# GRASP-based resource re-optimization for effective big data access in federated clouds

CrossMark

Francesco Palmieri [a],*, Ugo Fiore [b], Sergio Ricciardi [c], Aniello Castiglione [d]

[a] Department of Industrial and Information Engineering, Second University of Naples, Aversa (CE), Italy
[b] Information Services Centre, Federico II University, Napoli, Italy
[c] Department of Computer Architecture, Technical University of Catalonia, Barcelona, Spain
[d] Department of Computer Science, University of Salerno, Fisciano (SA), Italy

## HIGHLIGHTS

- Efficient re-optimization strategy for big data access in multi-tenant cloud infrastructures.
- Based on a Greedy Randomized Adaptive Search meta-heuristic working on a flexible federated cloud model.
- Performs end-to-end path rerouting and Virtual Machine migration to improve access to big data.
- Rebalances cloud resource usage so that more virtual machines can effectively access data sources.

## ARTICLE INFO

## ABSTRACT

Federated cloud organizations, spanning across multiple networked sites that provide both computing and storage resources, can be considered the state-of-the-art solutions for providing multi-tenant runtime services in modern distributed processing environments. In these scenarios, by re-optimizing the communication paths between virtual machines and big data sources, at evenly spaced interval or when required by circumstances, the overall communication and runtime resource utilization on the cloud infrastructure is re-balanced, so that more virtual machines can be allowed to access the needed big data sources with adequate bandwidth, thereby significantly improving the perceived performance and quality of service. The problem of re-optimization is tackled with a powerful meta-heuristic, the *greedy randomized adaptive search procedure* (GRASP), augmented by path re-linking. In order to evaluate the proposed approach, extensive simulations have been performed, leading to very interesting results, demonstrating the effectiveness and validity of the underlying ideas and their applicability to real large-scale federated cloud scenarios.

© 2015 Elsevier B.V. All rights reserved.

## 1. Introduction

Nowadays the proliferation of the data sources available on the Internet and the widespread deployment of network-based applications are fostering the emergence of new architectures referred as "big data", characterized by the need of capturing, combining and processing an always growing amount of heterogeneous and unstructured data coming from new mobile devices, emerging social media and human/machine-to-machine communication. The distributed nature of these resources implies remotely accessing and moving data at volumes and rates that push the frontiers of current networking technologies, so that the Internet, due to its known scalability limitations, becomes the major bottleneck for these network-based data-intensive applications. This results into an abrupt shift from the classic *application-centric* paradigm, characterized by static applications triggering on-demand data transfers towards their sites, to a novel *data-centric* model, in which applications themselves are dynamically moved through the network in order to be run in the most convenient places, where enough communication capacity is available to provide effective access to their data.

Clouds are an excellent accelerator for such a shift because they make the deployment of new large-scale data warehousing architectures, supporting data processing or storage activities,

* Corresponding author.
*E-mail addresses:* francesco.palmieri@unina.it (F. Palmieri), ugo.fiore@unina.it (U. Fiore), sergio.ricciardi@ac.upc.edu (S. Ricciardi), castiglione@ieee.org (A. Castiglione).

extremely easy and affordable. Accordingly, most of the services based on big data processing are currently delivered in cloud style, relying on the orchestrated usage of geographically sparse data centers to satisfy their huge run-time and storage needs. Thus, cloud providers are starting to offer virtual data center services through the cloud, according to the *Infrastructure as a Service (IaaS)* model, where such a virtualized infrastructure is built on a set of federated individual tenants scattered throughout the world and interconnected via virtual end-to-end connections realized over the Internet. In order to achieve this, cloud providers use the virtualization technology, providing the abstraction of computing, storage, and network resources from their physical counterparts (according to the server, storage, and network virtualization paradigms) so that:

- the physical server machines located in the federated data centers can run multiple *Virtual Machines (VMs)* capable of hosting their own operating systems and applications as well as even acting as virtual routers or switches;
- the storage resources available on the involved sites are managed in a totally flexible way by aggregating or partitioning them into virtual storage units or blocks providing the storage space needed by VMs. Such space can be dynamically extended or shrunk based on the changing VMs requirements.

Server and storage virtualization provides enormous advantages like elastic resource management, potentially unlimited scalability, reduced power usage, increased security/reliability and lower user downtimes. In particular, VMs can connect to big data repositories providing a virtually unlimited storage space, growing in an on-demand fashion across multiple data centers/sites and multiple physical storage systems, according to an horizontal scaling paradigm. However, implementing such a fully virtualized architecture requires the availability of enough transmission bandwidth between the sites hosting the VMs and those ones providing the virtual storage blocks they access, in order to avoid that network communication becomes a bottleneck, making access to remote big data repositories unpractical and hence invalidating the whole architectural framework. This implies setting up (or tiering down) on demand dedicated end-to-end traffic engineering interconnections between the involved sites, providing a certain amount of guaranteed bandwidth between the accessing VMs and the remote storage repositories. Such task can be easily accomplished by relying on the *Generalized Multi-Protocol Label Switching technology (GMPLS)* technology, supporting the creation of bandwidth guaranteed *Label Switched Paths (LSPs)* upon a mesh of IP-based and optical network infrastructures, that is the basic composition of the modern communication facilities such as the global Internet as well as dedicated high-performance transport backbones/links or a hybrid mix of both the aforementioned alternatives. The creation of such traffic-engineered paths is performed at the federated cloud middleware level under the control of the resource brokering/scheduling system. However, the dynamic on-demand creation or deletion of these LSPs has the drawback of unbalancing resource usage in the long run, introducing unwanted congestion effects on several network sections. Typically this will happen on the links which are most often requested, i.e., those providing the best connections towards frequently accessed big data repositories. Precisely, the paths for new data access requests are calculated one by one, according to a bandwidth-constrained route selection criterion, from the actual status of network resource usage which includes the existing LSPs. As the access pattern evolves over time, such solutions clearly become sub-optimal since the evolution process may lead to significant changes in the aforementioned status, resulting from VMs' varying demands for access to different repositories. This may cause bandwidth unavailability towards some repositories and/or run-time resource exhaustion on

some data centers so that no more accesses are allowed due to lack of network communication resources or run-time space, while a more efficient redistribution of LSPs and VMs would have allowed satisfaction of all the access requests.

Optimal behavior may be restored only by introducing a periodic offline re-optimization process in charge of rerouting some already set up LSPs over alternative paths, or moving VMs on different sites/data centers, by reclaiming the lost capacity and also achieving a load re-balance on the overall cloud infrastructure. However, re-optimization carries an unavoidable cost with it. Indeed, both LSP rerouting and live VM migration affect the service continuity and increase the network management overhead. Also, live VM migration across network boundaries implies an additional re-optimization effort due to the need of establishing a set of new LSPs towards the data repositories accessed by the involved VMs, and contextually tearing down the previously existing ones (according to a *make-before-break* mechanism), with the obvious consequence in terms of complexity on the overall resource management problem. Thus, the offline re-optimization activity has to be handled carefully, by using flexible adaptive and effective strategies, to be used only when necessary to recover stranded capacity, by maximizing the positive re-balancing effects on the resources usage.

Starting from the above considerations, we propose an efficient re-optimization strategy for LSP rerouting and VM migration based on a greedy randomized adaptive search meta-heuristic procedure (GRASP) [1], in conjunction with a final solution refinement procedure [2–5]. The whole re-optimization approach, essentially operating within the cloud resource management framework, requires a certain degree of network visibility at the cloud middleware layer and is based on a totally flexible federated cloud model, encompassing heterogeneous run-time, storage and communication resources, in which the capacity and computing power can be independently specified for each communication link and data center site, respectively. We evaluated the effectiveness of the proposal by simulating the whole re-optimization framework on realistic cloud topologies and big data access demands, and measuring the available communication bandwidth as well as the amount of access requests towards specific data repositories that the cloud is able to support before and after re-optimization. The results from these experiments showed that this idea may lead to significant improvements in cloud resource management by acquiring the needed data from the right places, at any time, and by using the right paths throughout the underlying communication networks. Obviously, the perceivable effects on the overall cloud economy grow with both the scale of the federated cloud and its load, and consequently with the amount of data to be produced and processed.

The remainder of this work is organized as follows: after a survey of relevant related work (Section 2), Section 3 quickly summarizes the fundamental points about the characteristics of data access traffic, as well as about meta-heuristics and GRASP. Following that, the system model is stated in Section 4 and the proposed strategy and its assumptions and constraints are discussed in Section 5. Section 6 describes the simulations performed and discusses the results. Section 7 presents the conclusions and some implementation perspectives that may lead to further performance improvements.

## 2. Related work

Resource management in distributed infrastructures is a general topic that has attracted much attention in recent years (see, for example, the survey in [6] or the contributions reported in [7] for more specific cloud-related issues). With the exponential growth of data volumes to be managed, leading to the well-known