

A quantitative approach to coordinated scaling of resources in complex cloud computing workflows

Laura Carnevali¹, Marco Paolieri², Benedetta Picano¹, Riccardo Reali¹,
Leonardo Scommegna¹, and Enrico Vicario¹

¹ University of Florence, Department of Information Engineering, Italy
{laura.carnevali, benedetta.picano, riccardo.reali,
leonardo.scommegna, enrico.vicario}@unifi.it

² University of Southern California, Department of Computer Science, USA
paolieri@usc.edu

Abstract. Resource scaling is widely employed in cloud computing to adapt system operation to internal (i.e., application) and external (i.e., environment) changes. We present a quantitative approach for coordinated vertical scaling of resources in cloud computing workflows, aimed at satisfying an agreed Service Level Objective (SLO) by improving the workflow end-to-end response time distribution. Workflows consist of IaaS services running on dedicated clusters, statically reserved before execution. Services are composed through sequence, choice/merge, and balanced split/join blocks, and have generally distributed (i.e., non-Markovian) durations possibly over bounded supports, facilitating fitting of analytical distributions from observed data. Resource allocation is performed through an efficient heuristic guided by the mean makespans of sub-workflows. The heuristic performs a top-down visit of the hierarchy of services, and it exploits an efficient compositional method to derive the response time distribution and the mean makespan of each sub-workflow. Experimental results on a workflow with high concurrency degree appear promising for feasibility and effectiveness of the approach.

Keywords: Cloud computing · coordinated scaling · stochastic workflow · end-to-end response time distribution · complex workflow structure

1 Introduction

Cloud Computing (CC) systems [15, 7] need to store, manage, and process enormous amounts of data continuously generated by a variety of sources within the Internet of Things (IoT) [27]. Excessive network traffic or heavy computational workload may lead to violations of Quality of Service (QoS) attributes granted through Service level Agreements (SLAs) [28]. Therefore, CC systems must autonomously adapt their operation in response to time-varying changes both in the software system itself and in its operating environment [31, 25]. Adaptation can be achieved through *autoscaling* systems [12], which dynamically change software configurations and provision hardware resources on demand, with the

goal of continuously satisfying cost objectives as well as non-functional Service Level Objectives (SLOs), i.e., specific measures agreed within an SLA. Scaling actions can be *horizontal*, if the system adds or removes containers or Virtual Machines (VMs) where services can be deployed, or *vertical*, if the system changes specifications of those containers or VMs, e.g., CPU cores or available memory.

Horizontal scaling can optimize resource provisioning for *individual* services orchestrated in larger applications [19, 1, 13, 3, 36], e.g., composite web services [8], Functions as a Service (FaaS) platforms [32, 22], microservice architectures [2]. In [19], bottlenecks in a multi-tier application are automatically detected and resolved, minimizing the number of web servers and database instances while guaranteeing a maximum response time. Dynamic scaling of the number of VMs in cloud services is performed based on the number of active sessions of each web server instance [13], using queueing theory to estimate demand [1], and leveraging also time series analysis to forecast load intensity [3]. Few approaches exploit *coordinated* scaling of resources to avoid undesired effects of *local* scaling like bottleneck shifting and oscillations [36], e.g., by reconfiguring services of small web applications together [34], by exploiting time-series analysis and queueing theory to determine the number of VM instances that minimizes energy consumption without violating SLAs [4], or by collectively providing application tiers with a number of servers or VMs that guarantees meeting contracted [36] or average response times [6]. Though horizontal scaling has received more attention [12] and has better support from cloud vendors [14], being easier to implement and manage, it performs coarse-grained adaptation through static replication of VMs or containers with fixed-size configurations, and it suffers from non-negligible lags to instance and start VMs or containers [38], which, despite lag-mitigating actions like dynamic VM cloning [23], may negatively affect time-critical applications.

Vertical scaling performs fine-grained resource adaptation by modifying attributes of VMs or containers [33, 20, 14], thus limiting resource over-provisioning and resulting preferable for applications with time-critical requirements. In [33], CPU voltage and frequency of VMs in multi-tenant cloud systems are individually adapted to meet SLOs, supporting migration to new VMs in case of overloading. Optimization of the amount of CPU and memory allocated to a cloud application is performed in [14] to meet requirements on mean response time, exploiting a performance model based on an inverse relationship between the application mean response time and the number of allocated CPU cores [20]. In [38], CPU power tuning and hotplugging are performed to improve CPU usage efficiency in a web server, with minimum SLA violation rate. Few approaches perform vertical scaling in a *coordinated* manner. In [21], soft resources of each server of a web application (e.g., number of server threads and database connections) are allocated based on measured throughput and concurrency. A resource-management framework is defined in [24] to manage shared resources among microservices, exploiting machine learning methods both to localize microservice instances responsible for SLO violations and to define methods to mitigate resource contention. Horizontal and vertical scaling are combined in [18] to determine a load distribution policy for co-located distributed applications by ex-

exploiting multi-class queueing networks and model predictive control, and in [16] to adapt the number of replicas and the CPU capacity of each microservice by using layered queueing networks to assess potential performance improvement of scaling actions.

The few approaches that address coordinated resource scaling [25] mainly consider simple cloud applications consisting of few services orchestrated as sequential workflows. Notably, no approach takes into account the e2e response time distribution in scaling decisions, which instead becomes relevant when SLAs are characterized by soft deadlines and penalty functions [26] defined as rewards calculated from such distribution.

In this paper, we present an efficient approach to perform *coordinated* vertical scaling of resources in complex stochastic workflows, with the aim of satisfying an SLO by improving the workflow *e2e response time distribution*. Specifically, workflows compose IaaS services running on dedicated clusters whose size must be determined in advance, reserving and statically assigning resources to services before execution. Services are composed through sequence, fork-join, and choice-merge patterns [29], and have generally distributed (GEN) response times possibly with bounded supports, facilitating representation of real-time constraints and fitting of analytical distributions from observed data. Each service is characterized by a job size [5], representing its mean makespan with a unitary amount of assigned resources; we assume the mean makespan to be inversely proportional to the amount of assigned resources [20, 14, 30]. We define a heuristic that uses a structured workflow model [10] to perform a top-down visit of the hierarchy of services, assigning resources so as to minimize the mean makespan of the workflow e2e response time and to satisfy the agreed SLO. To this end, the heuristic exploits an efficient compositional analysis method [11] to derive the response time distribution of each sub-workflow and to compute its mean makespan. Feasibility and effectiveness of the approach are assessed on a non-trivial synthetic workflow stressing computational complexity. Experimental results show that the heuristic is effective at improving the e2e response time distribution of the entire workflow, and very efficient, enabling its application at runtime in reaction to QoS changes.

In the framework of [9], our approach is defined by the following attributes: the *goal of resource adaptation* is to ensure that workflow execution fulfils non-functional requirements specified by percentiles of quality attributes, i.e., that the mean makespan of the workflow response time satisfies the agreed SLO; the *stage of system lifetime* at which resource adaptation is performed is the runtime stage with proactive mode (i.e., anticipating resource adaptation), though the approach can be applied also in reactive mode (i.e., after changes in quality attributes) as shown by the experimental results; the *composition level* at which resource adaptation is performed involves both services (i.e., abstract composition made of tasks orchestrated by some composition logic) and workflow (i.e., concrete composition where tasks of an abstract composition are mapped to implementations); the *scope of resource adaptation*, in terms of *number of systems* and *granularity* of adaptation, considers a single system and a single request; *adap-*

tation actions mainly consist of service tuning operations changing behavior of concrete services (e.g., reducing the mean makespan by increasing the amount of resources), though adaptation is performed in a coordinated manner; and, resource adaptation is performed by a *single authority*.

The rest of the paper is organized as follows. Section 2 recalls the hierarchical formalism for workflow modeling and the compositional method for evaluation of the workflow e2e response time distribution. Section 3 illustrates the proposed resource assignment method. Section 4 presents the experimental results achieved on a complex workflow. Finally, Section 5 draws conclusions and outlines possible extensions and improvements of the proposed approach.

2 Background: Workflow Modeling and Evaluation

We model workflows as recursive compositions of blocks specified by Stochastic Time Petri Nets (STPNs) [37]. Each STPN block has a single starting place, which receives a token when workflow execution starts, and a single final place, which receives a token when workflow execution eventually ends with probability 1 (w.p.1). As shown in Fig. 1, blocks model sequential, balanced split/join, and choice/merge workflow patterns [29, 39], with the following EBNF syntax:

$$\begin{aligned} \text{BLOCK} := & \text{ACT} \mid \text{SEQ}\{\text{BLOCK}_1, \dots, \text{BLOCK}_n\} \\ & \mid \text{AND}\{\text{BLOCK}_1, \dots, \text{BLOCK}_n\} \mid \text{XOR}\{\text{BLOCK}_1, \dots, \text{BLOCK}_n, p_1, \dots, p_n\} \end{aligned} \quad (1)$$

where ACT is an elementary activity (e.g., block A in Fig. 1b), $\text{SEQ}\{\text{BLOCK}_1, \dots, \text{BLOCK}_n\}$ models n sequential blocks $\text{BLOCK}_1, \dots, \text{BLOCK}_n$ (e.g., block S1 in Fig. 1b), $\text{AND}\{\text{BLOCK}_1, \dots, \text{BLOCK}_n\}$ models n concurrent blocks $\text{BLOCK}_1, \dots, \text{BLOCK}_n$ (e.g., block A1 in Fig. 1b), and $\text{XOR}\{\text{BLOCK}_1, \dots, \text{BLOCK}_n, p_1, \dots, p_n\}$ models n alternative blocks $\text{BLOCK}_1, \dots, \text{BLOCK}_n$ with probability p_1, \dots, p_n , respectively (e.g., block X1 in Fig. 1b). This workflow model can be represented as a *structure tree* [10] $S = \langle N, n_0 \rangle$, where N is the set of nodes (i.e., blocks) and $n_0 \in N$ is the root node (i.e., the entire workflow). In turn, each node $n_i \in N$ is a tuple $\langle C_i, \text{type}_i \rangle$, where C_i is the set of the children nodes of n_i and $\text{type}_i \in \{\text{ACT}, \text{SEQ}, \text{AND}, \text{XOR}\}$ is the type of the block modeled by node n_i , e.g., in Fig. 1a, node A1 models an AND block composing nodes I and J.

For complex workflows made of several concurrent activities with duration characterized by GEN Cumulative Distribution Functions (CDFs) possibly with bounded supports, the e2e response time distribution cannot be evaluated by transient analysis [17] of the STPN modeling the entire workflow. To address the issue, a compositional approach is defined in [11], which first performs a top-down visit of the structure tree to estimate the analysis complexity of blocks, then evaluates the response time distribution of the identified sub-workflows in isolation, and finally performs a bottom-up recomposition of the obtained results. In particular, for workflows defined by well-nested composite blocks as in this paper (i.e., composition of AND, SEQ, and XOR blocks), the exact e2e response time distribution can be evaluated by recursive numerical analysis [11].

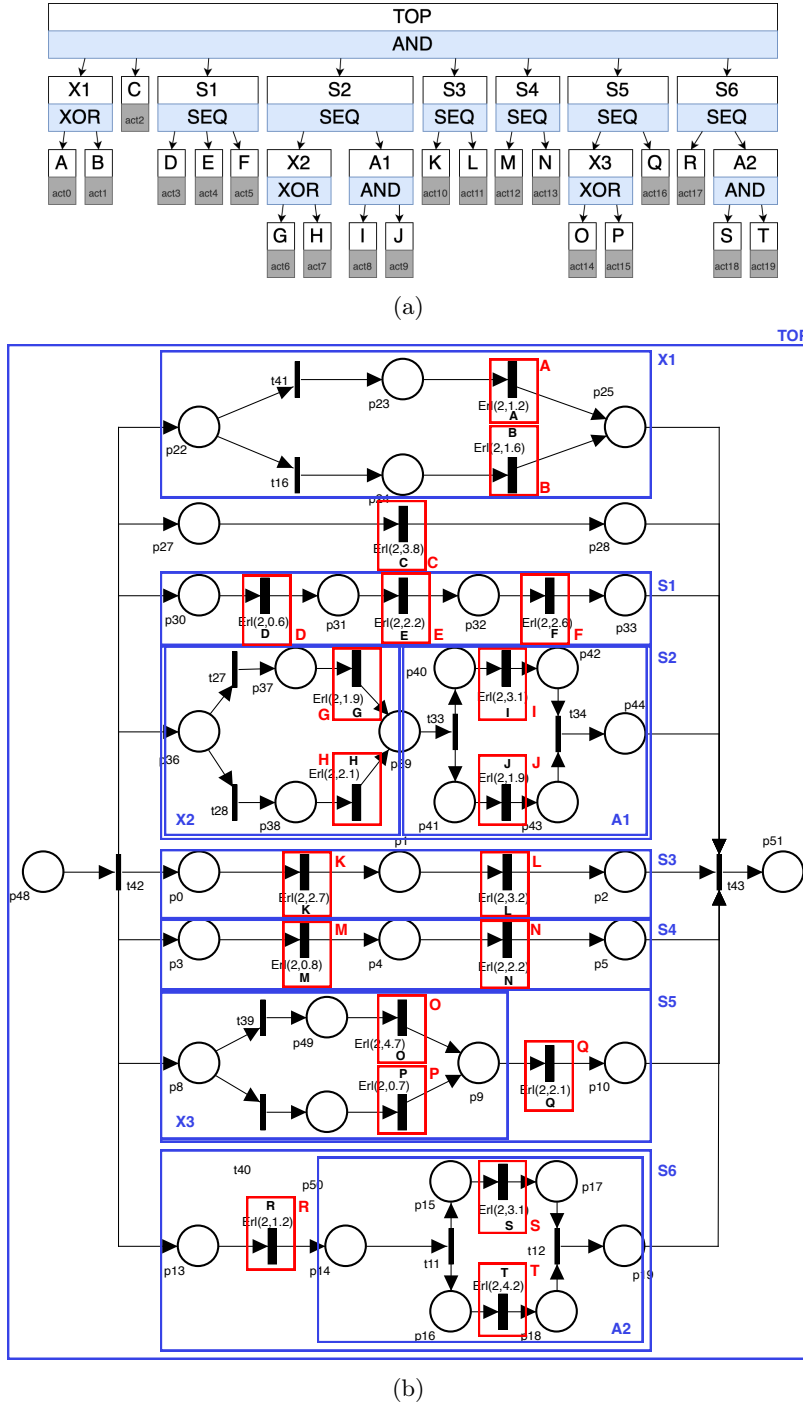


Fig. 1: (a) Structure tree and (b) STPN modeling a synthetic workflow.

3 Coordinated Resource Scaling Heuristics

In this section, we present a vertical scaling heuristics for coordinated allocation of resources to each workflow activity, improving the workflow e2e response time distribution. We illustrate how the workflow structure tree is visited (Section 3.1) and the scaling decisions for each node type (Section 3.2).

3.1 Heuristics Overview

The approach is based on the concept of *job size* of an activity: similarly to [20, 14, 30], we assume that, after allocating R resources for the activity, its mean makespan T_R is given by X/R , where $X := RT_R \forall R$ is the job size of the activity, i.e., the invariant amount of work to be completed with the assigned resources. For each node of the structure tree of a workflow, given a new resource allocation, the mean makespan can be estimated as a function of the invariant job sizes of the children of the node; this function depends on the node type, e.g., ACT, SEQ, AND, XOR.

The heuristics performs a top-down visit of the structure tree and splits the resources assigned to each node among its children by solving an optimization problem, until the resource allocation of all elementary activities is determined. By proceeding over successive locally refinements on the structure tree, the total response time of the workflow can be improved in a coordinated fashion. Note that, as long as the job sizes of the tasks are known, the method can identify an ideal resource allocation not only when the total amount of resources is redistributed, but also when it is incremented or decremented.

Let $S = \langle N, n_0 \rangle$ be a structure tree as defined in Section 2. We extend the definition of node $n_i \in N$ with the tuple $\langle C_i, \text{type}_i, R_i, T_i, X_i \rangle$, where $R_i \in \mathbb{R}_{>0}$ is the amount of resources initially assigned to n_i , $T_i \in \mathbb{R}_{\geq 0}$ is the mean makespan of n_i , and X_i is the job size of n_i . For each non-leaf node $n_i \in N$, the number of resources of n_i is the sum of the number of resources allocated to its children nodes, i.e., $\forall n_i \in N$ such that $C_i \neq \emptyset$, $R_i = \sum_{n_j \in C_i} R_j$. To coordinately adapt the resource provisioning of the activities of a workflow, the approach performs a *top-down* visit of the workflow structure tree.

- Initially, an arbitrary amount of resources R_0^{in} is assigned to the root node n_0 .
- For each non-leaf node n_k , the amount of input resources R_k^{in} is split by assigning an amount R_j^* to each child node n_j , i.e., $\sum_{n_j \in C_k} R_j^* = R_k^{\text{in}}$; the assignment depends on the node type.

By induction, the sum of the amounts of resources allocated to the leaf nodes is equal to the amount of resources of the root node, i.e., $\sum_{n_k | C_k = \emptyset} R_k^* = R_0^{\text{in}}$.

3.2 Scaling Decisions

We characterize the different resource scaling decisions based on the node types.

Elementary Activities. Let n_i be an elementary activity, i.e. $\text{type}_i = \text{ACT}$, R_i be the amount of resources assigned by the approach during the previous step,

and T_i the resulting mean makespan. Since we assume job sizes to be invariant, a new resource allocation R_i^* results in a new mean makespan equal to

$$T_i^* = \frac{R_i^*}{R_i} T_i. \quad (2)$$

Note that Eq. (2) determines a transformation of the parameters of the distribution of the node, which depends on the symbolic form of the distribution.

Sequential Activities. Let n_k be an activity with $\text{type}_k = \text{SEQ}$ and children $C_k = \{i, j\}$, and let R_k^{in} be the amount of resources to be split. The mean makespan of node n_k can be obtained as

$$T_k = T_i + T_j = \frac{X_i}{R_i} + \frac{X_j}{R_k^{\text{in}} - R_i} \quad (3)$$

which has a minimum when the following resources are allocated to node n_i :

$$R_i^* = \frac{\sqrt{X_i}}{\sqrt{X_i} + \sqrt{X_j}} R_k^{\text{in}}. \quad (4)$$

The result is obtained by imposing $\frac{dT_k}{dR_i} = 0$, and it can be extended by induction to the sequence of $K > 2$ activities, $\text{SEQ}(n_1, \dots, n_K)$:

$$R_i^* = \frac{\sqrt{X_i}}{\sum_{i=1}^K \sqrt{X_i}} R_k^{\text{in}}. \quad (5)$$

Concurrent Activities. Let n_k be an activity with $\text{type}_k = \text{AND}$, children $C_k = \{i, j\}$, and an amount R_k^{in} of resources to be split. The mean response time of n_k is:

$$T_k = \max(T_i, T_j) = \max\left(\frac{X_i}{R_i}, \frac{X_j}{R_k^{\text{in}} - R_i}\right) \quad (6)$$

Since response time is not defined as an explicit function of R_i , the minimum cannot be evaluated exploiting the Fermat theorem. Hence, we provide a heuristics evaluation of R_i^* which depends on a parameter $\alpha \in \mathbb{R}$ which modulates the weight of the job size in determining the solution. In particular, R_i^* is evaluated by imposing equality between the ratio of response times of node n_k children, and the resources provisioned to the children powered by $\alpha + 1$:

$$\frac{X_i}{R_i^{\alpha+1}} = \frac{X_j}{(R_k^{\text{in}} - R_i)^{\alpha+1}} \quad (7)$$

This leads to the allocation:

$$R_i^* = \frac{\alpha+1\sqrt{X_i}}{\alpha+1\sqrt{X_i} + \alpha+1\sqrt{X_j}} R_k^{\text{in}} \quad (8)$$

The solution is extended to m activities by induction:

$$R_i^* = \frac{\alpha+1\sqrt{X_i}}{\sum_{j \in C_k} \alpha+1\sqrt{X_j}} R_k^{\text{in}} \quad (9)$$

Note that, for $\alpha = 0$, the heuristics determines R_i^* by imposing equality between the response times of the considered node children.

Alternative Activities. Let n_k be an activity with $type_k = \text{XOR}$, children $C_k = \{i, j\}$ having probabilities p_i and $p_j = 1 - p_i$ to occur, and R_k^{in} the amount of resources to be split. The mean makespan of n_k is

$$T_k = p_i T_i + p_j T_j = p_i \frac{X_i}{R_i} + p_j \frac{X_j}{R_k^{\text{in}} - R_i} \quad (10)$$

which has the minimum

$$T_k^* = \frac{(\sqrt{p_i X_i} + \sqrt{p_j X_j})^2}{R_k^{\text{in}}} \quad (11)$$

when the following resources are allocated to node i :

$$R_i^* = \frac{\sqrt{p_i X_i}}{\sqrt{p_i X_i} + \sqrt{p_j X_j}} R_k^{\text{in}} \quad (12)$$

which is in turn obtained by exploiting the Fermat theorem. As the solution shows, the optimal allocation of a XOR node is a generalization of the optimal allocation for a SEQ node. Hence, the extension to the case of more than 2 activities, can be borrowed from the SEQ node type solution. Also note that the available resources R_k^{in} are split among the activities i and j : this assumption is useful to model workflows in microservice and service-oriented applications, where each service has a reserved amount of resources. In contrast, for microservices deployed with FaaS cloud solutions, resources are allocated on-demand for each service execution: in this case, cloud costs are accrued only for the resources of the selected service; the expected cost is in this case $p_i R_i + p_j R_j$ instead of $R_i + R_j$, resulting in a different optimal allocation.

4 Experimental Evaluation

In this section, we assess feasibility and effectiveness of the proposed heuristics. First, we show how different values of parameter α produce different resource allocations, with consequent different improvement of the workflow e2e response time CDF (Section 4.1). The experiment is performed for two different initial resource allocations. Then, we test the ability of the heuristics to meet an agreed SLO while minimizing the amount of allocated resources (Section 4.2).

Both experiments are performed on the synthetic workflow of Fig. 1, which consists of 20 elementary blocks combined through well-nested patterns, yielding a model with up to 10 concurrent activities. For each elementary activity, we assume that the response time CDF achieved with a given resource assignment is known, either from measurements of previous implementations or by contract. To easily manage the linear relation between response time and allocated resources assumed by the performance model of Section 3.1, and to facilitate the

ACT	r_{start}	r_0	$r_{1/4}$	$r_{1/2}$	r_1	r_2	r_4	ACT	r_{start}	r_0	$r_{1/4}$	$r_{1/2}$	r_1	r_2	r_4
A	1.00	0.52	0.64	0.73	0.85	0.97	1.07	A	0.37	0.21	0.32	0.43	0.60	0.80	0.99
B	1.00	0.55	0.68	0.78	0.90	1.03	1.14	B	0.37	0.22	0.34	0.45	0.63	0.85	1.05
C	1.00	0.25	0.41	0.57	0.84	1.23	1.65	C	0.37	0.01	0.21	0.33	0.59	1.01	1.52
D	1.00	2.85	2.51	2.29	2.01	1.75	1.54	D	0.37	1.13	1.27	1.34	1.41	1.44	1.42
E	1.00	1.49	1.31	1.19	1.05	0.91	0.81	E	0.37	0.59	0.66	0.70	0.74	0.75	0.74
F	1.00	1.37	1.20	1.01	0.97	0.84	0.74	F	0.37	0.54	0.61	0.65	0.68	0.69	0.68
G	1.00	0.72	0.71	0.70	0.68	0.65	0.62	G	3.75	2.49	2.22	2.02	1.76	1.48	1.26
H	1.00	0.84	0.83	0.82	0.79	0.75	0.72	H	0.37	0.92	0.82	0.75	0.65	0.55	0.46
I	1.00	0.63	0.66	0.68	0.69	0.69	0.68	I	3.75	3.14	2.63	2.28	1.83	1.40	1.08
J	1.00	1.03	0.98	0.94	0.88	0.81	0.75	J	0.37	0.51	0.62	0.68	0.74	0.77	0.75
K	1.00	0.54	0.67	0.77	0.89	1.03	1.15	K	0.37	0.21	0.34	0.45	0.63	0.85	1.05
L	1.00	0.50	0.61	0.70	0.82	0.95	1.05	L	0.37	0.20	0.31	0.41	0.58	0.78	0.97
M	1.00	1.61	1.67	1.69	1.69	1.68	1.65	M	0.37	0.64	0.84	0.99	1.19	1.38	1.51
N	1.00	0.97	1.00	1.02	1.02	1.01	0.99	N	0.37	0.39	0.51	0.60	0.72	0.83	0.91
O	1.00	0.48	0.47	0.46	0.44	0.42	0.40	O	3.75	2.05	1.93	1.83	1.67	1.47	1.29
P	1.00	1.70	1.66	1.63	1.56	1.48	1.40	P	0.37	2.29	2.16	2.04	1.86	1.64	1.44
Q	1.00	1.19	1.16	1.14	1.09	1.04	0.98	Q	0.37	1.01	0.95	0.90	0.82	0.72	0.63
R	1.00	1.49	1.52	1.52	1.52	1.49	1.45	R	0.37	1.46	1.44	1.40	1.32	1.20	1.09
S	1.00	0.71	0.71	0.70	0.68	0.66	0.63	S	3.75	2.68	2.51	2.33	2.04	1.67	1.34
T	1.00	0.53	0.56	0.57	0.59	0.59	0.59	T	0.37	0.20	0.31	0.41	0.55	0.70	0.80

(a) Balanced initial res. allocation. (b) Unbalanced initial res. allocation.

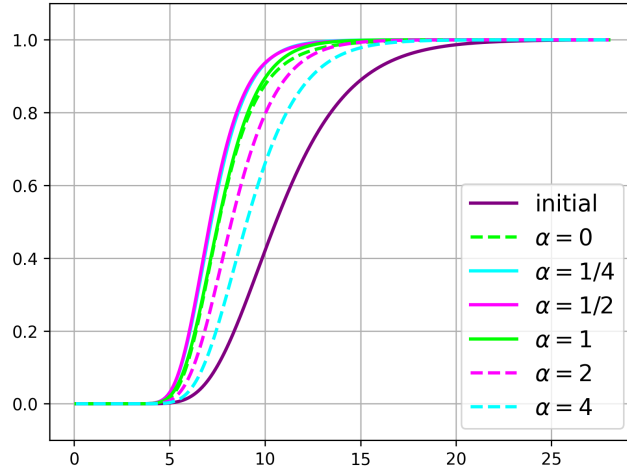
Table 1: Resources allocated to the activities of the workflow of Fig. 1, before (column r_{start}) and after heuristics execution with different values of α (columns r_α with $\alpha \in \{0, \frac{1}{4}, \frac{1}{2}, 1, 2, 4\}$): (a) balanced and (b) unbalanced initial allocation.

interpretation of the experimental results, without loss of generality, we consider Erlang CDFs for the response times of elementary activities. In particular, we consider Erlang CDFs with 5-phases and rates randomly selected in $[0, 5]$, so guaranteeing variability in expected response times of activities. We remark that any numerical CDF could be considered as well, or any analytical CDF in the class of expolynomial functions (also termed exponents [35]) supported by the compositional analysis technique of [11] exploited by the proposed heuristics.

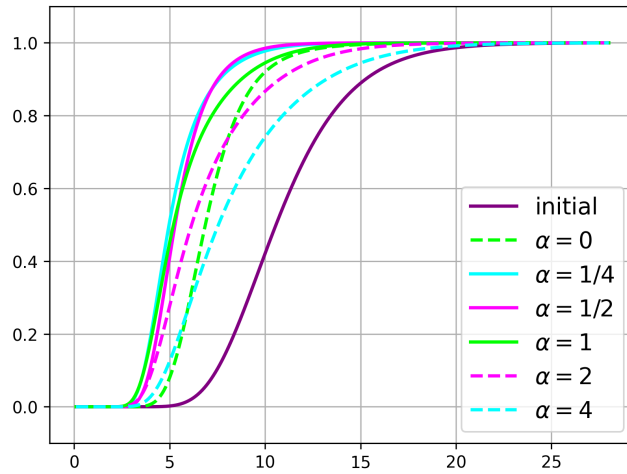
All the experiments reported in this section have been performed on a MacBook Pro 2021 equipped with an Apple M1 Pro processor.

4.1 Heuristics Sensitivity to Parameter α

We consider two different initial allocations of resources to the activities of the workflow of Fig. 1. In the first allocation, shown in column r_{start} of Table 1a, each activity is assigned 1 resource, for a total amount of 20 resources (*balanced initial resource allocation*). In the second allocation, shown in column r_{start} of Table 1b, activities G, I, O, and S are each assigned an amount of 3.75 resources, while each of the other resources is assigned an amount of 0.37 resources, for a total



(a) Balanced initial resource allocation.



(b) Unbalanced initial resource allocation.

Fig. 2: CDF of the e2e response time of the workflow of Fig. 1, obtained through the execution of the proposed heuristics with different values of parameter α , by assuming: (a) the balanced initial resource allocation shown in Table 1a, and (b) the unbalanced initial resource allocation shown in Table 1b.

amount of 20 resources again (*unbalanced initial resource allocation*). Tables 1a and 1b also show the resource allocation computed by the proposed heuristics.

We evaluate how the resource allocation provided by the heuristics improves the workflow e2e response time CDF for different values of $\alpha \in \{0, \frac{1}{4}, \frac{1}{2}, 1, 2, 4\}$. Results reported in Figs. 2a and 2b show an improvement for any considered

value of α and for any of the two initial allocations of resources, with better results achieved for $\alpha \in [0, 1]$ and nearly the best result obtained for $\alpha = \frac{1}{4}$. This result suggest that, being $\alpha = \frac{1}{n}$, the e2e response time CDF improves as n increases, up to a certain value beyond which the CDF gets worse.

If the initial resource allocation is balanced, different values of $\alpha \in [0, 1]$ produce nearly comparable e2e response time CDFs. This result suggests that, in this case, balancing the mean makespan of the children of AND nodes determines a good resource allocation, significantly improving the workflow e2e response time CDF. Conversely, if the initial resource allocation is unbalanced, better results are obtained for values of $\alpha \in (0, 1]$, with a worse result obtained for $\alpha = 0$ with respect to the case of balanced initial allocation of resources.

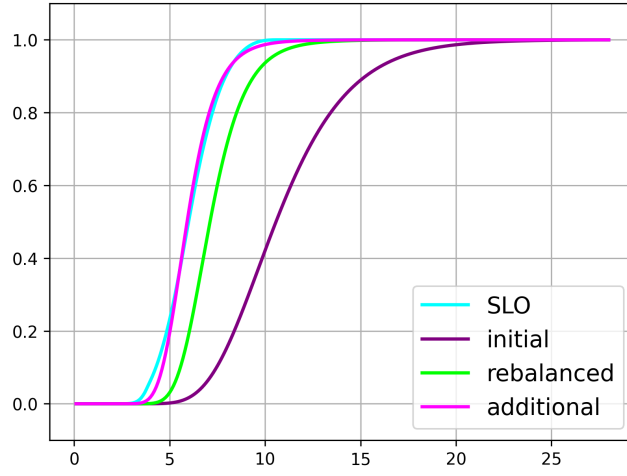
It is worth noting that, for both variants of the experiment, the proposed heuristics runs in 0.2 s on average, proving to be efficient even for complex workflows, which is an essential requirement for modern microservice architectures where hundreds of services are orchestrated as workflows of activities.

4.2 Heuristics Ability to Achieve SLO Guarantees

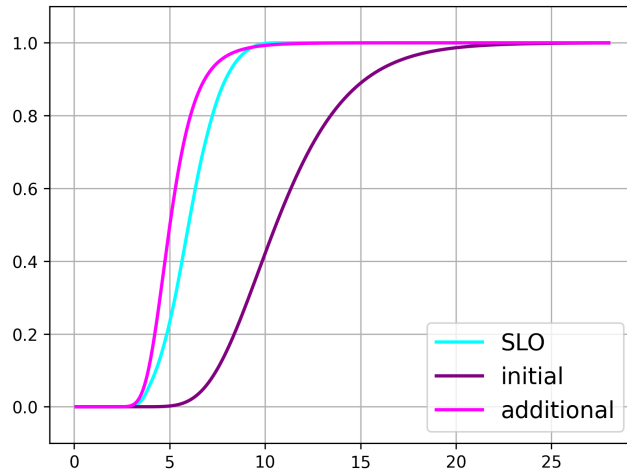
For the workflow of Fig. 1 with the initial balanced resource allocation of Table 1a, we consider an SLO expressed in terms of e2e response time CDF. In particular, the SLO is obtained as the e2e response time CDF of a randomly generated well-nested workflow, having expected response time equal to $T_{\text{SLO}} = 6.17$ s. As shown in Fig. 3, with this allocation of resources, the workflow e2e response time CDF (violet curve) does not meet the SLO (light blue curve). Rebalancing the existing resources by applying the proposed heuristics with $\alpha = \frac{1}{4}$ (which is the value yielding the best results in Section 4.1) produces an improvement of the e2e response time CDF, which however still violates the required SLO. Therefore, an additional amount of resources is needed not to violate the SLO.

To determine the new amount of resources, we consider two different strategies whose results are reported in Fig. 3. We perform an experiment (*reallocate resources first*) where we first compute a new resource allocation through our heuristics with $\alpha = \frac{1}{4}$ (green curve in Fig. 3a), by assuming that the total amount of allocated resources does not change, i.e., $R_0 = 20$. Then, we evaluate the additional amount of resources that is needed to meet the specified SLO, by exploiting the assumption of invariance of the job size of a workflow (as discussed in Section 3.1). In fact, by knowing the expected response time T_{SLO} of the SLO, the initial amount of allocated resources $R_0 = 20$, and the expected response time T_0 obtained after resource allocation (i.e., considering the resource allocation of column $r_{1/4}$ of Table 1a as initial resource allocation), the amount of resources required not to violate the SLO can be computed as $R^* = \frac{R_0 T_0}{T_{\text{SLO}}}$. In particular, the additional amount of resources turns out to be equal to 4.06. Finally, the amount R^* is allocated to the activities using the heuristics, and a new e2e response time CDF is computed (fuchsia curve in Fig. 3a).

Then, we perform a variant of the experiment (*add resources first*), where we directly evaluate the additional amount of resources needed to meet the specified



(a) Reallocate resources first.



(b) Add resources first.

Fig. 3: CDF of the e2e response time of the workflow of Fig. 1, obtained through two strategies: (a) reallocating resources through the heuristics, determining the amount of resources needed to satisfy the SLO, and allocating resources again through the heuristics; and, (b) determining the amount of resources needed to satisfy the SLO and allocating resources through the heuristics.

SLO as $R^* = \frac{R_0 T'_0}{T_{\text{SLO}}}$, where T'_0 is the workflow expected response time obtained by considering the initial resource allocation of column r_{start} of Table 1a. In particular, the additional amount of resources turns out to be equal to 16.15.

The allocation of the increased amount of resources through our heuristics yields a new e2e response time CDF (fuchsia curve in Fig. 3b).

Fig. 3 shows that the e2e response time CDF provided by the reallocate-resources-first strategy is stochastically larger than the one provided by the add-resources-first strategy, and is characterized by a larger expected response time. However, the add-resources-first strategy allocates a larger number of resources, which actually turn out to be over-provisioned, given that the obtained e2e response time CDF is stochastically lower than the SLO. Moreover, note that the time to calculate the new resource provisioning is 0.89s for the reallocate-resources-first strategy and 0.51s for the add-resources-first strategy, meaning that there is not a significant loss in performance when the heuristics is executed twice. Therefore, the reallocate-resources-first strategy is preferable.

5 Conclusions

We have presented a heuristics to perform coordinated scaling of resources in cloud computing workflows, with the aim of improving the e2e response time CDF. The heuristic is developed around the concept of job size of an activity, which is assumed to be invariant with respect to the amount of resources provisioned to the activity, and it is guided by the mean makespan indicators of sub-workflows. The method has been successfully tested on a complex workflow with a high degree of concurrency, and applied to the problem of identifying additional resources needed to guarantee a given SLO, proving to be not only effective at improving the workflow e2e response time CDF but also efficient.

Though scaling actions considered in this paper are vertical, the heuristics could be easily extended to perform horizontal scaling actions. In fact, it is sufficient to intend the involved resources as discrete, i.e., as containers or VMs with fixed capacities. In this case, the approach should be adapted so as to round up or down the identified amounts of resources to be allocated. Moreover, the proposed heuristics can be extended to manage workflow blocks that break the structure of well-formed nesting of activities, requiring to compute a makespan indicator and a (sub-optimal) resource assignment for such blocks. The heuristics could also be extended to efficiently derive the value of $\alpha \in [0, 1]$ that minimizes the makespan indicator of each block. Finally, the heuristics could also be improved by considering different performance models, so as to ensure the applicability of the method to contexts in which linearity between response time and amount of allocated resources may not be sufficient to properly characterize the behaviour of the system, e.g., due to the presence of not negligible VM start up times.

References

1. Ali-Eldin, A., Tordsson, J., Elmroth, E.: An adaptive hybrid elasticity controller for cloud infrastructures. In: IEEE Network Operations and Management Symposium. pp. 204–212. IEEE (2012)

2. Alshuqayran, N., Ali, N., Evans, R.: A systematic mapping study in microservice architecture. In: IEEE Int. Conf. on Service-Oriented Computing and Applications. pp. 44–51. IEEE (2016)
3. Bauer, A., Herbst, N., Spinner, S., Ali-Eldin, A., Kounev, S.: Chameleon: A hybrid, proactive auto-scaling mechanism on a level-playing field. IEEE Transactions on Parallel and Distributed Systems **30**(4), 800–813 (2018)
4. Bauer, A., Lesch, V., Versluis, L., Ilyushkin, A., Herbst, N., Kounev, S.: Chamulleon: Coordinated auto-scaling of micro-services. In: IEEE Int. Conf. on Distributed Computing Systems. pp. 2015–2025. IEEE (2019)
5. Berg, B., Dorsman, J.L., Harchol-Balter, M.: Towards optimality in parallel scheduling. Proc. ACM Meas. Anal. Comput. Syst. **1**(2), 40:1–40:30 (2017). <https://doi.org/10.1145/3154499>, <https://doi.org/10.1145/3154499>
6. Bi, J., Zhu, Z., Tian, R., Wang, Q.: Dynamic provisioning modeling for virtualized multi-tier applications in cloud data center. In: IEEE Int. Conf. on Cloud Computing. pp. 370–377. IEEE (2010)
7. Buyya, R., Yeo, C.S., Venugopal, S., Broberg, J., Brandic, I.: Cloud computing and emerging IT platforms: Vision, hype, and reality for delivering computing as the 5th utility. Future Generation Computer Systems **25**(6), 599–616 (2009)
8. Canfora, G., Di Penta, M., Esposito, R., Villani, M.L.: QoS-aware replanning of composite web services. In: IEEE Int. Conf. on Web Ser. pp. 121–129. IEEE (2005)
9. Cardellini, V., Casalicchio, E., Grassi, V., Iannucci, S., Presti, F.L., Mirandola, R.: Moses: A framework for qos driven runtime adaptation of service-oriented systems. IEEE Transactions on Software Engineering **38**(5), 1138–1159 (2011)
10. Carnevali, L., Paolieri, M., Reali, R., Vicario, E.: Compositional safe approximation of response time distribution of complex workflows. In: Proceedings of QEST 2021. vol. 12846, pp. 83–104. Springer (2021)
11. Carnevali, L., Paolieri, M., Reali, R., Vicario, E.: Compositional safe approximation of response time probability density function of complex workflows. ACM Transactions on Modeling and Computer Simulation (2023)
12. Chen, T., Bahsoon, R., Yao, X.: A survey and taxonomy of self-aware and self-adaptive cloud autoscaling systems. arXiv preprint arXiv:1609.03590 (2016)
13. Chieu, T.C., Mohindra, A., Karve, A.A., Segal, A.: Dynamic scaling of web applications in a virtualized cloud computing environment. In: IEEE Int. Conf. on e-Business Engineering. pp. 281–286. IEEE (2009)
14. Farokhi, S., Lakew, E.B., Klein, C., Brandic, I., Elmroth, E.: Coordinating CPU and memory elasticity controllers to meet service response time constraints. In: Int. Conf. on Cloud and Autonomic Computing. pp. 69–80. IEEE (2015)
15. Fox, A., Griffith, R., Joseph, A., Katz, R., Konwinski, A., Lee, G., Patterson, D., Rabkin, A., Stoica, I., et al.: Above the Clouds: A Berkeley View of Cloud Computing. Dept. Electrical Eng. and Comput. Sciences, University of California, Berkeley, Rep. UCB/EECS **28**(13), 2009 (2009)
16. Gias, A.U., Casale, G., Woodside, M.: Atom: Model-driven autoscaling for microservices. In: Int. Conf. on Distributed Comp. Sys. pp. 1994–2004. IEEE (2019)
17. Horváth, A., Paolieri, M., Ridi, L., Vicario, E.: Transient analysis of non-Markovian models using stochastic state classes. Perf. Eval. **69**(7-8), 315–335 (2012)
18. Incerto, E., Tribastone, M., Trubiani, C.: Combined vertical and horizontal autoscaling through model predictive control. In: Int. Conf. on Parallel and Distributed Computing. pp. 147–159. Springer (2018)
19. Iqbal, W., Dailey, M.N., Carrera, D., Janecek, P.: Adaptive resource provisioning for read intensive multi-tier applications in the cloud. Future Generation Computer Systems **27**(6), 871–879 (2011)

20. Lakew, E.B., Klein, C., Hernandez-Rodriguez, F., Elmroth, E.: Towards faster response time models for vertical elasticity. In: IEEE/ACM Int. Conf. on Utility and Cloud Computing. pp. 560–565. IEEE (2014)
21. Liu, J., Zhang, S., Wang, Q., Wei, J.: Coordinating fast concurrency adapting with autoscaling for slo-oriented web applications. IEEE Transactions on Parallel and Distributed Systems **33**(12), 3349–3362 (2022)
22. Lynn, T., Rosati, P., Lejeune, A., Emeakaroha, V.: A preliminary review of enterprise serverless cloud computing (function-as-a-service) platforms. In: IEEE Int. Conf. on Cloud Computing Technology and Science. pp. 162–169. IEEE (2017)
23. Nguyen, H., Shen, Z., Gu, X., Subbiah, S., Wilkes, J.: Agile: Elastic distributed resource scaling for infrastructure-as-a-service (2013)
24. Qiu, H., Banerjee, S.S., Jha, S., Kalbarczyk, Z.T., Iyer, R.K.: Firm: An intelligent fine-grained resource management framework for slo-oriented microservices. In: USENIX Symp. on Operating Systems Design and Implementation (2020)
25. Qu, C., Calheiros, R.N., Buyya, R.: Auto-scaling web applications in clouds: A taxonomy and survey. ACM Computing Surveys **51**(4), 1–33 (2018)
26. Rahman, J., Lama, P.: Predicting the end-to-end tail latency of containerized microservices in the cloud. In: Int. Conf. on Cloud Eng. pp. 200–210. IEEE (2019)
27. Rose, K., Eldridge, S., Chapin, L.: The Internet of Things: An Overview. The internet society (ISOC) **80**, 1–50 (2015)
28. Rosenberg, F., Leitner, P., Michlmayr, A., Celikovic, P., Dustdar, S.: Towards composition as a service-a quality of service driven approach. In: IEEE Int. Conf. on Data Engineering. pp. 1733–1740. IEEE (2009)
29. Russell, N., Ter Hofstede, A.H., Van Der Aalst, W.M., Mulyar, N.: Workflow control-flow patterns: A revised view. BPM Center Report BPM-06-22, BPMcenter.org **2006** (2006)
30. Salah, K., Elbadawi, K., Boutaba, R.: An analytical model for estimating cloud resources of elastic services. J. of Network and Sys. Manag. **24**, 285–308 (2016)
31. Salehie, M., Tahvildari, L.: Self-adaptive software: Landscape and research challenges. ACM Transactions on Autonomous and Adaptive Systems **4**(2), 1–42 (2009)
32. Shahradd, M., Balkind, J., Wentzlaff, D.: Architectural implications of function-as-a-service computing. In: IEEE/ACM Int. Symp. on microarchitecture. pp. 1063–1075 (2019)
33. Shen, Z., Subbiah, S., Gu, X., Wilkes, J.: Cloudscale: elastic resource scaling for multi-tenant cloud systems. In: ACM Symp. on Cloud Computing. pp. 1–14 (2011)
34. Stieß, S., Becker, S., Ege, F., Höppner, S., Tichy, M.: Coordination and explanation of reconfigurations in self-adaptive high-performance systems. In: Int. Conf. on Model Driven Eng. Languages and Systems: Companion Proc. pp. 486–490 (2022)
35. Trivedi, K.S., Sahner, R.: Sharpe at the age of twenty two. ACM SIGMETRICS Performance Evaluation Review **36**(4), 52–57 (2009)
36. Urgaonkar, B., Shenoy, P., Chandra, A., Goyal, P., Wood, T.: Agile dynamic provisioning of multi-tier internet applications. ACM Transactions on Autonomous and Adaptive Systems (TAAS) **3**(1), 1–39 (2008)
37. Vicario, E., Sassoli, L., Carnevali, L.: Using stochastic state classes in quantitative evaluation of dense-time reactive systems. IEEE Transactions on Software Engineering **35**(5), 703–719 (2009)
38. Yazdanov, L., Fetzer, C.: Vertical scaling for prioritized VMs provisioning. In: Int. Conf. on Cloud and Green Computing. pp. 118–125. IEEE (2012)
39. Zheng, Z., Trivedi, K.S., Qiu, K., Xia, R.: Semi-Markov models of composite web services for their performance, reliability and bottlenecks. IEEE Transactions on Services Computing **10**(3), 448–460 (2015)