

QoS Dependency Modeling for Composite Systems

Khayyam Hashmi, Zaki Malik, Abdelkarim Erradi, and Abdelmounaam Rezgui

Abstract—Software as a service is a well accepted software deployment and distribution model that has grown exponentially in the last few years. One of the biggest benefits of SaaS is the automated composition of different services in a composite system. It allows users to automatically find and bind these services (to maximize the productivity), meeting both functional and non-functional requirements. In this paper, we present a framework for modeling the dependency relationships among different Quality of Service parameters of the component services. Our proposed approach considers the different invocation patterns of component services, and presents a service composition framework to model the dependencies and uses the global QoS for service selection. We evaluate the efficiency of our proposed technique on the WSDream-QoS Dataset [1].

Index Terms—QoS, web service, dependency modeling

1 INTRODUCTION

IN recent years, the Service-Oriented Architecture (SOA) paradigm has gained momentum as a means to develop applications. In SOAs, loosely-coupled software artifacts (commonly referred to as services) may implement specialized functionalities which can be combined with other services from various business partners or public entities into composite services to provide some value-added functionality. Two major entities are involved in any SOA transaction: Service consumers, and Service providers. As the name implies, service providers provide a service on the network with the corresponding service description [2]. A service consumer needs to discover a service to perform its desired task among all the services published by the different providers. The consumer binds to the newly discovered service(s) for execution, where input parameters are sent to the service provider and output is returned to the consumer. In situations where a single service does not suffice, multiple services can form a composite system to deliver the required functionality [3].

In a running composite system, services may exhibit errors, undergo changes, or become unavailable, and may need to be replaced by other services (to achieve better performance or lower cost). The process of composing/replacing service(s) is usually time consuming, error-prone and

often non-optimal. Using automated composition techniques one could improve upon these results. Apart from the functional properties, within SOAs, Quality of Service (QoS) attributes express the non-functional aspects of a service, such as response time, cost, reliability or the supported security protocols etc. The overall performance of a composite system thus depends heavily on how the individual QoS values of the component services effect the composite solution. The orchestration of individual services in a composition also plays a significant role on the QoS values of the composition. Most existing approaches consider this as a local (service level) optimization problem and lack a coherent framework for the specification, and optimization of service compositions focusing on the global (system-wide) QoS properties of the system considering the orchestration and dependency relationships among component services.

Services in SOAs are *autonomous*, i.e., they are independently deployed, and the topology is *dynamic* where a service can leave the system, or fail without notification. In any SOA the various component services may be invoked using a different invocation model. Here, an invocation refers to *triggering a service (by calling the desired function and providing inputs) and receiving the response (return values if any) from the triggered service*. There are six major invocation relations defined in the literature for service compositions [4], [5]: Sequential Invocation, Parallel Invocation, Probabilistic Invocation, Circular Invocation, Synchronous Activation, and Asynchronous Activation. A brief overview of these follows.

Sequential Invocation: In sequential invocation, a service S invokes a service A. It is denoted as Sequential (S : A) (see Fig. 1a). Sequential invocation is also defined as a serial invocation.

Parallel Invocation: In parallel invocation, a service S simultaneously invokes multiple services. For example, if S has service A and service B which are independent and successors of S, S can invoke both A and B at the same time. It is denoted as Parallel (S : A, B) (see Fig. 1b).

- K. Hashmi is with the Department of Computer Science, Wayne State University, Detroit, MI 48202. E-mail: khayyam@wayne.edu.
- Z. Malik is with School of Information Security and Applied Computing, Eastern Michigan University, Ypsilanti, MI 48197. E-mail: zaki@vt.edu.
- A. Erradi is with Department of Computer Science and Engineering (CSE), Qatar University, Al Tarfa, Doha 2713, Qatar. E-mail: erradi@qu.edu.qa.
- A. Rezgui is with Department of Computer Science and Engineering, Mexico Tech University, Socorro, NM 87801. E-mail: rezgui@cs.nmt.edu.

Manuscript received 11 June 2015; revised 19 Apr. 2016; accepted 4 June 2016. Date of publication 7 July 2016; date of current version 7 Dec. 2018.
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.2016.2589244

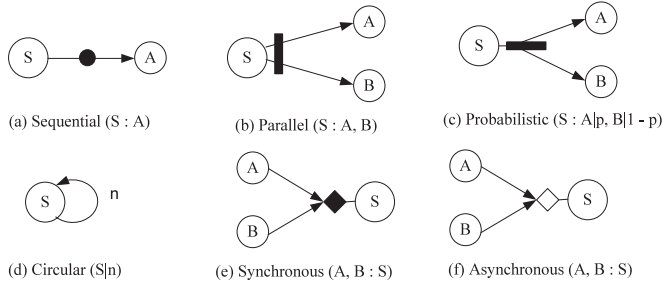


Fig. 1. Major SOA invocation models.

Probabilistic Invocation: In probabilistic invocation, a service S invokes service(s) with a probability (see Fig. 1c). For example, if S invokes service A with the probability p and service B with the probability $1-p$, it is denoted as Probabilistic ($S : A|p, B|1-p$). The probabilistic invocation is also defined as fork invocation.

Circular Invocation: In circular invocation, a service S invokes itself multiple (n) times. It is denoted as Circular ($S|n$). A circular invocation can thus be defined as sequentially invoking itself n times (see Fig. 1d).

Synchronous Invocation: In synchronous invocation, a service S is invoked only after all of its preceding services that were invoked in parallel have completed their execution. For example, if S has synchronous preceding (invoked in parallel) services A and B , both these services would need to complete before S can progress. It is denoted as Synchronous ($A, B : S$) (see Fig. 1e).

Asynchronous Invocation: In asynchronous invocation, a service S is invoked only after one of its preceding services that were invoked in parallel have completed their execution. For example, if S has asynchronous preceding (invoked in parallel) services A and B , either A or B 's completion would cause S to progress. It is denoted as Asynchronous ($A, B; S$) (see Fig. 1f).

It is highly likely that a composite system exhibits more than one type of service invocation. Since certain system properties are a composite function of its component services (e.g., the overall throughput in transactions per second of the system is dependent on the component service that has the least transaction per second-bottleneck), it is important to consider these invocation patterns while considering and modeling the dependencies of the different QoS attributes of the component services. In this paper, we propose a service composition approach that attempts to mitigate the fact that different users may experience different QoS values for a given service, and hence when the service is composed into a system, the observed QoS values could be different from the published QoS values. We use a collaborative filtering (CF) based approach to predict the observed QoS values for a service, then use the composition structure to model the dependency modeling relationship among different QoS values of component service, and then use a multi-objective Markov Decision Model to formulate a globally optimal solution.

The rest of the paper is organized as follows. Section 2 presents our motivating example. Section 3 presents our dependency modeling framework. Section 4 discusses the collaborative filtering approach. Section 5 details our multi-objective service composition approach. In Section 6 we present the results of our experiments and Section 6

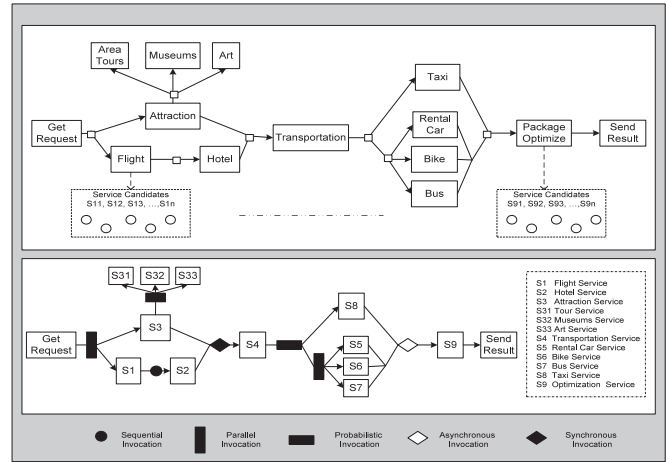


Fig. 2. Scenario with invocation models.

discusses the literature review. Section 8 presents the conclusion and our future directions.

2 SERVICE COMPOSITION OPTIMIZATION STRATEGIES

In this section we look at different service composition optimization goals, and the quantity of information available for the decision making process. It is assumed that all the services under discussion are functionally equivalent and QoS attributes are the primary selection parameters. We start by defining a travel reservation service example to motivate the problem and our proposed approach.

2.1 Sample Scenario

Assume that a user wants to attend a conference and needs to make travel arrangement for his journey. He needs to purchase an airline ticket and reserve a hotel for this travel. Moreover, he needs some transportation to go from the airport to the hotel and from the hotel to other venues "Site-seeing". Assume that the user would be using a SOA-based online service that is a *one-stop shop* providing all the five options (airline ticket, hotel, attractions and transportation) through different component services.

The online reservation system provides many services such as: Attraction service which are outsourced to three services (representing individual services): Art, Museums, and Area tours. This service provides arrangement to visit different areas through sub-contractor companies. For clarity, Fig. 2 shows the different options. User may select Art, Art and Museum, Art and Area tour, Museum or Area tour. The transportation service also works through subcontractors that provide car, bus or bike. These companies provide different services based on the distance between the places user plans to visit (user has the option to choose the mode of transportation). The system also provides services to calculate the distance between two points. The system employs two other services: one for airline ticket, and the other for hotel reservation.

Fig. 2 shows a sample structure of the travel reservation scenario. This structure depicts a combination of different invocation models that may be used to compose the system. Since the user is looking for a travel arrangement that include: booking a ticket, booking a hotel, transportation (car, bike or

bus) based on the distance between the places and visiting some attractive places, these services can be invoked in parallel. Booking a ticket and find attractions is an example of *parallel invocation*. There are three choices that provide attraction (Area tour, Museums and Art). Since user has to choose among these service instances, this is an example of *probabilistic invocation*. Similarly car, bike and bus services can be classified as probabilistic invocations. When the system provides the results of transportation to the reservation service, for generating a discount, this invocation is an example of *asynchronous invocation*. Finally, *synchronous invocation* appears when *package optimization* service waits for other services (hotel, ticket and attraction) to compute the final trip cost.

2.2 Local Optimization

The most basic scenario for service selection deals with the local resources and constraints. In the local optimization approach, the selection of a component Web service that performs a given task of the composite system is done as a stand alone component without considering any system-wide constraints or limitation of other components of the composition. When the composite system wants to perform a certain task (as a part of initial system composition or replacing a faulty service during execution process) the system gathers QoS information for each candidate Web service. After collecting the QoS information, the system constructs a QoS vector which is then evaluated to calculate the usefulness of each service. The candidate services are then ranked based on their usefulness or utility and the best available service is selected. This selection process could include user defined weights and constraints such as execution time, availability etc.

In our running scenario from Fig. 2 assume that the system is considering candidates S11, S12 and S13 for the flight service. Assume that the QoS vector used to make the decision consists of four QoS attributes $\langle \text{Throughput}, \text{Availability}, \text{Reliability}, \text{Execution Time} \rangle$. For the sake of simplicity, we assume that the service prices are static and all the above QoS attributes are equally important for the system. Assuming that the QoS vector for S11 is $\langle 91, 95, 95, 88 \rangle$, for S12 it is $\langle 94, 92, 92, 88 \rangle$ and S13's vector is $\langle 90, 98, 98, 90 \rangle$. For local optimization the system will select service S13. Although S11 has a lower value for execution time but S13 has much higher values for Availability and Reliability. The system is able to make this decision since it has all the local information regarding the candidate services to compare and make an informed decision. The selected service was the best possible choice in the local context, i.e., local maxima.

2.3 Global Optimization

Since local optimization is performed on a per service basis, it is possible that it may not be the optimal choice in the context of the composite system. In Fig. 2 let us assume that hotel service has a QoS vector of S2 $\langle 93, 95, 95, 92 \rangle$. One of the system wide constraints could be that both flight and hotel service should finish their execution in 180 ms. In order to meet this constraint, the system cannot choose service S13 $\langle 90, 98, 98, 90 \rangle$ since it does not satisfy the system constraint (i.e., $92 + 90 = 182 > 180$). Hence, it would have to select S11 $\langle 91, 95, 95, 88 \rangle$ which is clearly not the best local choice but turns out to be the best global choice for the

composite system. Global optimization is thus heavily dependent on the amount of system wide information and constraints for optimal decision making. Similarly, if we add another dimension of information, e.g., that we are looking to optimize the throughput (after meeting the execution time constraint) then we can see that service S12 $\langle 94, 92, 92, 88 \rangle$ will be a better choice since it has a higher throughput value; although it was the second best choice in the previous scenario, and the last choice in local optimization.

Let us look at another case of global optimization in our running example by adding structural information to our current discussion of service selection. Assume that we have two options for Attraction service: S3a $\langle 98, 95, 95, 180 \rangle$ and S3b $\langle 93, 96, 96, 180 \rangle$. Here service S3a is a better choice for Attraction service since it has far better throughput and comparable values of other QoS components. Now if we look at the structure of our system in Fig. 2 we see that Attraction service is invoked in parallel with Flight and Hotel services which in turn, are sequential in nature, i.e., Parallel (Get Request: S3, Sequential (S1:S2)). We can see that based on our current selection of S2 $\langle 93, 95, 95, 92 \rangle$ for hotel and S12 $\langle 94, 92, 92, 88 \rangle$ for flight the maximum throughput for the current composite solution with the Sequential invocation (S12:S2) is 93 (min (93,94)). This, in turn implies that the maximum throughput of the current composite solution with the Parallel invocation (Get Request: S3, Sequential (S12:S2)) cannot exceed 93, i.e., (min (98,93)). Based on this information our choice of S3a $\langle 98, 95, 95, 180 \rangle$ for hotel service is not a global optimal service for the composite system as S3b $\langle 93, 96, 96, 180 \rangle$ offers better QoS values for availability and reliability while still matching the maximum throughput value of 93 for the parallel invocation of the scenario. Hence, the service selection framework could benefit from the invocation pattern information for making optimal decisions to improve the overall QoS of the composition. To differentiate, we refer to such relationships among the same QoS components of different service of the composition as *dependency modeling* in the rest of the paper.

3 DEPENDENCY MODELING FRAMEWORK

In order to demonstrate our proposed approach for QoS dependency modeling we will be using the following quality components (note that other QoS attributes can be modeled in a similar manner).

- *Cost*: The cost $q_{co}(s, f_i)$ of an operation f_i of a service s is cost incurred (fee paid to the service provider) to invoke an operation in order to perform a particular task. These values are either published by the service providers or are available on demand. The cost is measured as a per invocation unit.
- *Execution Time*: Given an operation f_i of service s , the execution time $q_{et}(s, f_i)$ measures the difference between the time when the request is sent and the time when the results are received. The total execution time is the sum of the processing time $T_{process}(s, f_i)$ and the transmission time $T_{trans}(s, f_i)$. It is measured in milliseconds.
- *Reliability*: The reliability $q_{re}(s)$ of a service s is measured as the probability of a system to accurately

TABLE 1
QoS values for Different Invocation Models

	Sequential	Probabilistic	Circular	Parallel Synchronous	Parallel Asynchronous
Cost	$\sum_{i=1}^n q_{co}(f_i)$	$\sum_{i=1}^n p_i * q_{co}(f_i)$	$q_{ac}(f) * c$	$\sum_{i=1}^n q_{co}(f_i)$	$\sum_{i=1}^n q_{co}(f_i)$
Execution Time	$q_{et}(s, f_i) = T_{process}(s, f_i) + T_{trans}(s, f_i)$	$\sum_{i=1}^n q_{et}(s, f_i)$	$q_{ac}(s, f_i) * c$	$\max(q_{et}(s_1, f_i), \dots, q_{et}(s_n, f_i))$	$\max(q_{et}(s_1, f_i), q_{et}(s_2, f_i), \dots, q_{et}(s_n, f_i))$
Reliability	$\prod_{i=1}^n q_{re}(f_i)$	$\sum_{i=1}^n p_i * q_{re}(f_i)$	$q_{re}(f)^c$	$\prod_{i=1}^n q_{re}(f_i)$	$\max(q_{re}(f_1), q_{re}(f_2), \dots, q_{re}(f_n))$
Availability	$\prod_{i=1}^n q_{av}(f_i)$	$\sum_{i=1}^n p_i * q_{av}(f_i)$	$q_{av}(f)^c$	$\prod_{i=1}^n q_{av}(f_i)$	$\max(q_{av}(f_1), q_{av}(f_2), \dots, q_{av}(f_n))$
Throughput	$\min(q_{th}(s_1), q_{th}(s_2), \dots, q_{th}(s_n))$	$\sum_{i=1}^n p_i * q_{th}(s, f_i)$	$q_{th}(s)$	$\min(q_{th}(s_1), q_{th}(s_2), \dots, q_{th}(s_n))$	$\max(q_{th}(s_1), q_{th}(s_2), \dots, q_{th}(s_n))$

respond to a request (i.e., the operation is completed and a message indicating that the execution has been successfully completed is received by a service requestor) within the maximum expected time frame indicated in the Web service description. The reliability (or success rate) is dependent on the hardware and/or software configuration of Web services and the network connections between the service requesters and providers. The value of the reliability can be computed from data of past invocations using the expression $q_{re}(s) = (K - N_f(s))/K$, where $N_f(s)$ is the number of faulty executions of the service s , and K is the number of times service s has been invoked.

- **Availability:** The availability $q_{av}(s)$ of a service s is the probability that the service is accessible. The value of the availability of a service s is computed using the expression $q_{av}(s) = T_a(s)/\theta$, where T_a is the total amount of time for which service s is available over and observed θ amount of time. Availability is usually measured in the percentage. The higher the value is, the more a system is available.
- **Throughput:** The throughput $q_{th}(s)$ of a service s is measured as the number of completed requests per unit amount of time.

Table 1 shows how to calculate different QoS values for major SOA invocation models as discussed above.

For the composite system we would need to calculate the over all QoS values suitable for negotiation for the new requested service. This is mainly dominated by the structure of the composition. Apart from the simple invocation methods mentioned above we also need to look into complex invocation patterns. One of the key elements of a composite system is a complex loop. A complex loop can be defined as a *Circular Invocation* with a linear or non-linear complex execution path. Loops may contain different combinations of *Invocation Patterns*, i.e., nested loops, probabilistic invocations etc. Hence, the QoS values of a complex loop structure could be calculated by the probability of number of iterations of the loop, i.e., p_e , the probability of the exiting the loop p_e and QoS values of the execution path of the loop. We can always flatten a loop structure into repeatable blocks of linearly executing patterns. A loop may have multiple entry and exit points. We can assume that the actual entry and exit points of a loop structure could be ignored for the ease of calculation as they have minimal affect on the actual QoS values. However, the probability of exiting out of the loop p_e along with the individual probabilities of the different execution patterns within the loop are required for the QoS calculations.

Let us assume that the transition probability of execution of Service s_{i+1} after executing Service s_i is denoted by p_i

(in a linear execution path this will always be equal to 1). The probability of exiting a loop after executing Service s_i is denoted by $p_{e_{i,j}}$ (in a linear execution path this will be 0) where $p_i + \sum_{j=1}^m p_{e_{i,j}} = 1$ and $i \in [1, n]$. The cost of the loop is depicted by c_i and the total execution time is shown as t_i . After x executions of the loop where $x \in [0, +\infty]$ the loop can be transformed into x linear executions paths for the component services and we already know that the probability of executing the next service is $p_{e_{i,j}}$ where $j \in [1, m]$. So the probability for the combination of service for the execution for x where $x \in [0, +\infty]$ time and the execution of any service after Service S_1 is $(\prod_{i=1}^n p_i)^x p_{e_{1,j}}$ ($x \in [0, +\infty]$, $j \in [1, m]$). Similarly the probability for executing Service S_{k+1} after executing Service S_k where $k \in [2, n]$ is $(\prod_{i=1}^n p_i)^x (\prod_{i=1}^{k-1} p_i) p_{e_{1,j}}$ ($x \in [0, +\infty]$, $j \in [1, m]$). Let $p_0 = 1$ then the probability of executing a service after Service S_k is

$$p'_{e_{k,j}} = \sum_{x=0}^{+\infty} \left(\prod_{i=1}^n p_i \right)^x \left(\prod_{i=0}^{k-1} p_i \right) p_{e_{k,j}} \quad (1)$$

$$= \frac{(\prod_{i=1}^{k-1} p_i) p_{e_{k,j}}}{1 - \prod_{i=1}^n p_i} \text{ where } (k \in [1, n], j \in [1, m]). \quad (2)$$

Now the probability that service $S_1 \dots S_n$ is executed x times and then the loop is terminated at service S_k is $(\prod_{i=1}^n p_i)^x (\prod_{i=0}^{k-1} p_i) (1 - p_k)$. The cost for this execution accumulated by the services will be $x \sum_{i=1}^n c_i + \sum_{i=1}^k c_i$ and the execution time will be $x \sum_{i=1}^n t_i + \sum_{i=1}^k t_i$. The reliability for this execution will be $(\prod_{i=1}^n q_{re}(i))^x (\prod_{i=1}^k q_{re}(i))$. Hence we get the following equations for Cost (q_{co}), Total Execution Time (q_{et}), Reliability (q_{re}) and Availability (q_{av}) respectively

$$q_{co} = \sum_{k=1}^n \sum_{x=0}^{+\infty} \left(\left(\prod_{i=1}^n p_i \right)^x \left(\prod_{i=0}^{k-1} p_i \right) (1 - p_k) \left(x \sum_{i=1}^n c_i + \sum_{i=1}^k c_i \right) \right) \quad (3)$$

$$q_{et} = \sum_{k=1}^n \frac{(\prod_{i=0}^{k-1} p_i) (1 - p_k) (\sum_{i=1}^k t_i + \prod_{i=1}^n p_i \sum_{i=k+1}^n t_i)}{(1 - \prod_{i=1}^n p_i)^2} \quad (4)$$

$$q_{re} = \sum_{k=1}^n \sum_{x=0}^{+\infty} \left(\left(\prod_{i=1}^n p_i \right)^x \left(\prod_{i=0}^{k-1} p_i \right) (1 - p_k) \left(\prod_{i=1}^n r_i \right)^x \prod_{i=1}^k r_i \right) \quad (5)$$

$$q_{av} = \sum_{k=1}^n \frac{(\prod_{i=0}^{k-1} p_i) (1 - p_k) \prod_{i=1}^k r_i}{1 - \prod_{i=1}^n (p_i * r_i)} \quad (6)$$

$$q_{th} = \sum_{k=1}^n \frac{(\prod_{i=0}^{k-1} p_i)(1 - p_k) \prod_{i=1}^k r_i}{1 - \prod_{i=1}^n (p_i * r_i)}. \quad (7)$$

After calculating the individual QoS values for each invocation pattern we now need to calculate the max possible QoS value for each component of the QoS vector in order to establish a target QoS value that would be used for global optimization. Since our QoS calculation follow a single entry and single exit format we can use the refined process tree approach [6] to find the different execution paths and apply our previously described invocation patterns. In this regard, Algorithm 1 calculates the individual QoS values for all components of the QoS vector at a given node. In this algorithm we use depth first on all the nodes and create tree for all the path to the leaf nodes. Once we get to the leaf node we recursively backtrack and calculate the QoS value for that path. If we run into a conditional branch we add both the paths as child node and this allows us to still find all the possible execution paths for the composition. We need to focus on the QoS value for the paths that contain the service that we want to replace. This will ensure that for global optimization we are only considering the relevant paths. Initially, we will set the QoS values for service to be replaced at the maximum possible values to calculate the target QoS negotiation values for the service (refer to Section 2). We can see that the proposed depth first algorithm has low time and space complexity. With N nodes and E edges the first exploration of the graph takes $O(|V| + |E|)$ and in the second exploration for the graph we only visit all nodes once, (since the QoS calculation at every invocation point takes a linear amount of time) our execution takes $O(|V| + |E|)$ time.

Algorithm 1. QoS Calculation for a Service Execution

```

1: DepthFirstQoS (Node N, Tree F)
2: T = create a new tree;
3: Add T as child of F;
4: if N has no children then
5:   Get a path to parent node of N
6:   for All component of QoS vector do
7:     Apply the QoS pattern and multiply it with the
       probability to get the QoS value at current
       invocation point.
8:   end for
9: end if
10: for each child node of N do
11:   if isVisited = true then
12:     Add all the child nodes to T
13:   if current node = leaf node then
14:     Get a path to parent node of T
15:     for All component of QoS vector do
16:       Apply the QoS pattern and multiply it with
       the probability to get the QoS value at current
       invocation point.
17:     end for
18:   else
19:     Call DepthFirstQoS(current node, T)
20:   end if
21: end if
22: end for

```

Once we determine the dependency relationship we need to be able to predict the QoS values that we may be

able to get from our composition. That could be done by predicting the observed QoS values of individual services and then using them in the dependency calculation. We present our collaborative filtering approach in the next section that is used for this purpose.

4 COLLABORATIVE FILTERING

In ideal situation Web services should meet all the published performance criteria for all of its users. However it is rarely the case. We use a collaborative filtering based approach to mitigate the discrepancies between published and observed QoS values for the Web services. We can loosely classify the Web service attributes into two broad classes of certain and uncertain attributes. The value of certain attributes remains fixed or unvarying for all the users, i.e., service name, service provider, cost (published/mutually agreed upon). The value of uncertain attributes on the other hand may change or vary from user to user, or even for the same user, from one invocation to the other invocation, mostly based on the environmental factors. These attributes have a marked tendency to deviate from their published values. Response time is one such example of an uncertain attribute. Many factors can influence the response time of a service, e.g., internet speed, network congestion, slow hardware. These factors are usually out of the hands of the service provider. These uncertain attributes may influence the perceived overall QoS value of the systems. However, most of the service selection algorithms do not take the uncertainty of these attributes into account which leads to inaccurate/sub-optimal service matching, these system only consider the published QoS values of the services. However, as discussed above the observed QoS value may differ from the published QoS values for a user. It is well accepted that it is impractical (and nearly impossible) for a service consumer to invoke all the services (under consideration) to record their observed QoS values and then use them for service composition. We can instead utilize the experience of other (similar) users to predict the observed QoS values for services and use them as a metric for service composition. This requires us to investigate an approach that first of all would be able to search users that are similar (in terms of environment) to the current service consumer and have invoked the service under consideration with reported observed QoS values for the service. Second, use the experience of similar users to predict observed QoS values for the service consumer. These predicted QoS values could later be used in the negotiation process to minimize the affect of uncertainty of QoS values of the services. Note that it is common for the end user to receive lower values of QoS components than the promised/published QoS values (observed response time of 98 ms as compared the published response time of 60 ms) and very rare if not highly improbable to experience better than published QoS values. Hence, it is safe to assume that using observed QoS values will not over compensate a service provider for its published QoS values which in turn, could result in over promise of QoS to the consumer.

We use collaborative filtering on the observed QoS values to solve the problem of uncertainty. The two most commonly used memory based collaborative filtering

approaches are Pearson correlation and Vector cosine based similarity. Pearson coefficient [7] is symmetric, invariant to scale and location of variable and can also detect negative correlation making it suitable for our problem at hand. The similarity between two users i and j can be calculated using Pearson Correlation Coefficient as:

$$sim(i, j) = \frac{\sum_k (v_{i,k} - \bar{v}_i)(v_{j,k} - \bar{v}_j)}{\sqrt{\sum_k (v_{i,k} - \bar{v}_i)^2} \sqrt{\sum_k (v_{j,k} - \bar{v}_j)^2}}, \quad (8)$$

where $v_{i,k}$ denotes the QoS value that user i received from service k , \bar{v}_i denotes the average QoS that user i received from all the services invoked, and the summation is over all the services that have been invoked by both i and j .

In the real life scenarios we observe that it is rather rare to have a consumer that has a larger number of similar consumers or a service that has large number of similar services. Hence, we will be working with a very limited data set. The prediction confidence could be increased if we use both content and user based similarity to predict the observed QoS values. In our approach when we have a missing QoS value for a service and we do not have any similar users that have invoked that particular service, we find out similar services and use their QoS values to predict the missing QoS value, and vice versa. We assume that $S_i \neq \emptyset$ and $U_p \neq \emptyset$. We use f_u and f_i to assign weightage to both the filtering values

$$f_u = \sum_{i \in S_i} \frac{Sim(i, j)}{\sum_{i \in S_i} Sim(i, j)} \times Sim(i, j). \quad (9)$$

We use $\lambda(0 \leq \lambda \leq 1)$ to adjust the influence of both user (u) based and content based (i) filtering

$$W_u = \frac{f_u \times \lambda}{(f_u \times \lambda) + (f_i \times (1 - \lambda))} \quad \text{and} \quad (10)$$

$$W_i = \frac{f_i \times \lambda}{(f_i \times \lambda) + (f_u \times (1 - \lambda))},$$

where $W_u + W_i = 1$, hence the predicted QoS value would be

$$Q_{val} = W_u \times \left(\bar{v}_i + \frac{\sum_{j \in S_i} Sim(i, j)(v_{j,p} - \bar{v}_j)}{\sum_{j \in S_i} Sim(i, j)} \right) \quad (11)$$

$$+ W_i \times \left(\bar{v}_q + \frac{\sum_{q \in U_p} Sim(p, q)(v_{p,q} - \bar{v}_q)}{\sum_{q \in U_p} Sim(p, q)} \right).$$

Using the above mentioned approach implies that the initial set of services that are being considered for selection process have a very high chance of forming a service agreement. The observed and predicted QoS values help eliminate the uncertainty in perceived QoS values. The next step is to formulate a globally optimal composition using the dependency modeling and the services with better predicted observed values.

5 MULTI-CONSTRAINT MODELING

The initial filtering process provides us with a good platform of a subset of services to be considered for service

composition. For service composition we use the concept of Markov Decision Process (MDP). MDP is an Artificially Intelligent model for making decisions in environments where there is a higher percentage of uncertain outcomes. MDP has been successfully applied in different domains to solve decision making problems [8], [9], [10], [11]. We model each dependency relationship as a single constraint and use this model to solve the global optimization dependency problem for QoS parameters of component Web services.

In our system each Web service has a unique identifier ID and a QoS vector. Each QoS vector is a combination of five QoS values, i.e., $QoS = \langle Cost, ExecutionTime, Reliability, Availability, Throughput \rangle$ (it can easily be generalized to more QoS component values since $QoS = \langle QoS_1, QoS_2, \dots, QoS_n \rangle$). Markov Decision Process is defined as $\langle S, A, P, R, \gamma \rangle$. S represents the set of states of the system; A stands for the set of actions in the system and $A(s)$ is a subset of action available at state s , i.e., $A(s) \in A$ where $s \in S$; P is the probability of an action such that $P_a(s, s') = Probability(s_{t+1} = s' \mid s_t = s, a_t = a)$; R is the expected or immediate reward for the current transition and γ is used to factor in the importance of current reward and the future rewards. We can extend this basic single constraint MDP to a multi-constraint MDP by enhancing the reward function to consider the multiple constraints. We modify the reward function to compensate for multiple constraints. Let us assume that the system has c number of constraints. The reward function will be as follows:

$$R_a(s, s') = [R1_a(s, s'), R2_a(s, s'), \dots, Rc_a(s, s')]^T. \quad (12)$$

In addition to the above constraints we need to consider the fact that every composite system has an entry and exit point. We can argue that this entry and exit point may contain more than one service. Some systems can be invoked using multiple services and similarly their execution may take multiple routes to reach different end states. Hence, the updated model for Web service will have seven tuples $MDP = \langle S, S_s, S_e, A, P, R, \gamma \rangle$. Where S_s is a set of starting states for the composition and S_e is the set of end states of the composition where $S_s \in S$ and $S_e \in S$. Here an action corresponds to the execution of a Web service and the reward vector contains one reward for every QoS component of the Web service. For Web service ws and our QoS attributes under consideration the reward vector will be

$$ws(q)(s, s') = [ws(q_{co})(s, s'), ws(q_{et})(s, s'), ws(q_{re})(s, s'), \quad (13)$$

$$ws(q_{av})(s, s'), ws(q_{th})(s, s')]^T.$$

The solution of MDP is a decision constraint and a constraint is the procedure for selecting Web services. These constraints, represented by ζ are basically just mappings from states to actions, defined as $\zeta: S \rightarrow A$. Each constraint can represent a single composite solution of Web service. Hence, our system searches for a set of Pareto optimal constraints which optimize QoS attributes of the composite system. The set ζ^o of Pareto optimal constraints is defined by

$$\zeta = \left\{ \zeta^O \in \prod \mid \nexists \zeta \in \prod, s.t. V^{\zeta^O}(s) >_o V^\zeta(s), \forall s \in S \right\}, \quad (14)$$

where \prod defines the set of all constraints and dominance relation is represented by $>_o$. For $a = (a_1, a_2, \dots, a_n)$ and $b = (b_1, b_2, \dots, b_n)$, $a >_o b$ means that $a_i \geq b_i$ is satisfied for all i and $a_i > b_i$ for at least one i . Moreover $V^\zeta(s) = (V_1^\zeta(s), V_2^\zeta(s), \dots, V_m^\zeta(s))$ is the value of the vector s as per the constraint ζ , i.e.:

$$V^\zeta(s) = E_\zeta \left\{ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid s_t = s \right\}, \quad (15)$$

where E_ζ represents the expected value when the service follows constraint ζ , s_t at time t with the reward vector r_t .

Q-learning is calculated as:

$$Q_\zeta(s, a) \models E_\zeta \left\{ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid s_t = s, a_t = a \right\}. \quad (16)$$

5.1 Single Constraint Multi-Objective Service Composition

In our approach we use Q-learning for QoS objectives. We spawn a Q-learning service for every QoS objective. Based on the multiple objective criteria the importance (weight) of every QoS objective for a Web service is learned rather than using predefined weights. Every service i selects a Web service ws_i at each state S in such a way, that optimizes the relative QoS objective of the Web service. Once this process is completed then the Web service assign a weight $W_i(s)$ to their selected service and later negotiate among themselves to select the most suitable candidate to be executed at each state. The service having the maximum weight will be able to execute its option at each state. The objective is:

$$W_k(s) = \text{Max}_{i \in 1, \dots, n} W_i(s). \quad (17)$$

Therefore, in this case the Web service k is called the leader service and is allowed to invoke Web service $W_k(s)$. In the next round the Web services evaluate the results of the previous selections and adjust their $w_i(s)$ values based on positive or negative outcome of the previous round. Hence, we may have a different leader service in the next round.

Algorithm 2. Single Constraint Algorithm

```

1: Start at state s
2: randomly select a service k to be the leader
3:  $W_k \leftarrow 0$ 
4:  $a_k \leftarrow \text{argmax}_a Q_k(s, a)$ 
5: repeat
6:   for all services i except k do
7:      $W_i \leftarrow \text{max}_a Q_i(s, a) - Q_i(s, a_k)$ 
8:     if the highest  $W_i > W_k$  then
9:        $W_k \leftarrow W_i$ 
10:     $a_k \leftarrow \text{argmax}_a Q_i(s, a)$ 
11:     $k \leftarrow i$ 
12:   end if
13: end for
14: until converges
15: return b
```

W represents the difference in the predicted versus actual reward received by the Web service. Web services predict to receive a reward value P if their selected service was executed at s . If their service was not executed then instead of receiving the predicted reward P it receives the actual reward A . So $W = P - A$. In the case when a service's suggested selection was executed $P = A$ and if not then the service will receive a negative reward that is equal to $(P - A)$. So if service k ends up being the leader in a certain round then all service except k will update their W values using the following:

$$W_i(x) \rightarrow (Q_i(x, a_i) - (r_i + \gamma \max_{b \in a} Q_i(y, b))). \quad (18)$$

At this point the next state " s " and the reward vector r_i is influenced by the leader service rather than being a decision of each individual service. We present this in Algorithm 2.

5.2 Multiple Constraint Multi-Objective Service Composition

In our second approach we introduce the concept of convex hull to solve the multiple constraint problem for service selection for composite solutions. The convex hull is defined as the smallest convex set that contains all the points that lie on the boundary of this convex set. We combine this concept, which is similar to the Pareto front, into our Q-learning based approach based on the fact that both concepts are essentially maxima over different trade-off factors in the linear domains.

We use a value iteration based method to obtain a set of service selection constraints that is Pareto optimal:

$$V(s, a) = V(s, a)(1 - \alpha) + \alpha \left[\gamma \mathcal{R} \bigcup_{a'} V(s', a') + r(s, a) \right]. \quad (19)$$

$V(s, a)$ represents the set of vertices of the \mathcal{R} when action a is taken at state s for all possible Q-value vectors, r is the reward vector, discount value for the process is represented by γ , rate of learning is controlled by α and the operator \mathcal{R} represents the set of extracted vertices of the \mathcal{R} .

Algorithm 3. Multiple Constraint Algorithm

```

1:  $V(a, s) = \text{Randomly initialize } \forall s, a$ 
2: while convergence condition not met do
3:   for all  $s \in S, a \in A$  do
4:      $V(s, a) = V(s, a)(1 - \alpha) + \alpha [\gamma \mathcal{R} \bigcup_{a'} V(s', a') + r(s, a)]$ 
5:   end for
6: end while
```

Our solution follows the greedy exploration methodology and the dominance relationship between Q-vectors is used for selecting a particular action. In this approach we do not backtrack based on the maximal expected reward for each vector rather we use the set of maximal expected rewards for the given set of constraints as the basis for backtracking. This is illustrated in Algorithm 3.

6 STUDY AND RESULTS

To study the efficiency of our approach we implemented our motivational scenario. We compared the single constraint and multi-constraint approaches along with our invocation

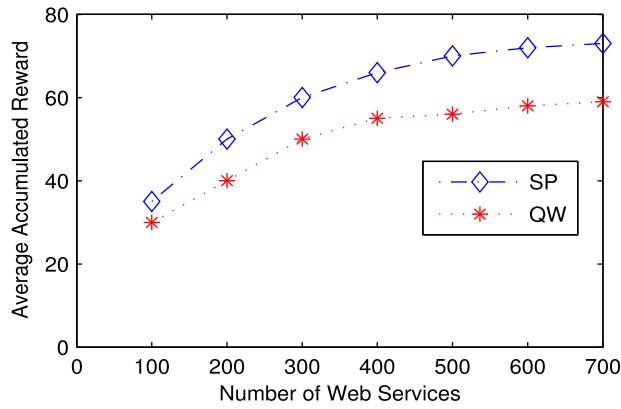


Fig. 3. Single constraint algorithm.

patterns and collaborative filtering based solution. Moreover, we compared our approach with similar approaches in the literature. The experiment environment consists of a Windows server 2008 (SP2)-based Quad core machine with 8.0 GB of ram. We used the WSDream-QoS dataset [1], which contains multiple Web services distributed in computer nodes located all over the world (i.e., distributed in 22 different countries). Planet-Lab is employed for monitoring the Web services. We take the published services and their corresponding QoS values and assign cost values to them.

6.1 Experiment 1: Single Constraint Algorithm

In the first experiment, we compared the single constraint algorithm with the Q-learning approach [12]. We compare the effectiveness of our approach in formulating a composite solution with no pre-defined user weight preferences against the user defined weight vector for the Q-learning approach. We use the accumulated reward for a composition to compare the quality of the discovered solution. The Q-learning approach uses a weight vector of $\omega = (q_{co} = 0.2, q_{et} = 0.2, q_{re} = 0.2, q_{av} = 0.2, q_{th} = 0.2)$ for this experiment, i.e., equal weightage for all the QoS components. Fig. 3 shows the results of our experiment of average results of 30 runs of both the algorithm with varying number of Web services. We can see that the single constraint approach out performs the Q-learning approach regardless of the number of Web services. The difference in quality of solutions increases when the number of Web services increase. This is attributed to the fact that our approach scales better at exploring the Pareto front than the Q-learning approach. Second, pre-defined weights guide the search process for the Q-learning based approach with may not be the best case solution since learning the weights on the fly could find solutions in other wise unexplored regions.

6.2 Experiment 2: Multi-Constraint Algorithm

In this set of experiments, we demonstrate that our algorithm finds Pareto optimal solutions considering the dependency relationships among the different QoS components. We assign multiple concrete services (50 and 100) respectively to the abstract services and observe the proposed solutions. In the first experiment, we assigned 50 concrete Web services to every abstract component Web service. We use three QoS attributes, i.e., Cost, Availability and Response time. The objective is to minimize the cost and response

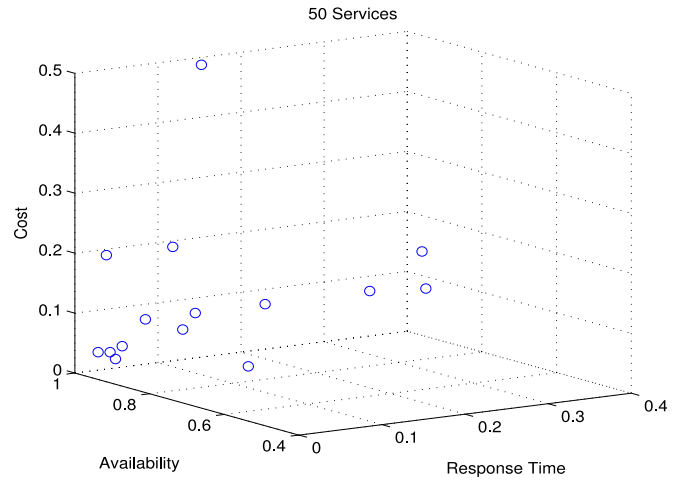


Fig. 4. Results of composition with 50 services.

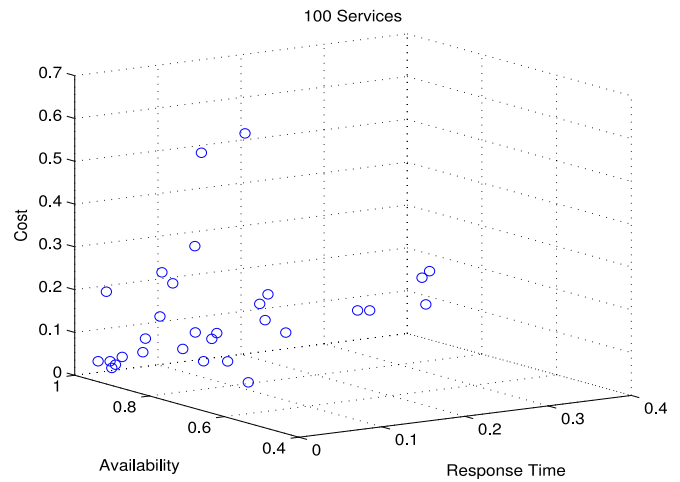


Fig. 5. Results of composition with 100 services.

time, while maximizing the availability of the system. Fig. 4 shows the results of the pareto-optimal solutions found. We can see that the proposed algorithm is able to find high quality solutions. Fig. 5 shows the results of our next experiment, where we increased the number of concrete services from 50 to 100.

In the next experiment, we show that our solution converges very effectively towards the solution with varying number of services. We tested our approach with four abstract Web services and varied the number of concrete services to 100, 200, 300 and 400 Web services. As we can see from Fig. 6, our solution takes longer time for finding optimal solutions when the number of concrete Web services are increased. With 100 concrete Web services, our solution converges at around 600 episodes mark, while when we increase the number of concrete Web service to 200 it take it about 800 episodes to converge. Similarly, it takes 1,200 and 1,400 episodes to find solutions for 300 and 400 concrete Web services respectively. This shows that our approach can handle a large number of Web services and still performs efficiently.

6.3 Experiment 3: Individual QoS Value Comparisons

In this set of experiments, we show the difference in performance of solution when we use just single constraint MDP

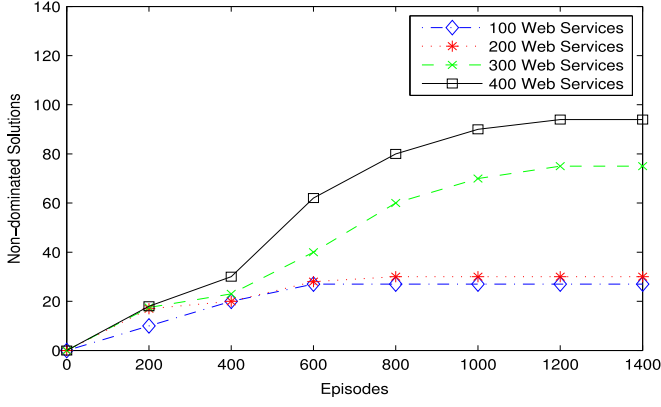


Fig. 6. Multiple constraint algorithm.

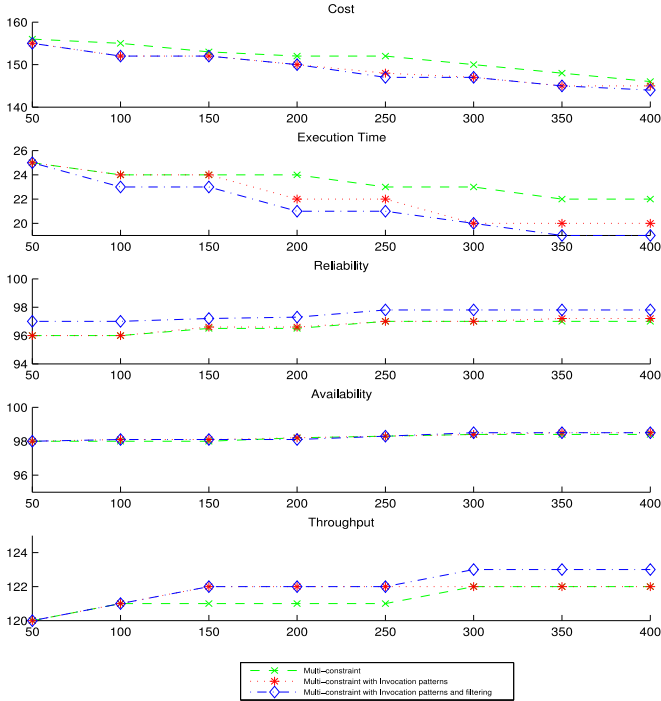


Fig. 7. QoS value comparison.

versus multiple constraint MDP and adding the information from dependency modeling and collaborative filtering. Fig. 7 shows the individual QoS value comparisons for our proposed approach. We can see the performance of multiple constraints, multiple constraints with invocation patterns and multiple constraints with invocation patterns with filtering on the different QoS values of Cost, Execution Time, Reliability, Availability and Throughput when tested on a pool of 50 to 400 service. We can see that Multi-policy with invocation patterns with filtering consistently yields better results for all the QoS values. We can see that introducing invocation patterns into multiple constraints algorithm yields significant improvement for QoS value of Cost, Reliability and Availability. These QoS values rely heavily on the orchestration pattern of service within a composition. Reliability benefits most from the introduction of collaborative filtering as compared to other methods. Availability was the least improved QoS value among our experiments as the services did not show any significant variation in their availability values. However, it is evident that adding

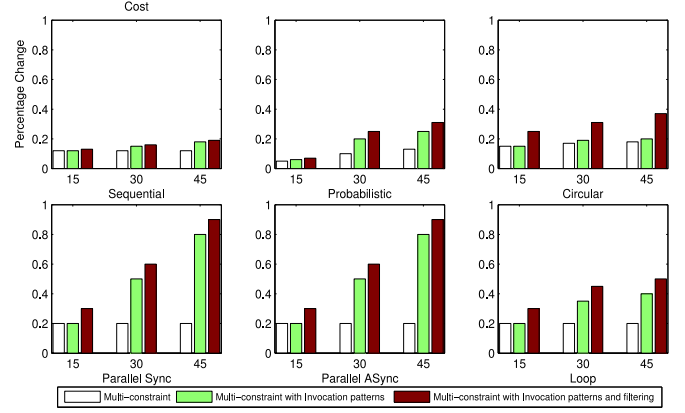


Fig. 8. Utility comparison of cost.

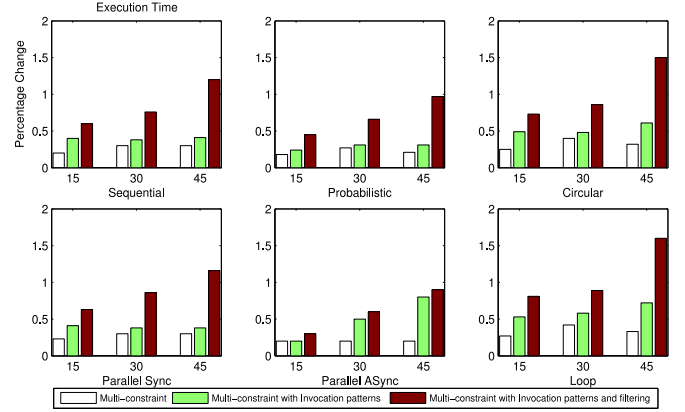


Fig. 9. Utility comparison of execution time.

the invocation information and predicting observed QoS values improves the overall utility of the system.

6.4 Experiment 4: Individual QoS Value Comparisons for Invocation Patterns

In this set of experiments we compare the impact of using our approaches for different QoS invocation patterns. We implemented three different SOA systems that use 15, 30 and 45 services respectively with different invocation patterns. The number of distinct execution paths for each of these systems is 16, 51, and 77 respectively. We replace one service from the composite system at different invocation points to measure the impact on the total utility of the system. Figs. 8, 9, 10, 11, and 12 show the utility gain for all three solutions with varying number of component services. We can see that for every invocation pattern MclpFi performs better when the number of component service increase. This gain is primarily related to better prediction since the more data points are available to calculate the similarity of services. Cost, Execution Time and Throughput are the components that benefit more from the proposed solution. Reliability and Availability also show slight gains in the utility value.

6.5 Experiment 5: Comparison with Existing Approaches

In this experiment we compare our solution(MclpFi) with similar approaches presented in the literature. We base our comparison on the utility values of these techniques and the

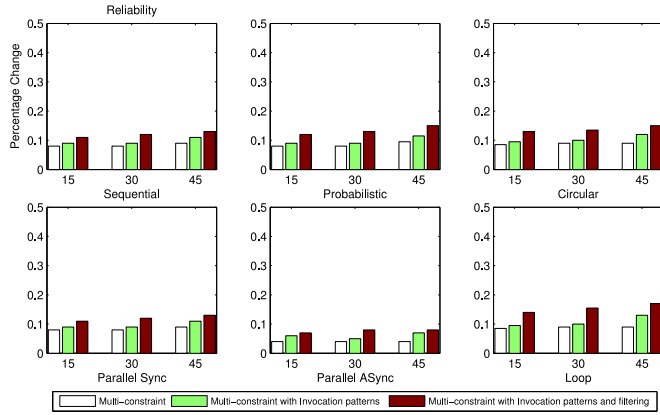


Fig. 10. Utility comparison of reliability.

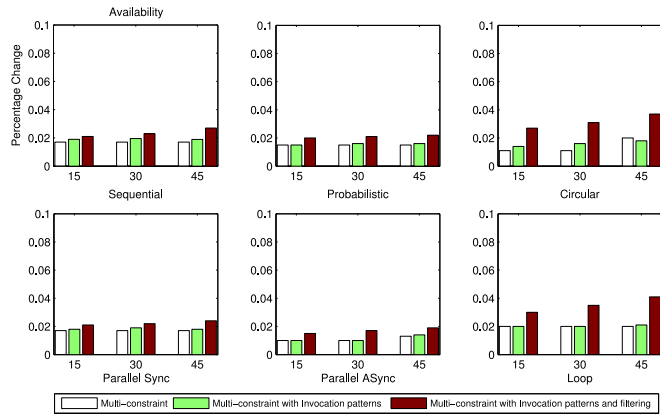


Fig. 11. Utility comparison of availability.

time it takes to reach a solution. SBA [13] uses a GA based approach with an offer and counter-offer based protocol for searching a mutually agreeable solution. The degree of overlap among the QoS values requested by the consumer and those offered by the provider are taken into account in SBA. We use the results for the maximum overlap (80 percent) in our comparisons. Similarly, NBA [14] uses a GA based approach with a very similar fitness function as used in our technique, but does not take into consideration any other parameters. SWC [15] also uses a GA based approach for the semantic composition of Web services. It uses the semantic equivalence in addition to the QoS values to determine the best offering for the composition. BLGAN [16] uses a Bayesian learning based approach with GA and incomplete information model to learn the reserve price of its opponent. GTFSN [17] presents a game theoretical model of signaling games for Service Level Agreement negotiation. The results are presented in Fig. 13. We can see that our solution is the quickest in improving on the initial solution. This can be attributed to the use of invocation patterns, and the solution improves exponentially as the values stabilize (high jumps around time 27 and 49). We can also see that our solution finds a solution within the 97 percent utility range in about 66 ms, while SWC (the second best) takes about three times more time. Other techniques fail to generate such a solution and NBA and SBA plateau around the 92 percent range while BLGAN and GTFSN only reach a solution of around 94 percent utility. The results suggest that our approach outperforms similar methods both in

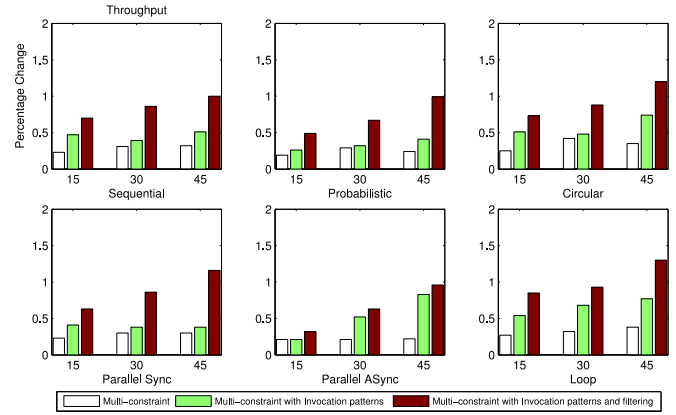


Fig. 12. Utility comparison of throughput.

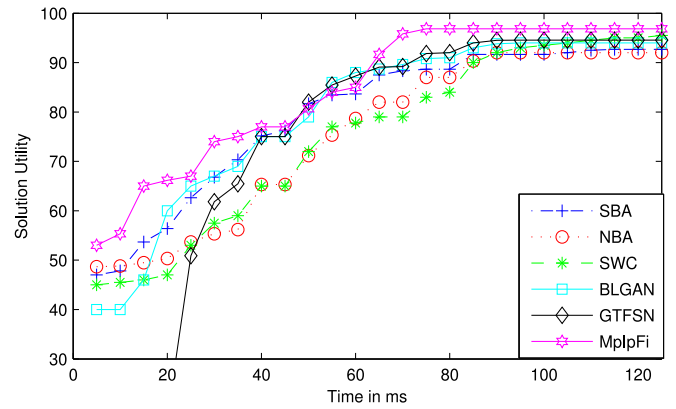


Fig. 13. Performance comparison of proposed approach.

terms of finding the optimal solution and the amount of time it takes to find that solution.

7 LITERATURE REVIEW

QoS based service selection has been an active research topic for service composition. Most of the research in this area has focused on local optimization of QoS attributes that may not result in a globally optimal solution [18], [19], [20]. Some early efforts of global optimization using integer programming, heuristic based searches and critical path have been presented in [21], [22], [23], [24].

A mixed integer programming based global optimization approach has been presented in [25]. In this approach, the authors decompose the global constraints into local constraints by mapping them to a set of pre-computed local QoS values. This approach provides locally efficient component services that are then combined to formulate the composition. A major shortcoming of this approach lies in the fact that all QoS dimensions are considered independently, and no dependency or correlations among these components is taken into consideration. Second, in some scenarios where the QoS requirements are very aggressive, the approach translates the global requirements into very constrained local requirements. Thus, the algorithm fails to find a solution where a solution adhering to the global constraints may exist in the system.

Canfora et al. [26] proposed a Genetic Algorithms (GAs) based QoS-aware composite service binding and rebinding approach. This approach allows the service orchestrators to apply non-linear QoS aggregation formulae as compared to

linearization for the traditional approaches. However, their solution focuses on collecting usage pattern data to predict the need of re-binding for different invocation instances. This increases the system overhead and the service needs to be a part of the composition for it to be optimized. Second the proposed solution assumes that there will always be a solution that will meet all the requirement. Yu et al. [27] present a broker based solution for QoS based service selection. They present user-defined utility function for both of their models, i.e., combinatorial and graph based model. They consider invocation patterns in their utility functions and present different heuristic based approaches to find near optimal solutions. However, they present two different solutions for sequential executions and complex structures, i.e., loop etc and rely solely on the published QoS values. Zeng et al. [28] present an Integer Programming based solution that finds optimal service composition solutions. However, in case of multi-path scenarios they only optimize the solution based on the path with the highest probability as well as their definition of critical path uses the worst case scenarios of all the service in the execution path, which may not be a good approximation and hence is not suitable for large systems or systems with dynamic service needs.

In [29] authors propose a discrete probability distribution based model to solve the uncertainty of QoS values of the services. They incorporate composition control along with probability distribution on a uniform framework. The proposed approach works better than either using the mean or median and is independent of the any normality or uniformity constraints. This allows to directly convert the observed behavior in inputs for the system and a single analysis could approximate the results that could have been obtained by exhaustive black box simulations of the composition. This method however does not consider the dependency modeling behavior of the system and assumes that the probability distribution is independent of the factors like location of users and any other differentiating factors. In [30] present a soft constraint based solution for QoS service selection. They model selection characteristics using soft constraints, assign preferences to the penalties and constraints. Then the cumulative ranking of the above mentioned factors to compute the composite rank of the solution. This allows the users more flexibility incase services do not meet the exact composition criteria. However, the proposed approach does not consider the effect of dependency modeling and leaves it to the user to use a simple ranking based solution.

Collaborative Filtering approaches, i.e., both memory based and model based algorithms have been widely used in recommendation systems [31], [32], [33], [34]. Memory based approaches use the historic values of a user to recommend similar items [35], [36] on the other hand model based approaches [37], [38], [39], [40] apply different machine learning and statistical methods to learn user rating behaviors. Memory based approaches are relatively easier to implement, do not need any training set and can easily accommodate new rating data from users. However, they do not scale well as the data set grows. On the other hand item based approaches usually outperform memory based approaches in terms of scalability and quality of recommendation but they suffer from the fact that the model needs to be updated/rebuild

when new data points are received. There have been efforts to use collaborative filtering in service composition. Margaritis et al. [41] use collaborative filtering for QoS requirements for WS-BPEL scenarios. They combine the collaborative filtering based score and the QoS requirements of the users to calculate an aggregate score for WS-BEPL scenarios. Karta [42] have used both Pearson correlation and vector based similarity approaches in collaborative filtering on MovieLens data set for service selection, comparing their work with multidimensional recommender system. However, this approach uses static QoS ontology and ranking registry data.

8 CONCLUSION

In this paper we presented a multi-objective QoS optimization approach that uses the structure of a service based system by incorporating different invocation patterns, and uses collaborative filtering to predict the QoS values of different services to formulate an optimized solution. Using the orchestration pattern of a composition allows valuable information that could be exploited to further optimize different QoS value such as Throughput and Execution time. Our proposed approach uses this information and searches for services that are best suited for the current orchestration, hence providing a very optimal solution for the composition at hand. Since this approach is highly reliant on the observed QoS values of a Web service, we employed collaborative filtering to fill in the gaps and predict the QoS values for the missing service interactions. The experiment results show the advantage of our approach when compared to other methodologies.

ACKNOWLEDGMENTS

This research was made possible by NPRP 7-481-1-088 grant from the Qatar National Research Fund (a member of The Qatar Foundation). The statements made herein are solely the responsibility of the authors.

REFERENCES

- [1] Y. Zhang, Z. Zheng, and M. R. Lyu, "WSEXPRESS: A QoS-aware search engine for web services," in *Proc. IEEE Int. Conf. Web Services*, 2010, pp. 91–98.
- [2] Z. Malik and A. Bouguettaya, "RATEweb: Reputation assessment for trust establishment among web services," *Int. J. Very Large Data Bases*, vol. 18, no. 4, pp. 885–911, 2009.
- [3] M. Papazoglou, *Web Services: Principles and Technology*, Pearson Education, 2008.
- [4] D. A. D'Mello and V. S. Ananthanarayana, "A tree structure for web service compositions," in *Proc. 2nd Bangalore Annu. Comput. Conf.*, 2009, pp. 1–4.
- [5] D. A. Menasce, "Composing web services: A QoS view," *IEEE Internet Comput.*, vol. 8, no. 6, pp. 88–90, Nov./Dec. 2004.
- [6] J. Vanhatalo, H. Völzer, and J. Koehler, "The refined process structure tree," *Data Knowl. Eng.*, vol. 68, no. 9, pp. 793–818, 2009.
- [7] M. R. McLaughlin and J. L. Herlocker, "A collaborative filtering algorithm and evaluation metric that accurately model the user experience," in *Proc. 27th Annu. Int. ACM SIGIR Conf. Res. Develop. Inf. Retrieval*, 2004, pp. 329–336.
- [8] N. Mastronarde, K. Kanoun, D. Atienza, P. Frossard, and M. van der Schaar, "Markov decision process based energy-efficient online scheduling for slice-parallel video decoders on multicore systems," *IEEE Trans. Multimedia*, vol. 15, no. 2, pp. 268–278, 2013.
- [9] H. E. Carter, et al., "A decision model to estimate the cost-effectiveness of intensity modulated radiation therapy (IMRT) compared to three dimensional conformal radiation therapy (3DCRT) in patients receiving radiotherapy to the prostate bed," *Radiotherapy Oncology*, vol. 112, no. 2, pp. 187–193, Aug. 2014.

- [10] M. Pesce, D. Munaretto, and M. Zorzi, "A Markov decision model for source video rate allocation and scheduling policies in mobile networks," in *Proc. IEEE 13th Annu. Mediterranean Ad Hoc Netw. Workshop*, 2014, pp. 119–125.
- [11] J. Patrick, "A Markov decision model for determining optimal outpatient scheduling," *Health Care Manag. Sci.*, vol. 15, no. 2, pp. 91–102, Jun. 2012.
- [12] H. Wang, X. Zhou, X. Zhou, W. Liu, W. Li, and A. Bouguettaya, "Adaptive service composition based on reinforcement learning," in *Proc. 8th Int. Conf. Service-Oriented Comput.*, 2010, pp. 92–107.
- [13] E. D. Nitto, M. D. Penta, A. Gambi, G. Ripa, and M. L. Villani, "Negotiation of service level agreements: An architecture and a search-based approach," in *Proc. 5th Int. Conf. Service-Oriented Comput.*, 2007, pp. 295–306.
- [14] X. Niu and S. Wang, "Genetic algorithm for automatic negotiation based on agent," in *Proc. 7th World Congr. Intell. Control Autom.*, Jun. 2008, pp. 3834–3838.
- [15] F. Lecue, U. Wajid, and N. Mehandjiev, "Negotiating robustness in Semantic Web service composition," in *Proc. 7th IEEE Eur. Conf. Web Services*, 2009, pp. 75–84.
- [16] M. K. Sim, Y. Guo, and B. Shi, "BLGAN: Bayesian learning and genetic algorithm for supporting negotiation with incomplete information," *IEEE Trans. Syst. Man Cybern. B Cybern.*, vol. 39, no. 1, pp. 198–211, Feb. 2009.
- [17] C. Figueroa, N. Figueroa, A. Jofre, A. Sahai, Y. Chen, and S. Iyer, "A Game Theoretic Framework for SLA Negotiation," HP Lab., Palo Alto, CA, Tech. Rep. HPL-2008-5, 2008.
- [18] D. Booth, H. Haas, F. McCabe, E. Newcomer, M. Champion, C. Ferris, and D. Orchard, "Web services architecture," Oct. 2004.
- [19] F. Curbera, M. Duftler, R. Khalaf, W. Nagy, N. Mukhi, and S. Weerawarana, "Unraveling the web services web: An introduction to SOAP, WSDL, and UDDI," *IEEE Internet Comput.*, vol. 6, no. 2, pp. 86–93, Mar./Apr. 2002.
- [20] R. Chinnici, J.-J. Moreau, A. Ryman, and S. Weerawarana, "Web services description language (WSDL) version 2.0 part 1: Core language," W3C Recommendation, Jun. 26, 2007.
- [21] R. Aggarwal, K. Verma, J. Miller, and W. Milnor, "Constraint driven web service composition in METEOR-S," in *Proc. IEEE Int. Conf. Services Comput.*, 2004, pp. 23–30.
- [22] R. Berbner, M. Spahn, N. Repp, O. Heckmann, and R. Steinmetz, "Heuristics for QoS-aware web service composition," in *Proc. IEEE Int. Conf. Web Services*, 2006, pp. 72–82.
- [23] D. Ardagna and B. Pernici, "Adaptive service composition in flexible processes," *IEEE Trans. Softw. Eng.*, vol. 33, no. 6, pp. 369–384, Jun. 2007.
- [24] L. Ai, M. Tang, and C. Fidge, "Partitioning composite web services for decentralized execution using a genetic algorithm," *Future Generation Comput. Syst.*, vol. 27, no. 2, pp. 157–172, 2011.
- [25] 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.
- [26] G. Canfora, M. Di Penta, R. Esposito, and M. L. Villani, "A framework for QoS-aware binding and re-binding of composite web services," *J. Syst. Softw.*, vol. 81, no. 10, pp. 1754–1769, 2008.
- [27] T. Yu, Y. Zhang, and K.-J. Lin, "Efficient algorithms for web services selection with end-to-end QoS constraints," *ACM Trans. Web*, vol. 1, no. 1, 2007, Art. no. 6.
- [28] L. Zeng, B. Benatallah, A. 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.
- [29] D. Ivanovic, P. Kaowichakorn, and M. Carro, "Towards QoS prediction based on composition structure analysis and probabilistic environment models," in *Proc. Int. Conf. Workshop Principles Eng. Service-Oriented Syst.*, 2013, pp. 11–20.
- [30] M. A. Zemni, S. Benbernou, and M. Carro, "A soft constraint-based approach to QoS-aware service selection," in *Proc. 8th Int. Conf. Service-Oriented Comput.*, 2010, pp. 596–602.
- [31] R. Burke, "Hybrid recommender systems: Survey and experiments," *User Modeling User-Adapted Interaction*, vol. 12, no. 4, pp. 331–370, 2002.
- [32] G. Linden, B. Smith, and J. York, "Amazon. com recommendations: Item-to-item collaborative filtering," *IEEE Internet Comput.*, vol. 7, no. 1, pp. 76–80, Jan./Feb. 2003.
- [33] H. Ma, I. King, and M. R. Lyu, "Effective missing data prediction for collaborative filtering," in *Proc. 30th Annu. Int. ACM SIGIR Conf. Res. Develop. Inf. Retrieval*, 2007, pp. 39–46.
- [34] P. Resnick, N. Iacovou, M. Suchak, P. Bergstrom, and J. Riedl, "GroupLens: An open architecture for collaborative filtering of netnews," in *Proc. ACM Conf. Comput. Supported Cooperative Work*, 1994, pp. 175–186.
- [35] A. Krzywicki, et al., "Collaborative filtering for people-to-people recommendation in online dating: Data analysis and user trial," *Int. J. Human-Comput. Studies*, vol. 76, pp. 50–66, 2015.
- [36] R. Jin, J. Y. Chai, and L. Si, "An automatic weighting scheme for collaborative filtering," in *Proc. 27th Annu. Int. ACM SIGIR Conf. Res. Develop. Inf. Retrieval*, 2004, pp. 337–344.
- [37] G.-R. Xue, et al., "Scalable collaborative filtering using cluster-based smoothing," in *Proc. 28th Annu. Int. ACM SIGIR Conf. Res. Develop. Inf. Retrieval*, 2005, pp. 114–121.
- [38] T. Hofmann, "Collaborative filtering via gaussian probabilistic latent semantic analysis," in *Proc. 26th Annu. Int. ACM SIGIR Conf. Res. Develop. Inf. Retrieval*, 2003, pp. 259–266.
- [39] T. Hofmann, "Latent semantic models for collaborative filtering," *ACM Trans. Inf. Syst.*, vol. 22, no. 1, pp. 89–115, 2004.
- [40] T. Hofmann and D. Hartmann, "Collaborative filtering with privacy via factor analysis," in *Proc. ACM Symp. Appl. Comput.*, 2005, pp. 791–795.
- [41] D. Margaris, C. Vassilakis, and P. Georgiadis, "An integrated framework for adapting WS-BPEL scenario execution using QoS and collaborative filtering techniques," *Sci. Comput. Programming*, vol. 98, pp. 707–734, 2015.
- [42] K. Karta, "An investigation on personalized collaborative filtering for web service selection," Honours Programme thesis, School Comput. Sci. Softw. Eng., Univ. Western Australia, Perth, WA, 2005.



Khayyam Hashmi is currently working toward the PhD degree in the Department of Computer Science, Wayne State University. He is a member of the Services Computing Research (SCORE) Laboratory. His research interests include service computing, distributed systems, Semantic Web, and software process engineering.



Zaki Malik is an assistant professor in the School of Information Security and Applied Computing, Eastern Michigan University, where he heads the Services Computing Research (SCORE) Laboratory. His research interests lie broadly in service computing, reliable distributed systems, security and privacy, and semantic integration.



Abdelmounaam Rezgui is an assistant professor in the Department of Computer Science, New Mexico Tech, University of Socorro. His research interests include high performance spatial cloud computing, spatio-temporal data intensive applications on the clouds, data and service quality, power saving techniques for cellular networks, service-oriented sensor networks, semantics-based data interoperability, and integration.



Abdelkarim Erradi is an assistant professor in the Department of Computer Science and Engineering, Qatar University. His research interests include cloud computing particularly, software composition for distributed systems, autonomic computing systems, policy-based management of middleware systems particularly Web services, and Web service technologies.