

# Dynamic Resource Management Across Cloud-Edge Resources for Performance-Sensitive Applications

Shashank Shekhar  
Vanderbilt University  
Nashville, TN 37235, USA  
Email: shashank.shekhar@vanderbilt.edu

Aniruddha Gokhale (Adviser)  
Vanderbilt University  
Nashville, TN 37235, USA  
Email: agokhale@vanderbilt.edu

**Abstract**—A large number of modern applications and systems are cloud-hosted, however, limitations in performance assurances from the cloud, and the longer and often unpredictable end-to-end network latencies between the end user and the cloud can be detrimental to the response time requirements of the applications, specifically those that have stringent Quality of Service (QoS) requirements. Although edge resources, such as cloudlets, may alleviate some of the latency concerns, there is a general lack of mechanisms that can dynamically manage resources across the cloud-edge spectrum. To address these gaps, this research proposes Dynamic Data Driven Cloud and Edge Systems ( $D^3CES$ ). It uses measurement data collected from adaptively instrumenting the cloud and edge resources to learn and enhance models of the distributed resource pool. In turn, the framework uses the learned models in a feedback loop to make effective resource management decisions to host applications and deliver their QoS properties.  $D^3CES$  is being evaluated in the context of a variety of cyber physical systems, such as smart city, online games, and augmented reality applications.

**Keywords**—Cloud Computing, Edge Computing, Fog Computing, Resource Management, DDDAS, CPS, IoT.

## I. INTRODUCTION

The elastic properties and cost benefits of the cloud has made it an attractive hosting platform for a variety of soft real-time cyber physical systems (CPS)/Internet of Things (IoT) applications, such as cognitive assistance, patient health monitoring and industrial automation. The stringent quality of service (QoS) considerations of these applications mandate both predictable performance from the cloud and lower end-to-end network latencies between the end user and the cloud. To date, security and performance assurance continues to be a hard problem to resolve in cloud platforms due to their virtualized and multi-tenant nature [12]. Although recent advances in fog and edge computing have enabled cloud resources to move closer to the CPS/IoT devices thereby mitigating the network latency concerns to some extent [3], there is still a general lack of mechanisms that can dynamically manage resources across the cloud-edge spectrum. This is a hard problem to resolve due to the highly dynamic behaviors of the edge and cloud. Consequently, any pre-defined and fixed set resource management policies will be rendered useless for hosting CPS/IoT applications in the cloud.

The dynamic data driven application systems (DDDAS) paradigm [10] addresses precisely these challenges. DDDAS prescribes an approach where applications are instrumented

adaptively so that their models can be learned and enhanced continuously, and in turn these models can be analyzed and used in a feedback loop to steer the applications along their intended trajectories. Previous work have focused on a specific application or applied DDDAS for resilience and security [2]. We propose to apply the DDDAS principle to the pool of resources spanning the cloud-edge spectrum to enable and enforce dynamic resource management decisions that deliver the required QoS properties of cloud-hosted applications. To that end we propose Dynamic Data Driven Cloud and Edge Systems ( $D^3CES$ ), which uses performance data collected from adaptively instrumenting the cloud and edge resources to learn and enhance models of the distributed resource pool, and in turn using these models in a feedback loop to make effective resource management decisions to host CPS applications and deliver their QoS properties.

Rest of this paper is organized as follows, Section II discusses the challenges faced in realizing  $D^3CES$ , Section III discuss the current state of the art, Section IV proposes a set of solutions that form the contours of this doctoral research and Section V provides the concluding remarks.

## II. KEY RESEARCH CHALLENGES AND SOLUTION NEEDS

Our research calls for an effective use of resources across the cloud data centers (CDCs) and the micro data centers (MDCs) that reside at the edge. The following sub-sections lists a non-exhaustive set of challenges along three dimensions that we are addressing in this doctoral research.

### A. Application-imposed Challenges

- 1) **Workload variations:** The workload generated by CPS/IoT applications may illustrate both transient and sustained variability which needs to be predicted and addressed.
- 2) **Stochastic execution semantics:** For some CPS/IoT applications, their uncertain and dynamic nature may require several instances of the same tasks to be executed to reach specified confidence levels. Each execution may take different execution times but impose certain QoS needs.
- 3) **Application structure:** Increasingly, cloud-based applications are realized as a collection of communicating microservices, which can be deployed independently

across the spectrum of resources. This gives rise to interesting challenges in whether part or entire service must be migrated closer to the edge.

- 4) **Reconciling application state:** When a cloud-hosted application is migrated to a MDC, often not all of its state may be transferred to the MDC and hence may have to be reconciled periodically with the state maintained at the CDC, which gives rise to interesting consistency versus availability tradeoffs.
- 5) **High degree of user mobility:** CPS/IoT systems, such as autonomous transport vehicles, unmanned aerial vehicles, and mobile devices, operate in a highly uncertain environments with dynamic movement profiles. Thus, a designated edge resource cannot serve such users for long durations of times.

### B. Cloud Provider-related Challenges

- 1) **Effective utilization of edge resources:** Although exploiting edge resources is an intuitive solution to addressing the network latency issues, the MDCs will also face the same challenges as a CDC, which stem from virtualization and multi-tenancy resulting in application performance interference [16].
- 2) **Workload consolidation and migration across MDCs:** Since the edge may comprise multiple MDCs, there is a need for effective and dynamic server workload consolidation across MDCs and CDCs.
- 3) **Distributed user base:** Collaborative applications such as online games may often involve a distributed set of users. Consequently, determining the MDC to migrate the application to and whether to migrate it to multiple MDCs remains an open question.
- 4) **Shared micro data centers:** In the simplest case, an edge-based MDC may be considered to be owned by the same provider that owns a CDC. In general, however, a MDC could be shared across different CDC providers. Assuring security and isolation guarantees in these scenarios is an open question.
- 5) **Energy savings and revenue generation:** In making use of the spectrum of resources across the cloud and the edge, a cloud provider will be concerned about maximizing revenues and conserving energy while ensuring that application SLOs are met.

### C. Measurement-related Challenge

- 1) **Collecting metrics under hardware heterogeneity:** The plethora of deployed hardware configurations with different architectures and versions makes it hard to collect various performance metrics. Modern architectures are making it easier to collect more finer grained performance metrics, however, much more research is needed in identifying effective approaches to control the hardware and derive the best performance out of them.
- 2) **Lack of benchmarks:** There is a general lack of open source and effective benchmarking suites that researchers can use to conduct studies and build models

of the cloud-edge spectrum of resources that then can be used in resource management.

## III. RELATED RESEARCH

Representative related efforts are outlined along the dimensions of challenges from Section II.

**Related Efforts along Application and Cloud Provider Dimensions:** The mJADES [17] project is founded on a Java-based architecture and is designed to run multiple concurrent simulations while automatically acquiring resources from an ad-hoc federation of cloud providers. However, this framework relies on heavy-weight virtualization whose setup and tear down cost is too high for short duration simulations. DEXSim [6] is a distributed execution framework for replicated simulations that provides two-level parallelism, i.e., at CPU core-level and at system-level. Nonetheless, the solution is tied to the HLA framework and is not generic.

Hidra [18] applies control techniques for fluctuating workloads and constrained resources to meet QoS for distributed real-time and embedded systems. However, its efficacy on modern cloud systems has not been tested. MAUI [9] and CloneCloud [7] are mechanisms to offload tasks from resource-constrained systems such as mobile devices to cloud systems with the purpose of saving battery and increasing speed of execution, but does not focus on QoS needs of latency sensitive applications. Fesehaye et al. [14] proposed a design for selecting between edge and cloud resources for interactive mobile cloud applications based on the number of hops, mobility and latency but use only simulations for validation.

**Related Efforts along the Measurements Dimension:** One of the preliminary task for designing adaptive resource management solutions is to develop a mechanism to discover the performance needs of different CPS/IoT applications in the cloud. This objective requires a benchmarking tool to gather the performance requirements. There are a number of popular benchmarking solutions for the cloud. This includes CloudSuite [13] that is a collection of contemporary applications that run in a cloud and encompasses web serving, data serving and analytics, and media streaming. YCSB [8] and BigDataBench [21] are popular benchmarks for cloud deployment of big data applications. However, these benchmarks do not cover the CPS/IoT scenarios that we are interested in. IoTABench [1] is a benchmark for analyzing Internet of Things (IoT) composed of synthetic load generator, benchmark queries and a test harness. This benchmark is presently composed of only one use case and its extensibility is yet to be observed.

## IV. RESEARCH CONTRIBUTIONS AND ONGOING RESEARCH

Based on the DDDAS paradigm, this doctoral research is realizing the  $D^3CES$  framework to address a range of challenges outlined in Section II. We describe our accomplishments and remaining tasks.

**Contribution 1: Benchmarking CPS/IoT applications and performance monitoring of cloud-edge resources:** Any systems research requires a solid framework to conduct empirical studies. To address Challenges II-C1 and II-C2, we have developed a benchmark framework that is available in open source at <https://github.com/shekharshank/indices>. The benchmark gathers both system and micro-architectural performance metrics while varying application workload and collocation patterns. We also plan to leverage micro-benchmarks for cloud data centers such as ibench [11] for faster and accurate profiling of applications.

**Contribution 2: Algorithms for elastic and scalable scheduling of CPS/IoT tasks in the cloud:** To address Challenges II-A1 and II-A2, where the CPS/IoT applications may need to execute large number of task instances (e.g., stochastic simulations) in the cloud environment with QoS requirements, we need feedback-based algorithms that provide the QoS guarantees while scaling across multiple servers.

To that end, we have designed the *Simulation as a Service (SIMaaS)* cloud middleware [19] that is based on the Linux container-based infrastructure as shown in Figure 1. The contributions include an admission control and a resource management algorithm that reduce the cost to the service provider and enhances the parallelization of the simulation jobs by fanning out more instances until the deadline of the tasks is met while simultaneously auto-tuning itself based on the feedback. Presently, we only have a centralized architecture. We are designing a distributed version of our middleware framework that can exploit the distributed cloud-edge resource pool.

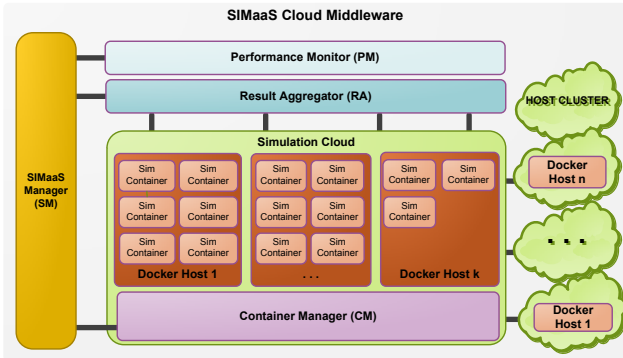


Fig. 1. System Architecture

**Contribution 3: Exploiting edge resources to move applications closer to the users and the data sources:** To address challenge II-B1, the cloud-based service must move closer to the users and data sources. This needs distributed coordination, synchronization and resource allocation challenges to be investigated in the context of CDCs and MDCs, and even utilizing spare capacity of other edge resources. This task also entails discovery and orchestration of those resources not provided out-of-the-box by the traditional cloud.

Our contributions to date have addressed the needs of

applications such as augmented reality that offload computer vision algorithms including SIFT [15] for processing at the cloud. For our current work we have assumed near constant workload and low mobility of the users, which is the case for image processing performed by a stationary camera. The user continuously sends data in the form of images to the cloud for processing and in turn receives responses within a specified time bound. Nonetheless, due to geographical distance and network variabilities, a central cloud may be unable to meet the QoS needs, requiring MDCs closer to the user to be leveraged. As highlighted in II-B1, performance interference caused by co-located applications extends itself from central cloud to edge and leads to delayed response time at the edge resources too. Our initial work [5] was targeted at the centralized cloud and did not address hardware heterogeneity related issues.

To address these challenges, we have formulated an optimization problem to minimize the cost to the cloud provider while meeting the QoS constraints imposed by the application. As part of our framework called INDICES [20], we perform the optimization at two different layers, local, i.e. at a MDC, and global, i.e. at a CDC. Figure 2 depicts the local cost estimation module where the performance of virtualized servers are monitored and used to calculate the cost of deployment of applications for decision making.

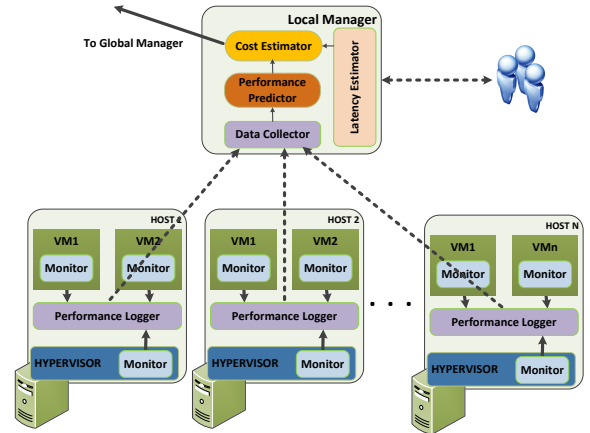


Fig. 2. INDICES Cost Estimation

**Planned Research and Timeline to Completion:** Several other challenges from Section II remain to be resolved. As part of our ongoing work, we are expanding our framework to consider variable workloads (Challenge II-B1) and user mobility (Challenge II-B5). This requires building a profile of users and CPS/IoT applications to forecast the load and expected network latency and bandwidth, and employ efficient resource management algorithms. Another key research direction relates to on-demand workload consolidation and service migration, instead of only an initial deployment framework (Challenges II-A4 and II-B2).

We are incorporating serverless computing and micro-services as part of our research. These micro services may

have dependencies on each other and their states have to be managed while distributing and migrating them across the central cloud and edge resources for optimal performance (Challenges II-A3 and II-A4). We intend to develop repositories with dynamic structure that can maintain up-to-date resource information. We also plan to leverage CPS/IoT devices themselves for ad-hoc edge cloud formation and deployment of micro-services. The micro services can be packaged as self-contained deployable units using Linux containers such as Docker or Unikernels to address the heterogeneity and orchestration issues. As part of the research, we will develop algorithms to perform global optimization to answer *if* and *when* should we migrate the services, identify the nearest edge cloud and how to do this efficiently.

We are also working on integrating our micro benchmark framework for accurate performance estimation in these newer contexts (Challenges II-C1 and II-C2). Notably, we plan to accommodate latest advances in server virtualization, such as cache allocation technology (CAT) [4] for mitigating performance interference due to co-location of applications.

We expect around a year's worth of research investigations to completing the planned research tasks and defending the dissertation during the Spring semester 2018.

## V. CONCLUSION

In this work, we identified the key challenges that inhibit the universal adoption of cloud, especially in the context of CPS/IoT applications. We highlighted the shortcomings in the state of the art and proposed  $D^3CES$  as a framework for adaptive resource management across cloud and edge resources to provide QoS guarantees to performance sensitive applications and described ongoing and future work to address the challenges.

## ACKNOWLEDGMENTS

This work is supported in part by the AFOSR DDDAS FA9550-13-1-0227 and NSF US Ignite CNS 1531079. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of NSF and AFOSR.

## REFERENCES

- [1] Martin Arlitt, Manish Marwah, Gowtham Bellala, Amip Shah, Jeff Healey, and Ben Vandiver. Iotabench: an internet of things analytics benchmark. In *Proceedings of the 6th ACM/SPEC International Conference on Performance Engineering*, pages 133–144. ACM, 2015.
- [2] Youakim Badr, Salim Hariri, Youssif AL-Nashif, and Erik Blasch. Resilient and trustworthy dynamic data-driven application systems (dddas) services for crisis management environments. *Procedia Computer Science*, 51:2623 – 2637, 2015.
- [3] Flavio Bonomi, Rodolfo Milito, Jiang Zhu, and Sateesh Addepalli. Fog computing and its role in the internet of things. In *Proceedings of the first edition of the MCC workshop on Mobile cloud computing*, pages 13–16. ACM, 2012.
- [4] Cache allocation technology improves real-time performance. <http://www.intel.com/content/www/us/en/communications/cache-allocation-technology-white-paper.html>.
- [5] Faruk Caglar, Shashank Shekhar, Aniruddha Gokhale, and Xenofon Koutsoukos. Intelligent, performance interference-aware resource management for iot cloud backends. In *2016 IEEE First International Conference on Internet-of-Things Design and Implementation (IoTDI)*, pages 95–105. IEEE, 2016.
- [6] Changbeom Choi, Kyung-Min Seo, and Tag Gon Kim. Dexsim: an experimental environment for distributed execution of replicated simulators using a concept of single simulation multiple scenarios. *Simulation*, 90(4):355–376, 2014.
- [7] Byung-Gon Chun, Sunghwan Ihm, Petros Maniatis, Mayur Naik, and Ashwin Patti. Clonecloud: elastic execution between mobile device and cloud. In *Proceedings of the sixth conference on Computer systems*, pages 301–314. ACM, 2011.
- [8] Brian F Cooper, Adam Silberstein, Erwin Tam, Raghu Ramakrishnan, and Russell Sears. Benchmarking cloud serving systems with ycsb. In *Proceedings of the 1st ACM symposium on Cloud computing*, pages 143–154. ACM, 2010.
- [9] Eduardo Cuervo, Aruna Balasubramanian, Dae-ki Cho, Alec Wolman, Stefan Saroiu, Ranveer Chandra, and Paramvir Bahl. Maui: making smartphones last longer with code offload. In *Proceedings of the 8th international conference on Mobile systems, applications, and services*, pages 49–62. ACM, 2010.
- [10] Frederica Darella. Dynamic data driven applications systems (dddas)—a transformative paradigm. In *Computational Science—ICCS 2008*, pages 5–5. Springer, 2008.
- [11] Christina Delimitrou and Christos Kozyrakis. Ibench: Quantifying interference for datacenter applications. In *Workload Characterization (IISWC), 2013 IEEE International Symposium on*, pages 23–33. IEEE, 2013.
- [12] Tharam Dillon, Chen Wu, and Elizabeth Chang. Cloud computing: issues and challenges. In *Advanced Information Networking and Applications (AINA), 2010 24th IEEE International Conference on*, pages 27–33. Ieee, 2010.
- [13] Michael Ferdman, Almutaz Adileh, Yusuf Onur Koçberber, Stavros Volos, Mohammad Alisafae, Djordje Jevdjic, Cansu Kaynak, Adrian Daniel Popescu, Anastasia Ailamaki, and Babak Falsafi. Clearing the clouds: a study of emerging scale-out workloads on modern hardware. In *Proceedings of the 17th International Conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS 2012, London, UK, March 3-7, 2012*, pages 37–48, 2012.
- [14] Debessay Fesehaye, Yunlong Gao, Klara Nahrstedt, and Guijun Wang. Impact of cloudlets on interactive mobile cloud applications. In *Enterprise Distributed Object Computing Conference (EDOC), 2012 IEEE 16th International*, pages 123–132. IEEE, 2012.
- [15] David G Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60(2):91–110, 2004.
- [16] Jason Mars, Lingjia Tang, Robert Hundt, Kevin Skadron, and Mary Lou Soffa. Bubble-up: Increasing utilization in modern warehouse scale computers via sensible co-locations. In *44th annual IEEE/ACM International Symposium on Microarchitecture*, pages 248–259. ACM, 2011.
- [17] Massimiliano Rak, Antonio Cuomo, and Umberto Villano. Mjades: Concurrent simulation in the cloud. In *Complex, Intelligent and Software Intensive Systems (CISIS), 2012 Sixth International Conference on*, pages 853–860. IEEE, 2012.
- [18] Nishanth Shankaran, Xenofon Koutsoukos, Chenyang Lu, Douglas C. Schmidt, and Yuan Xue. Hierarchical Control of Multiple Resources in Distributed Real-time and Embedded Systems. In *Proceedings of the Euromicro Conference on Real-Time Systems (ECRTS 06)*, Dresden, Germany, July 2006.
- [19] Shashank Shekhar, Hamzah Abdel-Aziz, Michael Walker, Faruk Caglar, Aniruddha Gokhale, and Xenofon Koutsoukos. A simulation as a service cloud middleware. *Annals of Telecommunications*, 71(3):93–108, 2016.
- [20] Shashank Shekhar, Ajay Chhokra, Anirban Bhattacharjee, Guillaume Aupy, and Aniruddha Gokhale. INDICES: Exploiting Edge Resources for Performance-aware Cloud-hosted Services. In *1st IEEE/ACM International Conference on Fog and Edge Computing (ICFEC) (to appear)*, Madrid, Spain, May 2017. IEEE.
- [21] Lei Wang, Jianfeng Zhan, Chunjie Luo, Yuqing Zhu, Qiang Yang, Yongqiang He, Wanling Gao, Zhen Jia, Yingjie Shi, Shujie Zhang, et al. Bigdatabench: A big data benchmark suite from internet services. In *High Performance Computer Architecture (HPCA), 2014 IEEE 20th International Symposium on*, pages 488–499. IEEE, 2014.