

# Workflow Scheduling in the Cloud-Edge Continuum

L. Zanussi, D. Tessera, L. Massari, M. Calzarossa

**Abstract** Scheduling in the cloud-edge continuum is a challenging problem. In fact, scheduling has to cope with the peculiarities of these complex ecosystems and satisfy at the same time the desired service levels. In this paper, we investigate the benefits of the cloud-edge continuum for deploying workflows with different characteristics, e.g., computation or communication-intensive. In detail, we formulate a multi-objective optimization problem solved using a Genetic Algorithm. This problem is aimed at identifying the scheduling plans that minimize two conflicting objectives, namely, the expected workflow execution time and monetary cost associated with the cloud and edge resources to be provisioned. Our experiments have shown that the plans that exploit both cloud and edge resources represent a good trade-off between the two objectives. In addition, the workflow characteristics strongly influence these plans. Similarly, the uncertainties that might affect the infrastructure performance are responsible of significant changes in the corresponding Pareto fronts.

## 1 Introduction

Edge computing brings processing capabilities closer to the sources of data and reduces network delays and bandwidth usage. On the contrary, cloud computing offers potentially unlimited processing and storage capabilities at the expense of increased network delays. Therefore, to fully exploit these complex ecosystems, a seamless integration of these technologies is necessary [11].

---

L. Zanussi, L. Massari, M. Calzarossa

Department of Electrical, Computer and Biomedical Engineering, University of Pavia, Italy  
e-mail: luca.zanussi01@universitadipavia.it, luisa.massari@unipv.it, mcc@unipv.it

D. Tessera

Department of Mathematics and Physics, Catholic University of Sacred Heart, Italy  
e-mail: daniele.tessera@unicatt.it

In this paper we investigate the benefits of the cloud-edge continuum for deploying distributed applications characterized by different requirements in terms of computation and communication. In particular, we focus on workflows and on their scheduling. For this purpose, we formulate a multi-objective optimization problem aimed at identifying the “best” scheduling plans that minimize two conflicting objectives, that is, expected workflow execution time and monetary cost. The results of our experiments clearly demonstrate that the combined usage of cloud and edge resources represents a good compromise between the two objectives even under performance variability of the infrastructure.

The rest of the paper is organized as follows. Section 2 briefly discusses the state of the art in the area of cloud-edge scheduling. Section 3 presents the proposed scheduling framework, while Section 4 explains the setup of the experiments and discusses their results. Finally, Section 5 summarizes the paper and outlines future research directions.

## 2 Related work

Scheduling policies for cloud environments have been extensively studied for more than a decade. Many diverse approaches have been presented in the literature. These approaches mainly differ in terms of optimization models, number and types of objectives as well as type of workloads being scheduled (see, e.g., [1, 4, 14, 16] for detailed surveys).

In the context of cloud-edge continuum, scheduling policies have been researched to a more limited extent (see, e.g., [2, 3, 8, 12, 15, 17, 18]). Most papers focus on independent tasks, whereas only few consider more complex workloads consisting of workflows with tasks characterized by precedence constraints. As discussed in [13], genetic-based optimizations are often adopted. For example, Ali et al. [3] propose a task scheduling model based on an extension of the NSGA-II algorithm. Makespan and overall cost are considered as objectives to be minimized. Ijaz et al. [15] focus on workflow scheduling by formulating a multi-objective optimization problem that considers both makespan minimization and reduction in energy consumption. In particular, a weighted bi-objective cost function is introduced for selecting the processing node that minimizes task completion time and energy consumption based on a user-defined weighting factor. A further reduction in energy consumption is obtained by applying deadline constrained frequency scaling.

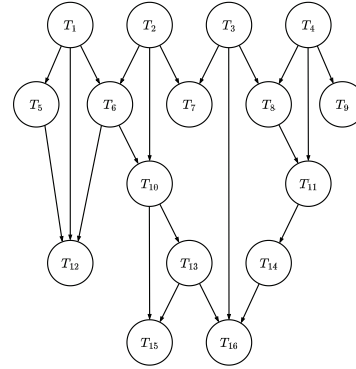
In this paper we study the benefits of the cloud-edge continuum for workflow scheduling by solving a multi-objective optimization problem that takes into account the performance variability that might affect these complex ecosystems.

### 3 Scheduling framework

The proposed scheduling framework is based on the formulation of an optimization problem aimed at minimizing two conflicting objectives, namely, the expected values of the workflow execution time and overall monetary cost under uncertain conditions. For solving this problem, we rely on a Genetic Algorithm applied in combination with a Monte Carlo simulation of all possible solutions. As a result, we obtain multiple “optimal” solutions, i.e., scheduling plans, distributed on the so-called Pareto front. In what follows we introduce the workload and architectural models and we discuss the details of the problem and its solution.

#### 3.1 Workload and architectural models

The workload considered in this study is represented by a workflow  $W$  consisting of  $n$  tasks. This workflow is modeled by a Directed Acyclic Graph (DAG) whose nodes  $T_i$  ( $i = 1, \dots, n$ ) correspond to the tasks and whose edges represent the control and data dependencies between the various tasks (see, e.g., Figure 1).



**Fig. 1** Example of a workflow with 16 nodes.

Each task is described by demands corresponding to the requirements in terms of computation, i.e., processing, communication, i.e., volume of data exchanged with other tasks, input, i.e., data volume transferred from external sources, and output, i.e., data volume transferred to storage devices.

The architectural model of the cloud-edge continuum is represented by multiple instances of physical and logical resources organized according to a layered structure. Some resources are located at the edge/fog layers, whereas some others are inside cloud data centers. Independently of its location, each resource is characterized by processing capacity, transfer rates towards other resources and transfer rates from external sources and towards storage devices.

A resource is also described by its monetary cost as well as the costs for the various types of data transfers.

For each task, we then obtain its computation, communication, input and output times which depend on the resources being allocated. From these values we compute the overall execution time of each task as well as the cost for its deployment. The times and costs of individual tasks together their dependencies are then used to derive the workflow execution time  $T_W$  and its overall cost  $C_W$ .

As discussed in [6, 10], the demands of individual tasks and the characteristics of the infrastructure are often affected by performance variability and uncertainty, thus we will model them as random variables. Hence, both  $T_W$  and  $C_W$  are random variables.

Another important aspect considered in the workload and architectural models deals with the characterization of the workflows in terms of computation and communication. In detail, we define the Communication to Computation Ratio as

$$CCR_W = \frac{\sum_i^n \overline{t_i^{comm}}}{\sum_i^n \overline{t_i^{comp}}}$$

where  $\overline{t_i^{comm}}$  and  $\overline{t_i^{comp}}$  denote the average communication and computation times of task  $T_i$  computed by considering all possible communication links and resources available in the infrastructure, respectively. A value of  $CCR_W$  greater than 1 denotes a communication-intensive workflow, while a value less than 1 denotes a computation-intensive workflow.

### 3.2 Scheduling problem

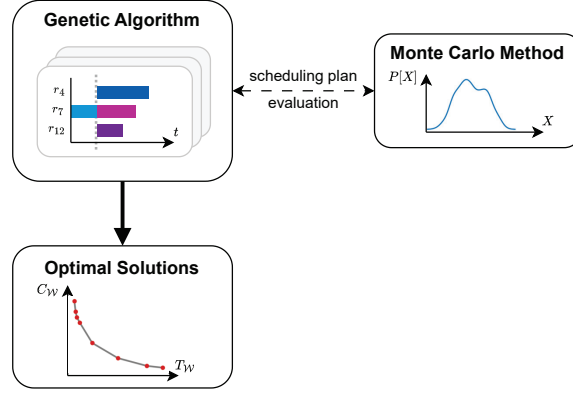
The scheduling problem is formulated as a multi-objective optimization problem as follows:

$$\text{minimize } (\mathbb{E}[T_W], \mathbb{E}[C_W])$$

where  $\mathbb{E}[T_W]$  and  $\mathbb{E}[C_W]$  denote the expected values of the workflow execution time and overall monetary cost. As already mentioned, both  $T_W$  and  $C_W$  are random variables since they are derived from the random variables describing the characteristics of the tasks and of the resources.

The main components of the solution framework are summarized in Figure 2. As can be seen, the solution of the optimization problem is based on the combined application of a meta-heuristic able to cope with multi-objective problems, i.e., the Genetic Algorithm (GA), and the Monte Carlo simulation, that provides the probabilistic evaluation of the objectives. As a result, “optimal” scheduling plans distributed on a Pareto front are obtained.

In detail, starting from the provisioned resources and the scheduling plans generated by the Genetic Algorithm, the Monte Carlo simulation provides their evaluation by properly combining samples of the random variables that describe the tasks and



**Fig. 2** Main components of the proposed solution framework.

the infrastructure. Hence, multiple realizations of the workflow execution time and cost are obtained and used to derive the empirical distributions of the execution time  $T_W$  and of the cost  $C_W$ .

## 4 Experimental results

To assess the benefits of deploying workflows in the cloud-edge continuum, we perform several experiments varying the characteristics of the workflow and the uncertainty affecting the infrastructure performance. In what follows, we describe the workflow and infrastructure used in the experiments as well as the choices made for solving the optimization problem. We also present the results of the experiments and discuss the main findings.

### 4.1 Workload and infrastructure characteristics

For the experiments, we consider the IoT data processing workflow proposed in [17] and displayed in Fig. 1. Each task performs some computation followed by communication with one or more tasks according to the precedence constraints depicted in the figure. On average, each task requires 2.24 Million Instructions and exchanges 5.14 GB of data with its neighboring tasks.

We also outline that the input data of the four entry tasks is generated by external sources associated with IoT devices, while the five exit tasks transfer their output data to storage devices associated with cloud infrastructure. In our experiments, we set the output data to 5GB, while we vary the input data.

The infrastructure used for the experiments consists of multiple instances of different resource types grouped in two layers, namely, four types of resources at the fog/edge layer and eight types at the cloud layer, whose characteristics are summarized in Tables 1 and 2. As can be seen, the resources at the fog/edge layer are

**Table 1** Characteristics of edge devices of the infrastructure considered in the experiments.

	Processing Capacity [MIPS $\times 10^3$ ]	Bandwidth [Mbps]	Latency [ms]	Pricing [USD/hr]	Instance Count
<i>Fog/Edge</i>	10	400	20	0.035	5
	8	100	100	0.020	10
	5	200	50	0.015	20
	2	50	200	0.005	20

**Table 2** Characteristics of cloud devices of the infrastructure considered in the experiments.

	Processing Capacity [MIPS $\times 10^3$ ]	Bandwidth [Mbps]	Pricing [USD/hr]
<i>Cloud A</i>	400	20,000	4.25
	100	4,000	1.00
	50	4,000	0.45
	25	3,000	0.20
<i>Cloud B</i>	240	24,000	2.00
	120	12,000	0.90
	60	6,000	0.40
	30	3,000	0.18

characterized by a limited processing capacity and lower bandwidth, while cloud resources are much more powerful and well connected in terms of bandwidth and latency (whose value is set to 2ms), although at higher costs. Note that the tables present the cost per hour of the resources. Nevertheless, in our experiments we use a per-minute billing, as typically adopted by cloud providers.

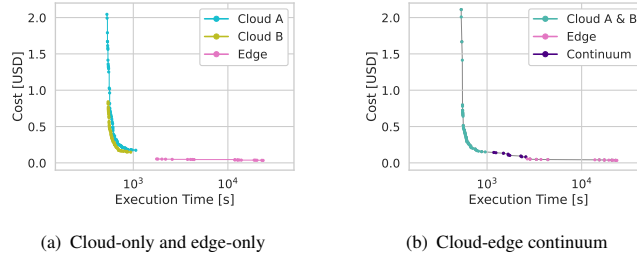
The solution of the multi-objective scheduling problem is based on the NSGA-II algorithm, an elitist non-dominated sorting Genetic Algorithm [9], properly customized as described in [5, 7]. Concerning the GA operators, we choose single point crossover, random mutation and a binary tournament selection method. In addition, we set the initial population to 100 individuals, i.e., scheduling plans, and the probabilities associated with crossover and mutation to 0.9 and 0.02, respectively.

To model the variability affecting the performance of the infrastructure, we build a uniform distribution using the nominal performance as its maximum value and computing its minimum value by subtracting from the nominal performance its value multiplied by a penalty factor corresponding to the desired variability. For example, a 10% variability applied to a processing capacity of 100 MIPS leads to a uniform distribution in the range [90, 100] MIPS, whereas for a 90% variability

we obtain the range  $[10, 100]$  MIPS. In our experiments, we assume both processing capacity and network bandwidth affected by variability.

## 4.2 Results

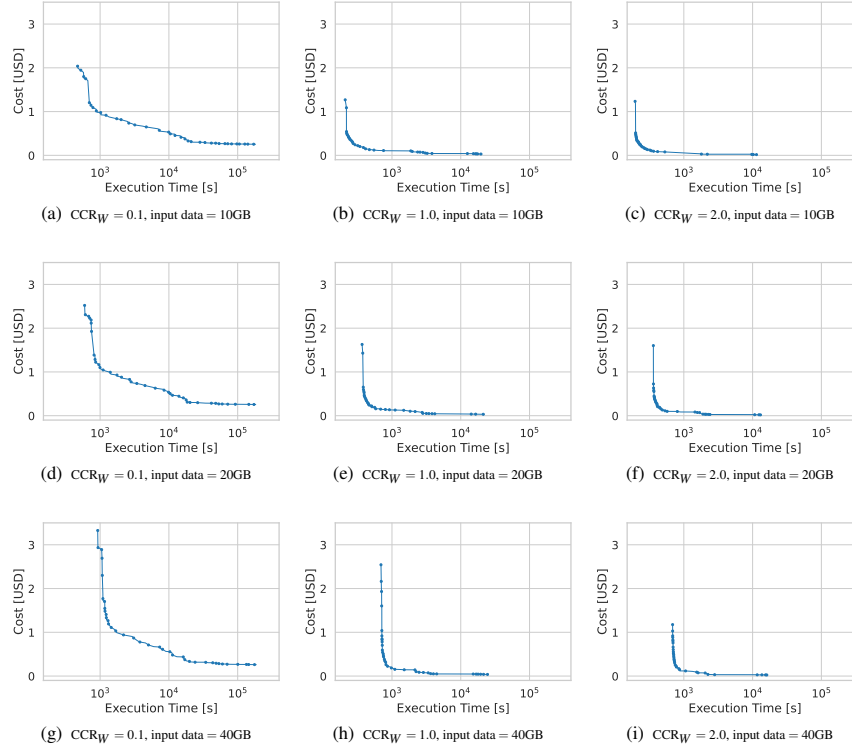
The first set of experiments considers a balanced workflow in terms of communication and computation, that is, with  $CCR_W = 1$ , whose entry tasks receive in total 15GB of input data from external sources. Each of the four experiments exploits different resources, namely, the resources of Cloud A, Cloud B, fog/edge layer and of the entire cloud-edge continuum. The main objective of these experiments is to investigate the impact of these infrastructures on the “optimal” scheduling plans identified by the Genetic Algorithm and the benefits of the cloud-edge continuum. Figure 3 displays the corresponding Pareto fronts. We can observe in Figure 3(a) that



**Fig. 3** Pareto fronts obtained for the balanced workflow, i.e., with  $CCR_W = 1$ , as a function of the resources being provisioned, i.e., belonging to Cloud A, Cloud B or to the fog/edge layer (a) and to the cloud-edge continuum (b). Log scale is used on the  $x$ -axis.

the two cloud-only Pareto fronts and the fog/edge front differ significantly. As expected, the monetary cost of the scheduling plans that exploit Cloud A resources are generally higher than their Cloud B counterparts. Nevertheless, the corresponding execution times only slightly benefit of these expensive resources. On the contrary, for the fog/edge resources, the costs are rather limited, but the execution times are much longer. There is also an evident gap between the cloud fronts and the fog/edge front which suggests the lack of solutions. As shown in Figure 4(b) the cloud-edge continuum is able to fill this gap. In fact, the GA identifies “hybrid” solutions, that is, scheduling plans that exploits the cloud-edge continuum. These plans represent a good compromise between time and cost.

Another set of experiments aims at analyzing the Pareto fronts obtained in the cloud-edge continuum as a function of the workflow characteristics, that is, balanced, computation or communication-intensive, as well as of the input data received from external sources. Figure 4 summarizes the behaviors of the corresponding Pareto fronts. We can easily notice that, for a given input data size, the fronts

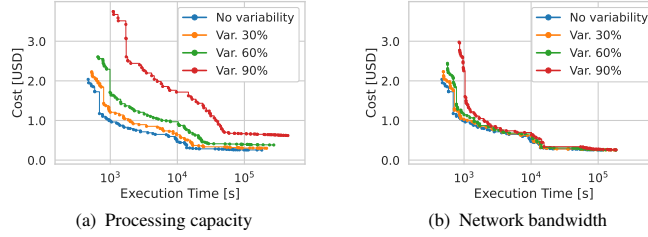


**Fig. 4** Pareto fronts as a function of the value of  $CCR_W$  describing the workflow and the input data size. Log scale is used on the  $x$ -axis.

change with the characteristics of the workflow because of the different resource requirements. For example, the scheduling plans of computation-intensive workflows are generally more expensive and characterized by a longer execution time. Concerning the effects of the data size, the diagrams suggest that, as the size increases, the cost tends to increase.

Finally, we analyze the effects of the variability affecting the infrastructure, namely, processing capacity and network bandwidth. Figure 5 displays the Pareto fronts obtained for the computation-intensive workflow, i.e., with  $CCR_W = 0.1$ , as a function of the variability. As expected, the variability of the processing capacity strongly affects the Pareto fronts, whereas the variability of the network bandwidth has more limited effects.





**Fig. 5** Pareto fronts obtained for the computation-intensive workflow as a function of the variability introduced in the processing capacity (a) and in the network bandwidth (b) of the infrastructure. Log scale is used on the  $x$ -axis.

## 5 Conclusion

Workflow scheduling in the cloud-edge continuum is a very challenging problem especially whenever multiple requirements have to be satisfied simultaneously. In this paper, we identified scheduling plans that minimize the workflow execution time and monetary cost by formulating a multi-objective optimization problem solved through a combined application of the Genetic Algorithm and the Monte Carlo simulation. Our experiments have demonstrated that workflows greatly benefit of the cloud-edge continuum especially in presence of large inputs from external data sources, such as IoT devices.

As future research directions, we plan to further study the scheduling problem in the cloud-edge continuum by introducing different types of variability and considering various types of constraints, e.g., related to data privacy and confidentiality.

## Acknowledgments

This work was partly supported by the Italian Ministry of University and Research (MUR) under the PRIN 2022 grant “Methodologies for the Parallelization, Performance Evaluation and Scheduling of Applications for the Cloud-Edge Continuum” (Master CUP: B53D23013090006, CUP: J53D23007110008, CUP: F53D23004300006) and by the European Union - Next Generation EU.

## References

1. Adhikari, M., Amgoth, T., Srirama, S.N.: A Survey on Scheduling Strategies for Workflows in Cloud Environment and Emerging Trends. *ACM Comp Surv.* **52**(4) (2019)
2. Agarwal, G., Gupta, S., Ahuja, R., Rai, A.: Multiprocessor task scheduling using multi-objective hybrid Genetic Algorithm in Fog-cloud computing. *Knowledge-Based Systems*

3. Ali, I., Sallam, K., Moustafa, N., Chakraborty, R., Ryan, M., Choo, K.K.R.: An Automated Task Scheduling Model Using Non-Dominated Sorting Genetic Algorithm II for Fog-Cloud Systems. *IEEE Trans. on Cloud Computing* **10**(4), 2294–2308 (2022)
4. Arunarani, A., Manjula, D., Sugumaran, V.: Task scheduling techniques in cloud computing: A literature survey. *Future Generation Computer Systems* **91**, 407–415 (2019)
5. Calzarossa, M.C., Della Vedova, M.L., Massari, L., Nebbione, G., Tessera, D.: Multi-objective optimization of deadline and budget-aware workflow scheduling in uncertain clouds. *IEEE Access* **9** (2021)
6. Calzarossa, M.C., Della Vedova, M.L., Tessera, D.: A methodological framework for cloud resource provisioning and scheduling of data parallel applications under uncertainty. *Future Generation Computer Systems* **93**, 212–223 (2019)
7. Calzarossa, M.C., Massari, L., Nebbione, G., Della Vedova, M.L., Tessera, D.: Tuning Genetic Algorithms for Resource Provisioning and Scheduling in Uncertain Cloud Environments: Challenges and Findings. In: *Proc. of the 27th Euromicro International Conference on Parallel, Distributed and Network-Based Processing (PDP)*, pp. 174–180 (2019)
8. De Maio, V., Kimovski, D.: Multi-objective scheduling of extreme data scientific workflows in Fog. *Future Generation Computer Systems* **106**, 171–184 (2020)
9. Deb, K., Pratap, A., Agarwal, S., Meyarivan, T.: A Fast and Elitist Multiobjective Genetic Algorithm: NSGA-II. *IEEE Trans. on Evolutionary Computation* **6**(2), 182–197 (2002)
10. Della Vedova, M.L., Tessera, D., Calzarossa, M.C.: Probabilistic Provisioning and Scheduling in Uncertain Cloud Environments. In: *Proc. of the 2016 IEEE Symposium on Computers and Communication - ISCC*, pp. 797–803 (2016)
11. Esposito, A., Aversa, R., Barbierato, E., Calzarossa, M., Di Martino, B., Massari, L., Mongiardo, I., Tessera, D., Venticinque, S., Zanussi, L., Zieni, R.: Methodologies for the Parallelization, Performance Evaluation and Scheduling of Applications for the Cloud-Edge Continuum. In: L. Barolli (ed.) *Advanced Information Networking and Applications - AINA*. Springer (2024)
12. Goudarzi, M., Wu, H., Palaniswami, M., Buyya, R.: An Application Placement Technique for Concurrent IoT Applications in Edge and Fog Computing Environments. *IEEE Trans. on Mobile Computing* **20**(4), 1298–1311 (2021)
13. Guerrero, C., Lera, I., Juiz, C.: Genetic-based optimization in fog computing: Current trends and research opportunities. *Swarm and Evolutionary Computation* **72** (2022)
14. Hosseinzadeh, M., Ghafour, M.Y., Hama, H.K., Vo, B., Khoshnevis, A.: Multi-objective task and workflow scheduling approaches in cloud computing: a comprehensive review. *Journal of Grid Computing* **18**, 327–356 (2020)
15. Ijaz, S., Munir, E., Ahmad, S., Rafique, M., Rana, O.: Energy-Makespan Optimization of Workflow Scheduling in Fog-Cloud Computing. *Computing* **103**(9), 2033–2059 (2021)
16. Masdari, M., ValiKardan, S., Shahi, Z., Azar, S.: Towards workflow scheduling in cloud computing: A comprehensive analysis. *Journal of Network and Computer Applications* **66**, 64–82 (2016)
17. Stavrinides, G.L., Karatza, H.D.: A hybrid approach to scheduling real-time IoT workflows in fog and cloud environments. *Multimedia Tools and Applications* **78** (2019)
18. Sun, Y., Lin, F., Xu, H.: Multi-objective optimization of resource scheduling in fog computing using an improved NSGA-II. *Wireless Personal Communications* **102**, 1369–1385 (2018)