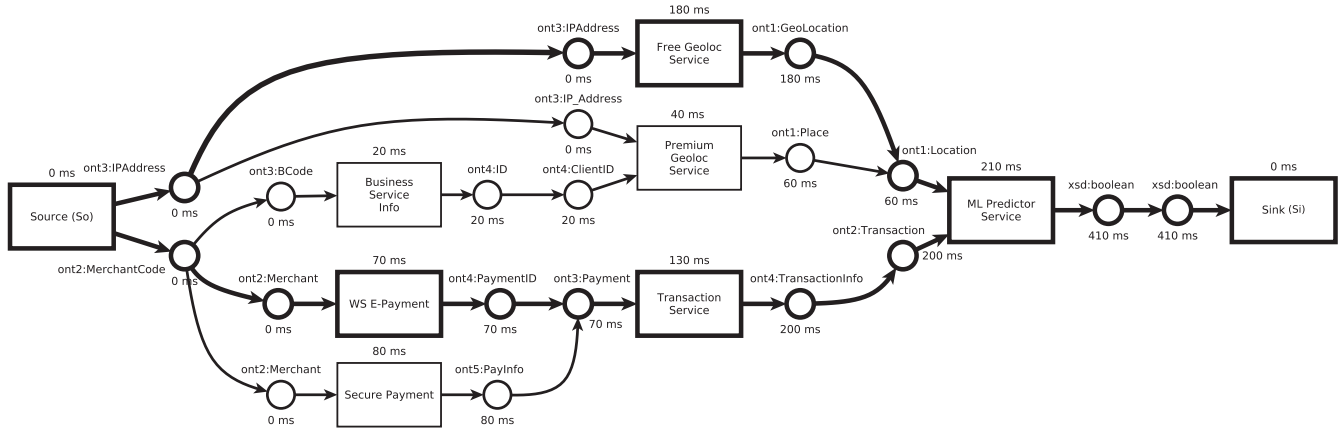


Fig. 1. Example of a *Service Match Graph* for a request  $R = \{\{ont3:IPAddress, ont2:MerchantCode\}, \{xsd:boolean\}\}$  to predict whether a business transaction is fraudulent or not. Each service is associated with an average response time. The optimal solution (*Service Composition Graph*), with an overall response time of 410 ms and four services (excluding  $So$  and  $Si$ ) is highlighted.

Another interesting graph-based approach has been presented in [7]. In this paper, the authors propose a service removal strategy that detects services that are redundant in terms of functionality and QoS. Results show that service removal techniques can be very effective to reduce the number of services before extracting the final composition, as anticipated by other similar approaches [27], [28], [29]. However, some important limitations of this work are: 1) The QoS is not always optimal, since the graph generated for the composition is not complete as it does not contain all the relations between services (it is acyclic) and 2) although the redundancy removal is an effective technique that can be used also to prune the search space, this strategy itself cannot provide optimal results in terms of number of services, and it should be combined with exhaustive search to improve the solutions obtained.

In summary, despite the large number of approaches for automatic QoS-aware service composition there is a lack of efficient techniques that are not only able to optimize the global end-to-end QoS, but also effectively minimize the number of services of the composition. This paper aims to provide an efficient graph-based approach that uses a hybrid local-global optimization algorithm in order to find optimal compositions both in terms of single QoS criteria and in terms of minimum number of services.

### 3 MOTIVATION

The aim of the automatic service composition problem, as considered in this paper, is to automatically select the best combination of available QoS-aware services in a way that can fulfil a user request that otherwise could not be solved by just invoking a single, existing service. This request is specified in terms of the information that the user provides (inputs), and the information it expects to obtain (outputs). The resulting composition should meet this request with an optimal, single criterion end-to-end QoS and using as less services as possible.

A motivating example of the problem is shown in Fig. 1. The figure represents a graph with all the relevant services for a request  $R$  where the inputs are  $\{ont3:IPAddress, ont2:MerchantCode\}$  and the output is  $\{xsd:boolean\}$ . The goal of this example is to obtain a composition to predict

whether a business transaction is fraudulent or not. Each service (associated to a response time QoS) is represented by squares. Inputs and outputs are represented by circles. The graph also contains edges connecting outputs and inputs. These edges represent valid semantic *matches* whenever an output of a service can be passed as an input of a different service. As can be seen, there are some inputs ( $ont1:Location, ont3:Payment$ ) that can be matched by more than one output, so there are many different ways to combine services to achieve the same goal.

Although finding the proper combination of services in terms of their inputs/outputs is essential to generate a solution, it is not enough to obtain good compositions, since there can exist different combinations of services with different QoS. Moreover, many different combinations of services may produce compositions with a different number of services but the same end-to-end QoS. For example, in Fig. 1 we can select *WS E-Payment* service or the *Secure Payment* service to process the electronic payment. However, the second service has a higher response time. Using this leads to a sub-optimal end-to-end QoS of 420 ms. However, there are other situations where the selection of different services leads to compositions with different size but same end-to-end QoS. For example, both *Free Geoloc Service* or the *Premium Geoloc Service* can be selected to translate an *IP* to a *Location*. Although the second one has a better average response time (40 ms), it requires an additional service to obtain the *ClientID* for verification purposes. However, selecting the *Premium Geoloc Service* or the *Free Geoloc Service* does not have an impact on the global QoS, since the *ML Predictor Service* has to wait longer to obtain the *Transaction* parameter (200 ms), but it has an impact on the total number of services of the solution.

The goal of this paper is to automatically generate, given a composition request, a graph like the one represented in Fig. 1 as well as to extract the optimal end-to-end QoS composition with the minimum number of services from that graph.

### 4 PROBLEM FORMULATION

We herein formalize the main concepts and assumptions regarding the composition model used in our approach, which consists of a semantic, graph-centric representation



of the service composition. These concepts are captured in three main models: 1) a service model, which is used to represent services and define how services can be connected or matched to generate composite services; 2) a graph-based composition model, which is used to represent both service interactions and compositions; and 3) a QoS computation model, which provides the operators required to compute the global QoS in a graph-based composition.

#### 4.1 Semantic Service Model

The automatic composition of services requires a mechanism to select appropriated services based on their functional descriptions, as well as to automatic match the services together by linking their inputs and outputs to generate executable data-flow compositions. To this end, we introduce here the main concepts that we use in this paper to support the automatic generation of compositions. This model is an extension of a previous model used in [30] to include QoS properties.

**Definition 1.** A Composition Request  $R$  is defined as a tuple  $R = \{I_R, O_R\}$ , where  $I_R$  is the set of provided inputs, and  $O_R$  the set of expected outputs. Each input and output is related to a semantic concept from the set  $C$  of the concepts defined in an ontology  $Ont$  ( $In_w, Out_w \subseteq C$ ). We say that a composition satisfies the request  $R$  if it can be invoked with the inputs in  $I_R$  and returns the outputs in  $O_R$ .

**Definition 2.** A Semantic Web Service (hereafter “service”) can be defined as a tuple  $w = \{In_w, Out_w, Q_w\} \in W$  where  $In_w$  is a set of inputs required to invoke  $w$ ,  $Out_w$  is the set of outputs returned by  $w$  after its execution,  $Q_w = \{q_w^1, \dots, q_w^n\}$  is the set of QoS values associated to the service, and  $W$  is the set of all services available in the service registry.

Each input and output is related to a semantic concept from the set  $C$  of the concepts defined in an ontology  $Ont$  ( $In_w, Out_w \subseteq C$ ). Each QoS value  $q_w^i \in Q_w$  has a concrete type associated to a set of valid values  $Q$ . For example, the QoS values of a service  $w$  with two different measures, an average response time of 20 ms and an average throughput of 1,000 invocations/second, is represented as  $Q_w = \{20 \text{ ms}, 1,000 \text{ inv/s}\}$ , where  $20 \text{ ms} \in Q_{RT}$  and  $1,000 \text{ inv/s} \in Q_{TH}$ .

Semantic inputs and outputs are used to compose the functionality of multiple services by matching their inputs and outputs together. In order to measure the quality of the match, we need a matchmaking mechanism that exploits the semantic I/O information of the services. The different matchmaking degrees that are contemplated are *exact*, *plugin*, *subsumes* and *fail* [31].

**Definition 3.** Given  $a, b \in C$ ,  $\text{degree}(a, b)$  returns the degree of match between both concepts (*exact*, *plugin*, *subsume* or *fail*), which is determined by the logical relationship of both concepts within the Ontology.

**Definition 4.** Given  $a, b \in C$ ,  $\text{match}(a, b)$  holds if  $\text{degree}(a, b) \neq \text{fail}$ .

In order to determine which concepts are matched by other concepts, we define a matchmaking operator “ $\otimes$ ” that given two sets of concepts  $C_1, C_2 \subseteq C$ , it returns the concepts from  $C_2$  matched by  $C_1$ .

**Definition 5.** Given  $C_1, C_2 \subseteq C$ , we define “ $\otimes : C \times C \rightarrow C$ ” such that  $C_1 \otimes C_2 = \{c_2 \in C_2 | \text{match}(c_1, c_2), c_1 \in C_1\}$ .

We can use the previous operator to define the concepts of full and partial matching between concepts.

**Definition 6.** Given  $C_1, C_2 \subseteq C$ , a full matching between  $C_1$  and  $C_2$  exists if  $C_1 \otimes C_2 = C_2$ , whereas a partial matching exists if  $C_1 \otimes C_2 \subset C_2$ .

**Definition 7.** Given a set of concepts  $C' \subseteq C$ , a service  $w = \{In_w, Out_w\}$  is invocable if  $C' \otimes In_w = In_w$ , i.e., there is a full match between the provided set of concepts  $C'$  and  $In_w$ , so the information required by  $w$  is fully satisfied.

This internal model used by the algorithm, which captures the core components required to perform semantic matchmaking and composition of services, is agnostic to how semantic services are represented. Thus, the algorithm is not bound to any concrete service description. Concretely, different service descriptions can be handled by the algorithm through the use of iServe importers for OWL-S, WSMO-lite, SAWSDL or MicroWSMO. For further details see [32].

#### 4.2 Graph-Based Composition Model

In a nutshell, a data-flow composition of services can be seen as a set of services connected together through their inputs and output, using the semantic model defined before, in a way that every service in the composition is invocable and the invocation of each service in the composition can transform a set of inputs into a set of outputs. These concepts can be naturally captured by graphs, where the vertices represent inputs, outputs and services, and the edges represent semantic matches between inputs and outputs. Here we define the notion of *Service Match Graph* and *Service Composition Graph*. The *Service Match Graph* is a graph that captures all the existent dependencies (matches) between all the relevant services for a composition request. The *Service Composition Graph* is a particular case of the *Service Match Graph* that represents a composition contained in the *Service Match Graph*.

The *Service Match Graph* represents the space of all possible valid solutions for a composition request  $R$ , and it is defined as a directed graph  $G_S = (V, E)$ , where:

- $V = W_R \cup I \cup O \cup \{So, Si\}$  is the set of vertices of the graph, where  $W_R \subseteq W$  is the set of relevant services,  $I$  is the set of inputs and  $O$  is the set of outputs.  $Si$  and  $So$  are two special services, called *Source* and *Sink* defined as  $So = \{\emptyset, I_R\}$ ,  $Si = \{O_R, \emptyset\}$ .
- $E = IW \cup WO \cup OI$  is the set of edges in the graph where:
  - $IW \subseteq \{(i_w, w) | i_w \in I \wedge w \in W\}$  is the set of input edges, i.e., edges connecting input concepts to their services.
  - $WO \subseteq \{(w, o_w) | w \in W \wedge o_w \in O\}$  is the set of output edges, i.e., edges connecting services with their output concepts.
  - $OI \subseteq \{(o_w, i_{w'}) | o_w, i_{w'} \in (I \cup O) \wedge \text{match}(o_w, i_{w'})\}$  is the set of edges that represent a semantic match between an output of  $w$  and an input of  $w'$ .

There are also some restrictions in the edge set to ensure that each input/output belongs to a single service:

- $\forall i \in I, d_{G_S}^+(i) = 1 \wedge ch_{G_S}(i) = \{w\}, w \in W$  (each input has only one outgoing edge which connects the input with its service)
- $\forall o \in O, d_{G_S}^-(o) = 1 \wedge par_{G_S}(o) = \{w\}, w \in W$  (each output has only one incoming edge which connects the output with its service)

Function  $d_{G_S}^+(v)$  returns the outdegree of a vertex  $v \in G_S$  (number of children vertices connected to  $v$ ), whereas  $d_{G_S}^-(v)$  returns the indegree of a vertex  $v$  (number of parent vertices connected to  $v$ ). The functions  $ch_G(v)$  and  $par_G(v)$  are the functions that returns the children vertices of  $v$  and the parent vertices of  $v \in G_S$ , respectively.

Fig. 1 shows an example of a *Service Match Graph* where each service is associated with its average response time. As can be seen, this graph contains many different compositions since there are inputs in the graph that can be matched by the outputs of different services. For example, the parent nodes of the input *ont1:Location* of the service *ML Service Predictor* ( $par_G(ont1:Location)$ ) in Fig. 1 are *ont1:GeoLocation* and *ont1:Place*, so the input is matched by two outputs  $d_{G_S}^-(ont1:Location) = 2$ .

A *Service Composition Graph*, denoted as  $G_C = (V, E)$ , represents a solution for the composition request where each input is exactly matched by one output. Formally, it is a subgraph of *Service Match Graph* ( $G_C \subseteq G_S$ ) that satisfies the following conditions:

- $\forall i \in I, d_{G_C}^-(i) = 1$  (each input is strictly matched by one output)
- $G_C$  is a Directed Acyclic Graph (DAG)

These conditions are important in order to guarantee that a solution is valid, i.e., each input is matched by an output of a service and each service is invocable (all inputs on the composition are matched with no cyclic dependencies). This definition of service composition is language-agnostic, so the resulting DAG is a representation of a solution for the composition problem which can be translated to a concrete language, such as OWL-S or BPEL.

### 4.3 QoS Computation Model

Before looking for optimal QoS service compositions, we need first to define a model to work with QoS over compositions of services which allow us to determine the best QoS that can be achieved for a given composition request on a service repository. When many services are chained together in a composition, the QoS of each individual service contributes to the global QoS of the composition. For example, suppose we want to measure the total response time of a simple composition with two services chained in sequence. The total response time is calculated as the sum of the response time of each service in the composition. However, if the composition has two services in parallel, the total time of the composition is given by the slowest services. Thus, the calculation of the QoS of a composition depends on the type of the QoS and on the structure of the composition.

In order to define the common rules to operate with QoS values in composite services, many approaches use a QoS computation model based on workflow patterns [33], which is adequate to measure the QoS of control-flow based compositions. However, this paper focuses on the automatic generation of optimal QoS-aware compositions driven by the

TABLE 1  
QoS Algebra Elements for Response Time and Throughput

QoS ( $Q$ )	$a \oplus b$	$a \ominus b$	$e$	$\phi$	Order ( $\preceq$ )
$Q_{RT} = \mathbb{R}_{\geq 0} \cup \{\infty\}$	$a + b$	$a - b$	0	$\infty$	$\leq$
$Q_{TH} = \mathbb{R}_{\geq 0} \cup \{\infty\}$	$\min(a, b)$	$\min(a, b)$	$\infty$	0	$\geq$

data-flow analysis of the service dependencies (input-output matches) that are represented as a *Service Match Graph*.

In this section we explain the general graph-centric QoS computation model that we use, based on the *path algebra* defined in [34]. This model is better suited to compute QoS values in a *Service Match Graph*, which, for extension, is also applicable to the particular case of the *Service Composition Graph*.

**Definition 8.**  $(Q, \oplus, \ominus, \preceq)$  is a QoS algebraic structure to operate with a set of QoS values, denoted as  $Q$ . This set is equipped with the following elements:

- $\oplus : Q \times Q \rightarrow Q$  is a closed binary operation for aggregating QoS values
- $\ominus : Q \times Q \rightarrow Q$  is a binary operation for subtracting QoS values
- $\preceq$  is a total order relation on  $Q$

This algebraic structure has the following properties:

- 1)  $Q$  is closed under  $\oplus$  (any aggregation of two QoS values always returns a QoS value)
- 2) The set  $Q$  contains an identity element  $e$  such that  $\forall a \in Q, a \oplus e = e \oplus a = a$
- 3) The set  $Q$  contains a zero element  $\phi$  such that  $\forall a \in Q, \phi \oplus a = a \oplus \phi = \phi$
- 4) The operator  $\oplus$  is associative
- 5) The operator  $\oplus$  is monotone for  $\preceq$  (preserves order). This implies that  $\forall a, b, c \in Q, a \preceq b \Leftrightarrow a \oplus c \preceq b \oplus c$
- 6) The operator  $\ominus$  is the inverse of  $\oplus$ :  $a \ominus b = c \Leftrightarrow a = c \oplus b$

Table 1 shows an example of the concrete elements in this algebra. Note that, for the sake of brevity, only the response time and throughput operators are represented in Table 1. However, other QoS properties such as cost, availability, reputation, etc, can also be defined by instantiating the corresponding operators. We denote  $Q_{RT}$  the set of QoS values for response time (in milliseconds),  $Q_{TH}$  the set of QoS values for throughput (invocations/second). The total order comparator  $\preceq$  is required to be able to order and compare different QoS values. Given two QoS values  $a, b \in Q$ ,  $a \preceq b$  means that  $a$  is equal or *better* than  $b$ , whereas  $b \preceq a$  means that  $a$  is equal or *worse* than  $b$ . The order depends on the concrete comparator defined on  $Q$ . For example,  $Q_{RT}$  uses the comparator  $\leq$  to order the response time, so  $a, b \in Q_{RT}$ ,  $a \preceq b \Leftrightarrow a \leq b$ . For example, given two response times 10, 20 ms  $\in Q_{RT}$ , 10 ms  $\preceq$  20 ms (10 ms is better than 20 ms) since 10 ms < 20 ms. However,  $Q_{TH}$  uses the comparator  $\geq$ , so  $a, b \in Q_{TH}$ ,  $a \preceq b \Leftrightarrow a \geq b$ . For example, given two throughput values 10 inv/s, 20 inv/s  $\in Q_{TH}$ , 20 inv/s  $\preceq$  10 inv/s (20 inv/s is better than 10 inv/s) since 20 inv/s > 10 inv/s. This order relation also affects the behavior of the min and max functions. The min function always selects the

best QoS value, whereas the max function always selects the worst QoS value.

**Definition 9.**  $F_Q(w) : W \rightarrow Q$  is a function that given a service  $w \in W$ , it returns its corresponding QoS value from  $Q_w$  with type  $Q$ . This function can be seen as a function to measure the QoS of a service.

For example, in Fig. 1,  $F_{Q_{RT}}(\text{Trans. Service}) = 130$  ms.

**Definition 10.**  $V_Q(w) : W \rightarrow Q$  is a function that given a service  $w$ , it returns its aggregated QoS value. This is defined as:

$$V_Q(w) = \begin{cases} \max_{i \in In_w} (V_Q^{\text{in}}(i)) \oplus F_Q(w) & \text{if } In_w \neq \emptyset \\ F_Q(w) & \text{if } In_w = \emptyset. \end{cases} \quad (1)$$

Informally, this function calculates the aggregated QoS of a service by taking the worst value of the QoS of its inputs plus the current QoS value of the service itself. Taking for example the service *Premium Geoloc Service* from Fig. 1,  $V_{Q_{RT}}(\text{Premium Geoloc Service})$  is computed as  $\max(V_{Q_{RT}}^{\text{in}}(\text{ont3:IP\_Address}), V_{Q_{RT}}^{\text{in}}(\text{ont4:ClientID})) \oplus 40$  ms, which is  $\max(0 \text{ ms}, 20 \text{ ms}) \oplus 40 \text{ ms} = 60 \text{ ms}$  (see Definition 12).

**Definition 11.**  $V_Q^{\text{out}}(o_w) : O \rightarrow Q$  is a function that given an output of a service  $w$ ,  $o_w \in O$ , it returns its aggregated QoS value. The aggregated QoS of an output is equal to the aggregated QoS of a service. Thus, it is defined as:

$$V_Q^{\text{out}}(o_w) = V_Q(w). \quad (2)$$

For example, the aggregated QoS of the output *ont1:Place* ( $V_{Q_{RT}}^{\text{out}}(\text{ont1:Place})$ ) is equal to the aggregated QoS of its service *Premium Geoloc Service* ( $V_{Q_{RT}}(\text{Premium Geoloc Service})$ ), which is equal to 60 ms.

**Definition 12.**  $V_Q^{\text{in}}(i_w) : I \rightarrow Q$  is a function that given an input of a service  $w$ ,  $i_w \in I$ , it returns its optimal aggregated QoS value. This function is defined as:

$$V_Q^{\text{in}}(i_w) = \begin{cases} \phi & \text{if } d_{G_S}^-(i_w) = 0 \\ V_Q^{\text{out}}(o_{w'}), o_{w'} \in \text{par}_G(i_w) & \text{if } d_{G_S}^-(i_w) = 1 \\ \min_{o_{w'} \in \text{par}_G(i_w)} (V_Q^{\text{out}}(o_{w'})) & \text{if } d_{G_S}^-(i_w) > 1. \end{cases} \quad (3)$$

Given an input  $i_w \in In_w$  of a service  $w$ , this function returns the accumulated QoS for that input. If the evaluated input is not matched by any output ( $d_{G_S}^-(i_w) = 0$ ), then the accumulated QoS of the input is undefined. If the evaluated input is matched by just one output ( $d_{G_S}^-(i_w) = 1$ ), then its accumulated QoS value is equal to the accumulated QoS of that output. If the evaluated input can be matched by more than one output ( $d_{G_S}^-(i_w) > 1$ ), i.e., there are many services that can match that input, then its accumulated QoS value is computed by selecting the optimal (best) QoS.

For example, the optimal aggregated QoS of the input *ont3:Payment* from *Transaction Service* ( $V_{Q_{RT}}^{\text{in}}(\text{ont3:Payment})$ ) is calculated as  $\min(V_{Q_{RT}}^{\text{out}}(\text{ont3:PaymentID}), V_{Q_{RT}}^{\text{out}}(\text{ont5:PayInfo})) = 70$  ms.

**Definition 13.** We define  $V_Q^G(g) : G \rightarrow Q$  as a function that given a Service Match Graph  $g = (V, E)$ , it returns its

optimal aggregated QoS value. This is defined as:

$$V_Q^G(g) = V_Q(Si), Si \in V. \quad (4)$$

Basically, the optimal QoS of a *Service Match Graph*  $G_S$  corresponds with the optimal aggregated QoS of its service  $Si \in G_S$ .

#### 4.4 Composition Problem

Given a composition request  $R = \{I_R, O_R\}$ , a set of semantic services  $W$ , a semantic model and a QoS algebra, the composition problem considered in this paper consists of generating the *Service Match Graph*  $G_S$  and selecting a composition graph  $G_C \subset G_S$  such that:

- 1)  $\forall G'_C, V_Q^G(G_C) \leq V_Q^G(G'_C)$ , i.e., the composition graph has the best possible QoS
- 2)  $W_R \subseteq V, |W_R|$  is minimized (the composition graph contains the minimum number of services).

### 5 COMPOSITION ALGORITHM

On the basis of the formal definition of the automatic QoS-aware composition problem, in this section we present our hybrid approach strategy for automatic, large-scale composition of services with optimal QoS, minimizing the services involved in the composition. The approach works as follows: given a request, a directed graph with the relevant services for the request is generated. Once the graph is built, an optimal label-correcting forward search is performed in polynomial time in order to compute the global optimal QoS. This information is used later in a multi-step pruning phase to remove sub-optimal services. Finally, a hybrid local/global search is performed within a fixed time limit to extract the optimal solution from the graph. The local search returns a near-optimal solution fast whereas the global search performs an incremental search to extract the composition with the minimum number of services in the remaining time. In this section we explain each step of the algorithm, namely: 1) generation of the *Service Match Graph*; 2) calculation of the optimal end-to-end QoS; 3) multi-step graph optimizations and 4) hybrid algorithm.

#### 5.1 Generation of the Service Match Graph

Given a composition request, which specifies the inputs provided by the user as well as the outputs it expects to obtain, and a set of available services, the first step consists of locating all the relevant services that can be part of the final composition, as well as computing all possible matches between their inputs and outputs, according to the semantic model presented in Section 4.1. The output of this step is a *Service Match Graph* that contains many possible valid compositions for the request, as the one represented in Fig. 3. In a nutshell, the generation of the graph is calculated by selecting all invocable services layer by layer, starting with  $So$  in the first layer (the source service whose outputs are the inputs of the request) and terminating with  $Si$  in the last layer (the sink service whose inputs are the outputs of the request) [35].

The pseudocode of the algorithm is shown in Fig. 2. The algorithm runs in polynomial time, selecting  $W_{\text{selected}} \subseteq W$  services at each step. At each layer, the algorithm finds a potential set of relevant services whose inputs are matched



```

1: function SERVICEMATCHGRAPH( $R = \{I_R, O_R\}, W$ )
2:    $C := I_R$ ;  $W' := W$ ;  $W_R := \{So, Si\}$ 
3:    $unmatchedIn := []$ ;  $availCon := I_R$ 
4:   repeat
5:      $W_{selected} = \emptyset$ 
6:      $W_{rel} := \{w \in W' \mid availCon \otimes In_w \neq \emptyset\}$ 
7:      $W_{rel} := W_{rel} \setminus W_R$ 
8:     for all  $w_i = \{In_{w_i}, Out_{w_i}\} \in W_{rel}$  do
9:        $U_{set} := unmatchedIn[w_i]$ 
10:       $M_{set} := C \otimes U_{set}$ 
11:       $unmatchedIn[w_i] := U_{set} \setminus M_{set}$ 
12:      if  $M_{set} = \emptyset$  then
13:         $W_{selected} = W_{selected} \cup w_i$ 
14:         $availCon := availCon \cup Out_{w_i}$ 
15:       $W' := W' \setminus W_{selected}$ 
16:       $W_R := W_R \cup W_{selected}$ 
17:       $C := C \cup availCon$ 
18:       $availCon := \emptyset$ 
19:   until  $W_{selected} = \emptyset$ 
20:   return COMPUTE-GRAPH( $W_R$ )

```

Fig. 2. Algorithm for generating a Service Match Graph from a composition request  $R$  and a set of services  $W$ .

by some outputs generated in the previous layer using the  $\otimes$  operator (L.6). Then, for each potential eligible service, the algorithm checks whether the service is invocable or not (i.e., all its inputs are matched by outputs of previous layers) by checking if all the unmatched inputs of the service are matches. All the inputs that are matched are removed from the unmatched set of inputs for the current service (L.11). If the service is invocable (has no unmatched inputs), it is selected and its outputs are added to the set of the available concepts. In case the service still has some unmatched inputs, these inputs are stored in a map to check it again in the next layer. For example, the first eligible services for the request shown in Fig. 3 are the services in the layer  $L1$ , which correspond with the services whose inputs are fully matched by  $I_R$  (the set of output concepts produced in  $L0$ ). The second eligible services are those services (placed in  $L2$ ) whose inputs are fully matched by the outputs of the previous layers, and so on. The algorithm stops when no more services are added to the set of selected services. Finally, *COMPUTE-GRAPH* computes all possible matches between the outputs and the inputs of the selected services. The output of this process is a complete *Service Match Graph* that can contain cycles, as the one depicted in Fig. 3.

## 5.2 Optimal End-to-End QoS

Once the *Service Match Graph* is computed for a composition request, the next step is to calculate the best end-to-end QoS

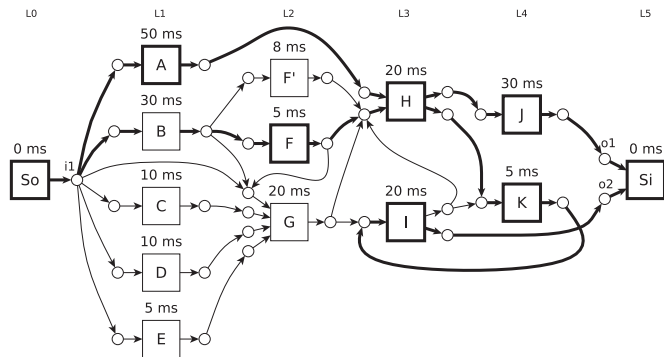


Fig. 3. Graph example with the solution with optimal QoS and minimum number of services highlighted.

```

1: function QoS-UPDATE( $G_S = \{V, E\}$ )
2:   /*  $qos$  is a table indexed by inputs ( $i$ )
3:   associated to their aggregated QoS ( $q$ ) */
4:    $qos[i, q] \leftarrow []$ 
5:   for all  $i \in I, I \subset V$  do
6:      $qos[i] \leftarrow \phi$ 
7:    $queue \leftarrow So$ 
8:   while  $queue \neq \emptyset$  do
9:     /* Queue sorted by aggregated QoS */
10:     $w \leftarrow POP(queue)$ 
11:     $updated = \{\}$ 
12:    for all  $o_w \in Out_w$  do
13:      for all  $i_{w'} \in ch_G(o_w)$  do
14:        if  $V_Q(w) \prec qos[i_{w'}]$  then
15:           $qos[i_{w'}] \leftarrow V_Q(w)$ 
16:           $updated \leftarrow updated \cup w'$ 
17:    for all  $w \in updated$  do
18:      if cost  $w$  has been improved then
19:         $queue \leftarrow INSERT(w, queue)$ 
20:   return  $qos$ 

```

Fig. 4. Dijkstra-based algorithm to compute the best QoS for each input and output in the *Service Match Graph*  $G_S$ .

achievable in the *Service Match Graph*. The optimal end-to-end QoS can be computed in polynomial time using a shortest path algorithm to calculate the best aggregated QoS values for each input and output of the graph, i.e., the best QoS values at which the outputs can be generated and the inputs are matched. In order to compute the optimal QoS, we use a generalized Dijkstra-based label-setting algorithm computed forwards from  $So$  to  $Si$  [36], based on the algebraic model of the QoS presented in Section 4. The optimality of the algorithm is guaranteed as long as the function defined to aggregate the QoS values ( $\otimes$ ) is monotonic, in order to satisfy the principle of optimality. A proof can be found in [37].

Fig. 4 shows the pseudocode of the generalized Dijkstra-based label-setting algorithm. The algorithm starts assigning infinite QoS cost to each input in the graph in the table  $qos$ . An infinite cost for an input means that the input is still not resolved. The first service to be processed is  $So$ . Each time a service  $w$  is processed from the queue, the best accumulated QoS cost of each input  $i_{w'}$  matched by the outputs of the service  $w$  is recalculated. If there is an improvement (i.e., a match with a better QoS is discovered) the affected service is stored in  $updated$  to recompute its new aggregated QoS. Finally, for each service  $w \in updated$ , we recompute its aggregated QoS using the updated values of each affected input. If the QoS has been improved, the service is added to the queue to expand it later.

## 5.3 Graph Optimizations

Finding the composition with the minimum number of services is a very hard combinatorial problem which, in most cases, has a very large search space, mainly determined by the size of the *Service Match Graph*. In order to improve the scalability with the number of services, we apply a set of admissible optimizations to reduce the search space. At each pass, the algorithm analyzes different criteria to identify services that are redundant or can be substituted by better ones, so the size of the graph decreases monotonically. The different passes that are sequentially applied are: 1) elimination of services that do not contribute to the outputs of the request; 2) pruning of services that lead to suboptimal QoS; 3) combination of interface (inputs/outputs) and QoS

equivalent services; and 4) replacement of interface and QoS dominated services. These optimizations are an extension of the optimizations presented in [30] to support QoS.

The **first pass** selects the set of reachable services in the *Service Match Graph*. Starting from the inputs of  $Si$ , it selects all those services whose outputs match any inputs of  $Si$ . This step is repeated with the new services until the empty set is selected. Those services that were not selected do not contribute to the expected outputs of the composition and can be safely removed from the graph.

The **second pass** prunes the services of the graph that are suboptimal in terms of QoS, i.e., they cannot be part of any optimal QoS composition. To do so, we compute the maximum admissible QoS bound for each input in the graph. In a nutshell, the maximum bound of the inputs of a service  $w$  can be calculated by selecting the maximum QoS bound among the bounds of all inputs matched by the outputs of the service  $w$  and subtracting the QoS of  $w$ . This can be recursively defined as:

$$\begin{aligned} \max_Q^i(i_w) &= \begin{cases} V_Q(w) \ominus F_Q(w) & \text{if } Out_w = \emptyset \\ \max_{\forall o_w, \forall i_{w'} \in ch_G(o_w)} (\max_Q^i(i_{w'}) \ominus F_Q(w)) & \text{if } Out_w \neq \emptyset. \end{cases} \end{aligned}$$

The value of  $\max_Q^i$  for each input in the graph can be easily calculated by propagating the bounds from  $Si$  to  $So$ . For example, in Fig. 1, we start computing the maximum bound of the inputs of  $Si$  ( $xsd:boolean$ ). Since  $Si$  has no outputs,  $\max_Q^i(xsd:boolean)$  is calculated as  $V_Q(Si) \ominus F_Q(Si) = 410 \text{ ms} - 0 \text{ ms}$ . Then, we select all the services whose outputs match  $xsd:boolean$ . In this case there is just one service, *ML Predictor Service*. The bounds of its inputs are now computed by subtracting out the  $F_Q(\textit{ML Predictor Service})$  from the maximum bound of the inputs that this service matches. Since there is just one input matched ( $xsd:boolean$  from  $Si$ ) whose bound is  $410 \text{ ms}$ , we have  $\max_Q^i(i) = 410 - 210 \text{ ms} = 200 \text{ ms}$  for each input  $i$  of the service. In the next step, we have three services that match the new calculated inputs (*Free Geoloc Service*, *Premium Geoloc Service* and *Transaction Service*). The maximum bounds of the inputs of these services are  $200 - 180 \text{ ms} = 20 \text{ ms}$ ,  $200 - 40 \text{ ms} = 160 \text{ ms}$  and  $200 - 130 \text{ ms} = 70 \text{ ms}$  respectively. Note that, since the maximum bound of *Transaction Service* is  $70 \text{ ms}$ , the service *Secure Payment* is out of the bounds (its output QoS is  $80 \text{ ms}$ ), so it can be safely pruned.

The **third and the forth pass** analyze service equivalences and dominances in the *Service Match Graph*. It is very frequent to find services from different providers that offer similar services with overlapping interfaces (inputs/outputs). In scenarios like this, it is easy to end up with large *Service Match Graph* that make very hard to find optimal compositions in reasonable time. One way to reduce the complexity without losing information is to analyze the interface equivalence and dominance between services in order to *combine* those that are equivalent, or replace those that are dominated in terms of the interface they provide and the QoS they offer. In a nutshell, we check three objectives to compare services: the *amount* of information they need to be invoked (inputs), the *amount* of information they

```

1: function LOCAL-SEARCH( $G_S = \{V, E\}$ )
2:   return LSBT( $G_S, In_{Si}, \{Si\}$ )
3:
4: function LSBT( $G_S, unresolved, services$ )
5:   if  $unresolved = \emptyset$  then return  $G_S$ 
6:    $servs \leftarrow \text{RANK-RESOLVERS}(unresolved)$ 
7:   for each  $w \in servs$  do
8:      $resolved \leftarrow \{\}$ 
9:      $matched \leftarrow Out_w \otimes unresolved$ 
10:    for each input  $\in matched$  do
11:      if  $\neg \text{CYCLE}(G_S, w, input)$  then
12:         $resolved \leftarrow resolved \cup input$ 
13:    if  $resolved \neq \emptyset$  then
14:       $unresolved \leftarrow unresolved \setminus resolved$ 
15:      if  $w \notin services$  then
16:         $unresolved \leftarrow unresolved \cup In_w$ 
17:       $G'_S \leftarrow \text{RESOLVE}(G_S, w, resolved)$ 
18:       $services \leftarrow services \cup w$ 
19:       $result \leftarrow \text{LSBT}(G'_S, unresolved, services)$ 
20:      if  $result \neq fail$  then return  $result$ 
  return fail

```

Fig. 5. Local search algorithm to extract a composition from a graph.

return (outputs), and their QoS. If a set of services are equal in all objectives, they are equivalent and they can be combined into an abstract service with several possible implementations. If a service is equal in all objectives and at least better in one objective (it requires less information to be invoked, produces more information or has a better QoS), then the service *dominates* the other service. A more detailed description of the interface and dominance optimizations is described in [30].

Note that optimizations are applied right before all semantic matches are computed in the *Service Match Graph*, since the optimizations are based on the analysis of the I/O matches among services. For this reason, they cannot be applied during the calculation of the graph (this would require to precompute in advance missing relations during the graph generation, which does not provide any benefit as this is what the *Service Match Graph* generation algorithm already does). On the other hand, optimizations are applied sequentially to save computation time, since the number of services in the graph decreases monotonically in each step. In order to take advantage of this, faster optimizations are applied first so that the slower optimizations in the pipeline can work with a reduced set of services.

## 5.4 Hybrid Algorithm

Each service in the composition graph may have different services that match each input, thus there may exist multiple combinations of services that satisfy the composition request with the same or different QoS. The goal of the hybrid search is to extract good solutions from the composition graph, optimizing the total number of involved services in the composition and guaranteeing the optimal QoS. Thus, for each input we select just one service of the graph to match that input, until the best combination is found. The hybrid search performs a local search to extract a good solution and in the remaining time, it tries to improve the solution by running a global search.

Fig. 5 shows the pseudocode of the **local search** strategy. The algorithm starts with a composition graph, the inputs of the service  $Si$  marked as unresolved (the expected outputs of the request) and the service  $Si$  selected to be part of the solution. An *unresolved input* is an input that can be



```

1: function CYCLE( $G_S = \{V, E\}, w, i_{w'}$ )
2:    $W_{visited} \leftarrow \{w'\}$ 
3:    $W_{new} \leftarrow \{w'\}$ 
4:   while  $W_{new} \neq \emptyset$  do
5:      $W_{reached} \leftarrow \{\}$ 
6:     for all  $w_n \in W_{new}$  do
7:       for all  $o_{w_n} \in Out_{w_n}$  do
8:         for all  $i_{w'_n} \in ch_{G_S}(o_{w_n})$  do
9:           if  $d_{G_S}^{-1}(i_{w'_n}) = 1 \wedge w'_n \notin W_{visited}$  then
10:            if  $w'_n = w$  then return true
11:             $W_{reached} \leftarrow W_{reached} \cup w'_n$ 
12:    $W_{new} \leftarrow W_{reached}$ 
13:    $W_{visited} \leftarrow W_{visited} \cup W_{new}$ 
14: return false

```

Fig. 6. Naïve breadth-first-search algorithm to check whether using the service  $w$  to resolve the input  $i_{w'}$  of a service  $w'$  leads to a cycle.

matched by many different outputs but no decision has been made yet. Using the list of the unresolved inputs to be matched, the method RANK-RESOLVERS returns a list of services that match any of the unresolved inputs. Services are ranked according to the number of unresolved inputs that match, so the service that matches more inputs is considered first to be part of the solution. Then, for each input that the selected service can match, the method CYCLE performs a forward search to check if resolving the selected input with that service leads to a cycle. For example, in Fig. 3, if we select the service  $K$  to match the input of  $I$  after having decided to resolve the input of  $K$  with the service  $I$ , we end up with an invalid composition, so  $K$  is an invalid resolver for  $I$  and it must be discarded. Once all resolvable inputs are collected in *resolved*, the method RESOLVE creates a copy of the current graph where the inputs in *unresolved* are matched only by the selected service, i.e., any other match between any output from a different service to that input is removed from the graph. If the selected service was not already selected, then all its inputs are then marked as unresolved and a recursive call to LSBT is performed to select a new service to resolve the remaining inputs, until a solution is found. If a dead end is reached (a solution that has no services to resolve the remaining inputs without cycles) the algorithm backtracks to a previous state to try a different service (L.7).

An implementation of the CYCLE method is provided in 6. The algorithm performs a *look-ahead* check in a breadth-first fashion to determine whether matching the selected input  $i$  with an output of the service  $w$  leads to a cyclic dependency. This is done by traversing only the resolved matches, i.e., inputs that are matched by just one output of a service, until the selected service  $w$  is reached, proving the existence of a cycle. A more memory efficient implementation of the cycle algorithm can be done using the *Tarjan's strongly connected components* algorithm [38], stopping at the first strongly connected component detected.

After the local search is used to find a good solution, the **global search** is performed in the remaining time to obtain a better solution by exhaustively exploring the space of possible solutions. In a nutshell, this algorithm works as follows: Given a *Service Match Graph*  $G_S$ , with some unresolved inputs, which initially are the inputs of the service  $S_i$ , the algorithm selects an input to be resolved and for each service candidate that can be used to resolve that input, it generates a copy of the graph  $G_S$  but with the input

```

1: function GLOBAL-SEARCH( $G_S$ )
2:    $qos[i, q] \leftarrow QoS-UPDATE(G_S)$ 
3:    $max \leftarrow COMPUTE-V_Q(S_i, qos)$ 
4:    $W_{sel} \leftarrow \{S_i\}$ 
5:   /*  $I_{un}$  is a key-value table where the keys are
6:    unresolved inputs and the values their QoS bounds */
7:   for  $i_{S_i} \in In_{S_i}$  do
8:      $I_{un}[i_{S_i}] \leftarrow [qos[i_{S_i}], max]$ 
9:   /* Queue sorted by  $|W_{sel}|$  */
10:   $queue \leftarrow INSERT(\langle G_S, I_{un}, qos, W_{sel} \rangle, queue)$ 
11:  while  $queue \neq \emptyset$  do
12:     $\langle G_S, I_{un}, qos, W_{sel} \rangle \leftarrow POP(queue)$ 
13:    if  $I_{un} = \emptyset$  then return  $G_S$ 
14:     $input \leftarrow SELECT(I_{un})$ 
15:     $[min, max] \leftarrow I_{un}[input]$ 
16:    for all  $w \in RESOLVERS(input, [min, max])$  do
17:      if  $\neg CYCLE(G_S, w, input)$  then
18:         $G'_S \leftarrow RESOLVE(G_S, w, \{input\})$ 
19:         $I'_{un} \leftarrow REMOVE(i, I_{un})$ 
20:         $qos' \leftarrow qos$ 
21:        if  $COMPUTE-V_Q(w, qos') \succ min$  then
22:           $qos' \leftarrow QoS-UPDATE(G'_S)$ 
23:        if  $w \notin W_{sel}$  then
24:           $W'_{sel} \leftarrow W_{sel} \cup w$ 
25:           $max' \leftarrow max \oplus F_Q(w)$ 
26:          for  $i_w \in In_w$  do
27:             $min' \leftarrow qos'[i_w]$ 
28:             $I'_{un}[i_w] \leftarrow [min', max']$ 
29:           $queue \leftarrow INSERT(\langle G'_S, I'_{un}, qos', W'_{sel} \rangle, queue)$ 
return fail

```

Fig. 7. Global search algorithm to extract the optimal composition.

resolved (i.e., the selected service is the only one that matches the unresolved input). The algorithm enqueues each new graph to be expanded again, and repeats the process by extracting the graph with the minimum number of services from the queue, until it eventually finds a graph with no unresolved inputs.

Fig. 7 shows the pseudocode of the global search algorithm. The algorithm starts computing the optimal QoS of the graph with the method *QoS-UPDATE*. This method returns a key-value table  $qos[i, q]$  where each key corresponds with an input  $i$  of the graph, and each value  $q$  its optimal aggregated QoS  $q = V_Q^{in}(i)$ . Then, the inputs of the service  $S_i$  of the graph are added to  $I_{un}$  to mark them as unresolved (L.8). In order to minimize the number of possible candidates for each unresolved input, we compute and propagate a range of valid QoS values, called QoS bounds, and defined as an interval  $[min, max]$ . These bounds determine the range of valid accumulated QoS values of the outputs that can be used to match each of the unresolved inputs without exceeding the optimal end-to-end QoS of the final composition. The *min* value is the optimal QoS for the input, i.e., there is no output in the graph that can match the input with a lower QoS, whereas the *max* value is the maximum QoS value supported. If this bound is exceeded, the total aggregated QoS of the composition worsens. For example, in Fig. 1, the bounds of the input *ont4:ClientID* of the service *Premium Geoloc Service* are  $[20, 160 \text{ ms}]$ . If we exceed the min bound (20ms), the output QoS of the service gets worse ( $> 60 \text{ ms}$ ), which also affects the optimal QoS of the input *ont1:Location*. However, as long as the max bound is not exceeded ( $\leq 160 \text{ ms}$ ), the optimal accumulated QoS of the *ML Predictor Service* would not be affected.

The method *COMPUTE-V<sub>Q</sub>* is used to compute the value of the  $V_Q$  function (Eq. (1)) using the best QoS values of

inputs, stored in  $qos$  ( $qos[i] = V_Q^{in}(i)$ ). A tuple  $\langle G_S, I_{un}, qos, W_{sel} \rangle$ , where  $G_S$  is the current graph,  $I_{un}$  are the unresolved inputs of  $G_S$ ,  $qos$  is the best aggregated QoS values for each input in  $G_S$  and  $W_{sel}$  is the set of the selected services, defines the components of a partial solution. Each partial solution is stored in a priority queue, which is sorted by the number of services  $W_{sel}$ . This allows an exploration of the search space in a breadth-first fashion, so the solution with the minimum number of services is always expanded first. At each iteration, a partial solution is extracted from the queue to be refined (L.12). If the partial solution has no unresolved inputs, the solution is complete, and has the minimum number of services. If the partial solution still has some unresolved inputs, it is refined by selecting an unresolved input with the method *SELECT*. This method selects the input to be resolved, using a *minimum-remaining-values* heuristic. This heuristic selects always the input with less resolvers (services candidates) in order to minimize the branching factor. The list of services that can match the selected input with a total aggregated QoS value within the  $[min, max]$  bound is calculated with the method *RESOLVERS*. For each valid service, the algorithm performs a *look-ahead* search to check whether using the current service to resolve the selected input leads to an unavoidable cycle. If so, the service is prematurely discarded to save computation time and space. If it does not lead to a cycle, then a copy of the graph ( $G'_S$ ) with the selected input resolved is generated, and the input is also removed from the set of unresolved inputs. Using the optimal aggregated QoS values for the inputs of the graph, stored in  $qos$ , the algorithm computes the aggregated QoS value of the service  $w$ . If this value is worse than the  $min$  bound ( $COMPUTE\_V_Q(w, qos') \succ min$ ), then the aggregated QoS value of some inputs and outputs of the graph may be affected. Thus, a repropagation of the QoS values for each input and output is computed again over the new graph  $G'_S$  (L.22). For example, if the *Business Service Info* increments its response time to 40 ms, a repropagation is required to recompute the accumulated QoS of all the services that may be affected. In this case, the *Premium Geoloc Service* increments its accumulated QoS cost from 60 to 80 ms, as well as the optimal QoS of the *ont1:Location*.

Finally, if the current service is not part of the current solution, its inputs are added to the unresolved table, and a new bound for each input is computed. The  $min$  bound corresponds with the optimal value, which is stored in  $qos'$ . In order to compute the  $max$  bound, we need to *subtract* the QoS of the selected service ( $F_Q(w)$ ) from the  $max$  bound of the resolved input, using the operator  $\ominus$  (L.25). This new partial solution is inserted in the queue to be expanded later on.

## 6 EVALUATION

In order to evaluate the performance of the proposed approach, we conducted two different experiments. In the first experiment, we evaluated the approach using the datasets of the Web Service Challenge 2009-2010 [39]. The goal of this first experiment was to evaluate the performance and scalability of the proposed approach on large-scale service repositories. In the second experiment, we tested the algorithm with five random datasets in order

TABLE 2  
Validation with the WSC 2009-2010

	D-01	D-02	D-03	D-04	D-05
#Services in the dataset	572	4,129	8,138	8,301	15,211
Validation with Response Time					
Optimal Response Time (ms)	500	1,690	760	1,470	4,070
#Graph services	81	141	154	331	238
#Graph services (opt)	21	57	15	160	126
Local Search					
#Services	5	20	10	40	32
Time (s)	0.613	0.988	2.608	7.767	2.920
Global Search					
#Services	5	20	10	-	32
Time (s)	0.617	1.580	2.613	-	24.971
Validation with Throughput					
Optimal Throughput (inv/s)	15,000	6,000	4,000	4,000	4,000
#Graph services	81	141	154	331	238
#Graph services (opt)	10	43	90	156	69
Local Search					
#Services	5	20	15	62	31
Time (s)	0.343	1.173	1.933	8.571	2.560
Global Search					
#Services	5	20	10	-	32
Time (s)	0.345	1.246	2.085	-	119.322

to better analyze the differences of the performance between the local and the global search. All tests were executed with a time limit of 5 min. Solutions produced by our algorithm are represented as *Service Composition Graphs* (no BPEL was generated).

### 6.1 Web Service Challenge 2009-2010 Datasets

The datasets of the Web Service Challenge 2009-2010 range from 572 to 15,211 services with two different QoS properties: response time and throughput. Table 2 shows the results obtained for each dataset and for each QoS property. The response time is the average time (measured in milliseconds) that a service takes to respond to a request. The throughput, as defined in the WSC, is the average ratio of invocations per second supported by a service.

Row *#Graph services* shows the number of services of the composition graph and *#Graph services (opt)* the number of services after applying the graph optimizations. As can be seen, the optimizations reduce, on average, by 64 percent the number of services in the initial composition graph. This indicates that equivalence and dominance analysis of the QoS and the functionality of services is a powerful technique to reduce the search space in large scale problems. Rows *Local search* and *Global search* show the number of services of the solution obtained with each respective method as well as the total amount of time spent in the search. The global search found the best solution for each dataset and for each QoS property, except for the dataset 04, where the composition with the minimum number of services could not be found due to combinatorial explosion. However, in those cases, the local search strategy is able to find an alternative solution very fast. Note also that, in many cases, the local search obtains the best solution (comparing it with the global search) except for the throughput in datasets 03 and 05.

We have compared our approach with the top-3 of the Web Service Challenge 2010 [40]. Table 3 shows this comparison following the same format and the same rules of the Web Service Challenge. The format, rules and other details of the challenge are described in [40]. Third and forth columns show the response time and the throughput obtained for each dataset. Note that, since all these algorithms

TABLE 3  
Comparison with the Top 3 WSC 2010

		R.Time	Through.	Min. Serv.	Time (ms)
D-01	CAS [4]	500	15,000	5	78
	RUG [6]	500	15,000	10	188
	Tsinghua [5]	500	15,000	9	109
	<b>Our approach</b>	500	15,000	5	956
D-02	CAS [4]	1,690	6,000	20	94
	RUG [6]	1,690	6,000	40	234
	Tsinghua [5]	1,690	6,000	36	140
	<b>Our approach</b>	1,690	6,000	20	2,171
D-03	CAS [4]	760	4,000	10	78
	RUG [6]	760	4,000	11	234
	Tsinghua [5]	760	4,000	18	125
	<b>Our approach</b>	760	4,000	10	4,693
D-04	CAS [4]	1,470	4,000	73	156
	RUG [6]	1,470	4,000	133	390
	Tsinghua [5]	1,470	4,000	133	188
	<b>Our approach</b>	1,470	4,000	40	16,338
D-05	CAS [4]	4,070	4,000	32	63
	RUG [6]	4,070	4,000	4,772	907
	Tsinghua [5]	4,070	4,000	4,772	531
	<b>Our approach</b>	4,070	4,000	30	122,242

minimize a single QoS, these values are computed by executing the algorithm twice, one for each QoS. Unfortunately, the results provided by the WSC organization in [40] show only the minimum number of services for both executions (fifth column). Thus, the number of services obtained for both the response time and throughput is unknown, which makes it hard to compare with our results. Even so, using the same evaluation criteria, our approach obtains the optimal QoS for the response time and the throughput, and also improves the number of services in D-04 (40 versus 73) and D-05 (30 versus 32) with respect to the solutions obtained by the winner of the challenge (the minimum number of services obtained for each dataset is highlighted). The last column shows the total execution time of each algorithm. The total time includes the time spent to obtain the solution for the response time and for the throughput.

Our approach takes, in general, more time to obtain a solution. However, it should be noted that we show the best results achieved by the hybrid approach, i.e., if the global search improves the solution of the local search, we show that solution along with the time taken by the global search. Anyway, the local search always provide a first good solution very fast. For example, as can be seen in Table 2, the optimal solution for D-05 has 30 services and has been obtained in 119.322 s, but the local search obtained a solution with 31 services in 2.56 s, still better than the solution with 32 services obtained by [4] (Table 3). Moreover, it should also be noted that the problem of finding the optimal composition with minimum number of services and optimal QoS is much harder than just optimizing the QoS objective function, which is the problem solved by the participants of the WSC 2010. Although the problem is intractable and requires exponential time, it can be optimally solved for many particular instances in a reasonable amount of time using adequate optimizations even in large datasets as shown in Tables 2 and 5. This is one of the main reasons why a combination of a local and global search can achieve good results in a wide variety of situations, in contrast with pure greedy strategies or with pure global optimization algorithms.

TABLE 4  
Detailed Comparison with [7]

		D-01	D-02	D-03	D-04	D-05
Chen et al.	R. Time	500	1,690	760	1,470	4,070
	Services	8	21	10	42	33
Our approach	R. Time	500	1,690	760	1,470	4,070
	Services	5	20	10	40	32
Chen et al.	Throughput	15,000	6,000	4,000	2,000	4,000
	Services	5	20	21	40	30
Our approach	Throughput	15,000	6,000	4,000	4,000	4,000
	Services	5	20	10	62	30

We also compare the results obtained with Chen and Yan [7], who offer a detailed analysis of their results. This comparison is shown in Table 4. Solutions are compared according to their QoS and number of services. A solution is better if 1) its overall QoS is better or 2) has the same QoS but less services. The results show that our algorithm always gets same or better results. Concretely, it finds solutions with optimal QoS and less services in D-01, D-02, D-04 and D-05 (response time), and D-03 (throughput). It also finds a solution with a better QoS (4000 inv/s versus 2000 inv/s) in D-04 (throughput).

## 6.2 Randomly Generated Datasets

Although the global search is able to obtain solutions with a lower number of services, a first look at the results with the WSC dataset might suggest that the difference of both strategies is not very significant, as most of the obtained solutions have the same number of services. However, this may be due to a bias in the repository, since all the datasets of the WSC are generated using the same random model. In order to better evaluate and characterize the performance of the hybrid algorithm, we generated a new set of five random datasets that range from 1,000 to 9,000 services. These datasets are available at [https://wiki.citius.usc.es/inv:downloadable\\_results:ws-random-qos](https://wiki.citius.usc.es/inv:downloadable_results:ws-random-qos). Table 5 shows the solutions obtained.

We found that in these datasets, the solutions obtained with the global search strategy are, on average,  $\approx 16$  percent

TABLE 5  
Validation with Random Datasets

	R-01	R-02	R-03	R-04	R-05
#Services in the dataset	1,000	3,000	5,000	7,000	9,000
Validation with Response Time					
Optimal Response Time (ms)	1,430	975	805	1,225	1,420
#Graph Services	54	168	285	383	499
#Graph Services (opt)	22	50	54	56	99
Local Search	#Services	7	18	20	15
	Time (s)	0.183	0.403	0.422	0.515
Global Search	#Services	7	14	15	16
	Time (s)	0.243	0.767	4.088	0.740
Validation with Throughput					
Optimal Throughput (inv/s)	1,000	2,500	1,500	2,000	2,500
#Graph Services	54	168	285	383	499
#Graph Services (opt)	19	46	133	116	103
Local Search	#Services	7	17	24	19
	Time (s)	0.072	0.143	0.606	0.732
Global Search	#Services	7	12	12	15
	Time (s)	0.155	0.310	2.479	1.485



smaller than the ones obtained with the local search, whereas the differences in search time are less pronounced than in the previous experiment. These findings suggest that the performance of each strategy highly depends on the underlying structure of the service repository, which is mostly determined by the number of services and the existing matching relations.

In order to test whether these differences are statistically significant or not, we conducted a nonparametric test using the *binomial sign test* for two dependent samples with a total of 20 datasets (5 WSC w/response time + 5 WSC w/throughput + 5 Random w/response time + 5 Random w/throughput). The null hypothesis was rejected with  $p\text{-value} \approx 0.01$  [41], meaning that both strategies (local and global search) find significantly different solutions. Thus, a hybrid strategy can perform better in many different scenarios, since it achieves a good tradeoff between quality and execution time.

This evaluation shows that, on one hand, the combination of local and global optimization is a general and powerful technique to extract optimal compositions in diverse scenarios, as it brings the best of both worlds. This is specially important when only a little or nothing is known concerning the structure of the underlying repository of services. On the other hand, the results obtained with the Web Service Challenge 2009-2010 show that the hybrid strategy performs better than the state-of-the-art, obtaining solutions with less services and optimal QoS.

## 7 CONCLUSIONS

In this paper we have presented a hybrid algorithm to automatically build semantic input-output based compositions minimizing the total number of services while guaranteeing the optimal QoS. The proposed approach combines a set of graph optimizations and a local-global search to extract the optimal composition from the graph. Results obtained with the Web Service Challenge 2009-2010 datasets show that the combination of graph optimizations with a local-global search strategy performs better than the state-of-the-art, as it obtained solutions with less services and optimal QoS. Moreover, the evaluation with a set of randomly generated datasets shows that the hybrid strategy is well suited to perform compositions in diverse scenarios, as it can achieve a good tradeoff between quality and execution time.

## ACKNOWLEDGMENTS

This work was supported by the Spanish Ministry of Economy and Competitiveness (MEC) under grant TIN2014-56633-C3-1-R and the Galician Ministry of Education under the project CN2012/151. Pablo Rodríguez-Mier is supported by an FPU Grant from the MEC (ref. AP2010-1078).

## REFERENCES

- [1] G. Alonso, F. Casati, H. Kuno, and V. Machiraju, *Web Services, series Data-Centric Systems and Applications*. New York, NY, USA: Springer, 2004.
- [2] J. Rao and X. Su, "A survey of automated web service composition methods," in *Proc. Int. Workshop Semantic Web Services Web Process Composition*, 2004, pp. 43–54.
- [3] A. Strunk, "QoS-aware service composition: A survey," *Proc. IEEE 8th Eur. Conf. Web Services*, 2010, pp. 67–74.
- [4] W. Jiang, S. Hu, Z. Huang, Z. Liu, and Q. D. Handler, "Two-phase graph search algorithm for QoS-Aware automatic service composition," in *Proc. IEEE Int. Conf. Service-Oriented Comput. Appl.*, 2010, pp. 1–4.
- [5] Y. Yan, B. Xu, Z. Gu, and S. Luo, "A QoS-driven approach for semantic service composition," in *Proc. IEEE Int. Conf. Commerce Enterprise Comput.*, 2009, pp. 523–526.
- [6] M. Aiello, E. E. Khoury, A. Lazovik, and P. Ratelband, "Optimal QoS-aware web service composition," in *Proc. IEEE Int. Conf. Commerce Enterprise Comput.*, 2009, pp. 491–494.
- [7] M. Chen and Y. Yan, "Redundant service removal in QoS-aware service composition," in *Proc. IEEE Int. Conf. Web Services*, 2012, pp. 431–439.
- [8] Z. Zheng, H. Ma, M. R. Lyu, and I. King, "Collaborative web service QoS prediction via neighborhood integrated matrix factorization," *IEEE Trans. Services Comput.*, vol. 6, no. 3, pp. 289–299, Jul.–Sep. 2013.
- [9] S. Dustdar and W. Schreiner, "A survey on web services composition," *Int. J. Web Grid Services*, vol. 1, no. 1, pp. 1–30, 2005.
- [10] P. Bertoli, M. Pistore, and J. Hoffmann, "Web service composition as planning, revisited: In between background theories and initial state uncertainty," in *Proc. 22nd Nat. Conf. Artif. Intell.*, pp. 1013–1018, 2007.
- [11] S. R. Ponnekanti and A. Fox, "SWORD: A developer toolkit for web service composition," in *Proc. 11th World Wide Web Conf.*, 2002.
- [12] M. Carman, L. Serafini, and P. Traverso, "Web service composition as planning," in *Proc. Workshop Planning Web Services*, 2003, pp. 1636–1642.
- [13] E. Sirin, B. Parsia, D. Wu, J. Hendler, and D. Nau, "HTN planning for web service composition using SHOP2," *Web Semantics: Sci. Services Agents World Wide Web*, vol. 1, no. 4, pp. 377–396, 2004.
- [14] E. Sirin and B. Parsia, "Planning for semantic web services," in *Proc. Semantic Web Services Workshop 3rd Int. Semantic Web Conf.*, 2004, pp. 33–40.
- [15] M. Klusch, A. Gerber, and M. Schmidt, "Semantic web service composition planning with OWLS-Xplan," in *Proc. AAAI Fall Symp. Semantic Web Agents*, 2005.
- [16] R. Akkiraju, B. Srivastava, A. A. Ivan, R. Goodwin, and T. Syeda-Mahmood, "SEMAPLAN: Combining planning with semantic matching to achieve web service composition," in *Proc. IEEE Int. Conf. Web Services*, 2006, pp. 37–44.
- [17] O. Hatzi, D. Vrakas, M. Nikolaidou, N. Bassiliades, D. Anagnostopoulos, and I. Vlahavas, "An integrated approach to automated semantic web service composition through planning," *IEEE Trans. Services Comput.*, vol. 5, no. 3, pp. 319–332, 3rd Quarter 2012.
- [18] S. Oh, D. Lee, and S. R. T. Kumara, "Web service planner (WSPR): An effective and scalable web service composition algorithm," *Int. J. Web Service Res.*, vol. 4, no. 1, pp. 1–22, 2007.
- [19] L. Zeng, B. Benatallah, A. H. Ngu, M. Dumas, J. Kalagnanam, and H. Chang, "QoS-Aware middleware for web services composition," *IEEE Trans. Softw. Eng.*, vol. 30, no. 5, pp. 311–327, May, 2004.
- [20] T. Yu and K.-J. Lin, "Service selection algorithms for web services with end-to-end QoS constraints," *Inform. Syst. e-Business Manage.*, vol. 3, no. 2, pp. 103–126, 2005.
- [21] R. Berbner, M. Spahn, N. Repp, O. Heckmann, and R. Steinmetz, "Heuristics for QoS-aware web service composition," in *Proc. Int. Conf. Web Services*, 2006, pp. 72–82.
- [22] M. Alrifai and T. Risse, "Combining global optimization with local selection for efficient QoS-aware service composition," in *Proc. 18th Int. Conf. World Wide Web*, 2009, pp. 881–890.
- [23] G. Zou, Q. Lu, Y. Chen, R. Huang, Y. Xu, and Y. Xiang, "QoS-Aware dynamic composition of web services using numerical temporal planning," *IEEE Trans. Services Comput.*, vol. 7, no. 1, pp. 18–31, Jan.–Mar. 2014.
- [24] G. Canfora, M. D. Penta, R. Esposito, and M. L. Villani, "An approach for QoS-aware service composition based on genetic algorithms," in *Proc. Genetic Evol. Comput. Conf.*, 2005, pp. 1069–1075.
- [25] H. Wada, J. Suzuki, Y. Yamano, and K. Oba, "E<sup>3</sup>: A multiobjective optimization framework for SLA-aware service composition," *IEEE Trans. Services Comput.*, vol. 5, no. 3, pp. 358–372, Jul. 2012.
- [26] W. Jiang, S. Hu, and Z. Liu, "Top K query for QoS-Aware automatic service composition," *IEEE Trans. Services Comput.*, vol. 7, no. 4, pp. 681–695, Oct. 2014.
- [27] L. Barakat, S. Miles, I. Poernomo, and M. Luck, "Efficient multi-granularity service composition," in *Proc. IEEE Int. Conf. Web Services*, 2011, pp. 227–234.

- [28] F. Wagner, F. Ishikawa, and S. Honiden, "QoS-aware automatic service composition by applying functional clustering," in *Proc. IEEE Int. Conf. Web Services*, 2011, pp. 89–96.
- [29] P. Rodríguez-Mier, M. Mucientes, and M. Lama, "Automatic web service composition with a heuristic-based search algorithm," in *Proc. IEEE Int. Conf. Web Services*, 2011, pp. 81–88.
- [30] P. Rodríguez Mier, C. Pedrinaci, M. Lama, and M. Mucientes, "An integrated semantic web service discovery and composition framework," *IEEE Trans. Services Comput.*, 2015 (DOI 10.1109/TSC.2015.2402679).
- [31] M. Paolucci, T. Kawamura, T. R. Payne, and K. Sycara, "Semantic matching of web services capabilities," in *Proc. 1st Int. Semantic Web Conf.*, 2002, pp. 333–347.
- [32] C. Pedrinaci, D. Liu, M. Maleshkova, D. Lambert, J. Kopecky, and J. Domingue, "iServe: A linked services publishing platform," *In Ontology Repositories and Editors for the Semantic Web*, 2010.
- [33] J. Cardoso, A. Sheth, J. Miller, J. Arnold, and K. Kochut, "Quality of service for workflows and web service processes," *Web Semantics*, vol. 1, pp. 281–308, 2004.
- [34] B. Carré, *Graphs and Networks*. New York, NY, USA: Oxford Univ. Press, 1979.
- [35] P. Rodríguez-Mier, M. Mucientes, J. C. Vidal, and M. Lama, "An optimal and complete algorithm for automatic web service composition," *Int. J. Web Service Res.*, vol. 9, no. 2, pp. 1–20, 2012.
- [36] P. Rodríguez-Mier, M. Mucientes, and M. Lama, "A dynamic QoS-Aware semantic web service composition algorithm," in *Proc. 10th Int. Conf. Service-Oriented Comput.*, 2012, pp. 623–660.
- [37] J. L. Sobrinho, "Algebra and algorithms for QoS path computation and hop-by-hop routing in the Internet," in *Proc. 20th Annu. Joint Conf. IEEE Comput. Commun. Soc.*, 2002, vol. 10, no. 4, pp. 541–550.
- [38] R. Tarjan, "Depth-first search and linear graph algorithms," *SIAM J. Comput.*, vol. 1, no. 2, pp. 146–160, 1972.
- [39] S. Kona, A. Bansal, M. B. Blake, S. Bleul, and T. Weise, "WSC-2009: A quality of service-oriented web services challenge," in *Proc. IEEE Int. Conf. Commerce Enterprise Comput.*, 2009, pp. 487–490.
- [40] T. Weise. (2010). Web Service Challenge 2010 [Online]. Available: [http://www.it-weise.de/documents/files/W2010WSC\\_pres.pdf](http://www.it-weise.de/documents/files/W2010WSC_pres.pdf)
- [41] I. Rodríguez-Fdez, A. Canosa, M. Mucientes, and A. Bugarín, "STAC: A web platform for the comparison of algorithms using statistical tests," in *Proc. IEEE Int. Conf. Fuzzy Syst.*, 2015.



**Pablo Rodríguez-Mier** received the BSc and MSc degrees in computer science from the Universidade de Santiago de Compostela (USC) in 2011. He is currently working toward the PhD degree at CiTIUS, USC, Spain. His research interests include heuristic search and automatic web service composition.



**Manuel Mucientes** is an associate professor in computer science and artificial intelligence within the CiTIUS of the Universidade de Santiago de Compostela. He has authored or coauthored more than 70 papers in international journals, book chapters, and conferences. His current research interests include evolutionary algorithms, genetic fuzzy systems, motion planning and control in robotics, visual SLAM, web services, and process mining.



**Manuel Lama** received the PhD degree in computer science from the University of Santiago de Compostela in 2000. He is an associate professor at CiTIUS, Universidade de Santiago de Compostela. His research interests focuses on discovery and composition of web services, semantic annotation, and process mining.