# Chapter 40
# Energy–Saving QoS Resource Management of Virtualized Networked Data Centers for Big Data Stream Computing

**Nicola Cordeschi**
*"Sapienza" University of Rome, Italy*

**Danilo Amendola**
*"Sapienza" University of Rome, Italy*

**Mohammad Shojafar**
*"Sapienza" University of Rome, Italy*

**Enzo Baccarelli**
*"Sapienza" University of Rome, Italy*

## ABSTRACT

*In this chapter, the authors develop the scheduler which optimizes the energy-vs.-performance trade-off in Software-as-a-Service (SaaS) Virtualized Networked Data Centers (VNetDCs) that support real-time Big Data Stream Computing (BDSC) services. The objective is to minimize the communication-plus-computing energy which is wasted by processing streams of Big Data under hard real-time constrains on the per-job computing-plus-communication delays. In order to deal with the inherently nonconvex nature of the resulting resource management optimization problem, the authors develop a solving approach that leads to the lossless decomposition of the afforded problem into the cascade of two simpler sub-problems. The resulting optimal scheduler is amenable of scalable and distributed adaptive implementation. The performance of a Xen-based prototype of the scheduler is tested under several Big Data workload traces and compared with the corresponding ones of some state-of-the-art static and sequential schedulers.*

## 1. INTRODUCTION

Energy-saving computing through Virtualized Networked Data Centers (VNetDCs) is an emerging paradigm that aims at performing the adaptive energy management of virtualized Software-as-a-Service (SaaS) computing platforms. The goal is to provide QoS Internet services to large populations of clients, while minimizing the overall computing-plus-networking energy consumption (Cugola & Margara, 2012; Baliga, Ayre, Hinton, & Tucker, 2011; Mishra, Jain, & Durresi, 2012). As recently pointed out

in (Mishra et al. 2012; Azodomolky, Wieder, & Yahyapour, 2013; Wang et al. 2014), the energy cost of communication gear for current data centers may represent a large fraction of the overall system cost and it is induced primarily by switches, LAN infrastructures, routers and load balancers.

However, actual virtualized data centers subsume the (usual) Map/Reduce-like batch processing paradigm and they are not designed for supporting networking-computing intensive real-time services, such as, for example, emerging Big Data Stream Computing (BDSC) services (Cugola et al. 2012). In fact, BDSC services retain the following (somewhat novel and unique) characteristics (Cugola et al. 2012; Scheneider, Hirzel, & Gedik, 2013; Qian, He, Su, Wu, Zhu, & Zhang, 2013; Kumbhare, 2014):

1. The incoming data (i.e., the offered workload) arrive continuously at volumes that far exceed the storage capabilities of individual computing machines. Furthermore, all data must be timely proceed but, typically, a few of data require to be stored. This means that the (usual) storing-then-computing batch paradigm is no longer feasible;
2. Since BDSC services acquire data from massive collections of distributed clients in a stream form, the size of each job is typically unpredictable and also its statistics may be quickly time-varying; and,
3. The offered workload is a real-time data stream, which needs real-time computing with latencies firmly limited up to a few of seconds (Qian et al. 2013; Kumbhare, 2014). Imposing hard limits on the overall per-job delay requires, in turn, that the overall VNetDC is capable to quickly adapt its resource allocation to the current (a priori unpredictable) size of the incoming big data.

In order to attain energy saving in such kind harsh computing scenario, the joint balanced provisioning and adaptive scaling of the networking-plus-computing resources is demanded. This is the focus of this work, whose main contributions may be so summarized. First, the contrasting objectives of low consumptions of both networking and computing energies in delay and bandwidth-constrained VNetDCs are cast in the form of a suitable constrained optimization problem, namely, the Computing and Communication Optimization Problem (CCOP). Second, due to the nonlinear behavior of the rate-vs.-power-vs.-delay relationship, the CCOP is not a convex optimization problem and neither guaranteed-convergence adaptive algorithms nor closed-form formulas are, to date, available for its solution. Hence, in order to solve the CCOP in exact and closed-form, we prove that it admits a loss-free (e.g., optimality preserving) decomposition into two simpler loosely coupled sub-problems, namely, the CoMmunication Optimization Problem (CMOP) and the ComPuting Optimization Problem (CPOP). Third, we develop a fully adaptive version of the proposed resource scheduler that is capable to quickly adapt to the a priori unknown time-variations of the workload offered by the supported Big Data Stream application and converges to the optimal resource allocation without the need to be restarted.

## 1.1 RELATED WORK

Updated surveys of the current technologies and open communication challenges about energy-efficient data centers have been recently presented in (Mishra et al. 2012; Balter, 2013). Specifically, power management schemes that exploit Dynamic Voltage and Frequency Scaling (DVFS) techniques for performing resource provisioning are the focus of (Chen & Kuo, 2005; Kim, Buyya, & Kim, 2007; Li, 2008). Although these contributions consider hard deadline constraints, they do not consider, indeed, the

## Related Content

A Novel Multi-Scale Feature Fusion Method for Region Proposal Network in Fast Object Detection
Gang Liuand Chuyi Wang (2020). *International Journal of Data Warehousing and Mining (pp. 132-145).*
www.irma-international.org/article/a-novel-multi-scale-feature-fusion-method-for-region-proposal-network-in-fast-object-detection/256166

Artificial Immune Systems: Using the Immune System as Inspiration for Data Mining
Jonathan Timmisand Thomas Knight (2002). *Data Mining: A Heuristic Approach (pp. 209-230).*
www.irma-international.org/chapter/artificial-immune-systems/7591

TODE: An Ontology-Based Model for the Dynamic Population of Web Directories
Sofia Stamou, Alexandros Ntoulasand Dimitris Christodoulakis (2008). *Data Mining with Ontologies: Implementations, Findings, and Frameworks (pp. 1-17).*
www.irma-international.org/chapter/tode-ontology-based-model-dynamic/7569

Early Warning System for SMEs as a Financial Risk Detector
Ali Serhan Koyuncugil (2009). *Data Mining Applications for Empowering Knowledge Societies (pp. 220-238).*
www.irma-international.org/chapter/early-warning-system-smes-financial/7554

Boat Detection in Marina Using Time-Delay Analysis and Deep Learning
Romane Scherrer, Erwan Aulnette, Thomas Quiniou, Joël Kasarherou, Pierre Kolband Nazha Selmaoui-Folcher (2022). *International Journal of Data Warehousing and Mining (pp. 1-16).*
www.irma-international.org/article/boat-detection-in-marina-using-time-delay-analysis-and-deep-learning/298006