# Report on the Sixth International Workshop on Cloud Data Management (CloudDB 2014)

Shuai Ma
SKLSDE Lab
Beihang University
Beijing, China
mashuai@buaa.edu.cn

Xiaofeng Meng
School of Information
Renmin University of China
Beijing, China
xfmeng@ruc.edu.cn

Fusheng Wang
Dept of Biomedical Informatics
Emory University
Atlanta, USA
fusheng.wang@emory.edu

## 1. INTRODUCTION

The workshop series on Cloud Data Management (CloudDB) was held successfully from 2009 to 2013, co-located with the ACM Conference on Information and Knowledge Management (CIKM) [1, 3–6]. The sixth International Workshop on Cloud Data Management was held in Chicago, IL, USA on March 31, 2014, co-located with the 30th IEEE International Conference on Data Engineering (ICDE) [2]. CloudDB is dedicated to address the challenges in managing big data in the cloud computing environment, and identifying information of value to business, science, government, and society, and it continues serving as a premier forum for researchers and practitioners to present research progress and share ideas in the cloud data management area.

Data management is one of the most important research areas in cloud computing. The huge volumes of data in cloud computing environments pose big infrastructure challenges, including data storage at Petabyte scale, massive parallel query execution, facilities for analytical processing, and online query processing. Meanwhile, the emergence of large data centers and computer clusters has created a new business model, cloud-based computing, for the provision of large-scale computer facilities, where businesses and individuals can rent storage and computing capacities, rather than make significant capital investments to construct. Cloud-based data storage and management is a rapidly expanding business. Whilst these emerging services have substantially reduced the cost of data storage and delivery, there is significant complexity involved in ensuring that they can sustain consistent and reliable operations under peak loads. A cloud-based environment has technical requirements to manage data center virtualization, lower cost and boost reliability by consolidating systems in the cloud. In addition, cloud systems ideally should be geographically dispersed, both to reduce their vulnerability to natural disasters and other catastrophes and to bring data and computation closer to a possibly global user base. This trend brings rise to new and complex technical challenges in the areas of distributed data interoperability and mobility.

This year, the program committee accepted ten papers from nineteen submissions, by authors from Australia, Brazil, China, Germany, Japan, USA and Switzerland. The program covered a variety of topics, including quality of services, query processing, system architecture, and benchmarking. In addition, the program included two invited keynote talks from leading cloud computing researchers.

Many people contributed to the success of this year's CloudDB. First, we would like to thank all authors for submitting their contributions and all attendees for being generous and warm-hearted in all interactions. We would also like to express our deepest gratitude to the program committee members who worked hard in reviewing papers and providing suggestions for improvements. We also give special thanks to our keynote speakers, Geoffrey Fox and Xiaodong Zhang. Finally, we would express our great appreciation to Beihang University, Renmin University of China and Emory University for their supports.

## 2. KEYNOTE TALKS

The two keynote talks covered topics on big data processing systems and on big data uses and architecture integrating high performance computing and the Apache stack, respectively.

- The first keynote talk was delivered by Xiaodong Zhang from the Ohio State University on "Building Big Data Processing Systems under New Computing Model". Firstly, the speaker introduced the implications of big data processing from a system perspective. (a) Conventional database systems are not designed for big data. (b) Big data users require cost-

effective solutions for their analytics because conventional solutions are not scalable and affordable. (c) System designers and practitioners highly demand various new software tools for big data processing and analytics. (d) Computing paradigm for data processing has been shifted from a scale-up model for high performance to the one for high throughput as the main role of computers becomes data centers. Then, the speaker described how the system community addressed the above mentioned issues with a case study on major technical advancements in Apache Hive, widely adopted by many organizations for big data analytics. In short, the speaker presented a community-based effort and showed how academic research laid a foundation for Hive to improve its daily operations in production systems.

- The second keynote talk was delivered by Geoffrey Charles Fox from Indiana University on "Multi-faceted Classification of Big Data Uses and Proposed Architecture Integrating High Performance Computing and the Apache Stack". The speaker firstly gave a nice introduction of the NIST collection of 51 use cases and their scope over industry, government and research areas. Then, the speaker proposed that in many cases it was wise to combine the well known commodity best practice (often Apache) Big Data Stack (with 120 software subsystems) with high performance computing technologies. Finally, the speaker identified key layers where HPC Apache integration was particularly important: File systems, Cluster resource management, File and object data management, Inter process and thread communication, Analytics libraries, Workflow and Monitoring.

## 3. RESEARCH PAPERS

The ten accepted papers were divided into four sessions on quality of services, query processing, system architecture and benchmarking, chaired by Dr. Ablimit Aji from HP Labs and Prof. Raymond Chi-Wing Wong from the Hong Kong University of Science and Technology.

### 3.1 Quality of Services

The proliferation of cloud computing has attracted the deployment of many applications based on the cloud platforms. This obviously raises the issue that how the cloud providers could support high quality services and eventually lead to the complete satisfaction of all sorts of requirements from the cloud users, e.g., in a pay-as-you-go fashion.

(1) Paper "Towards Improvements on the Quality of Service for Multi-Tenant RDBMS in the Cloud" by Leonardo O. Moreira, Victor A. E. Farias, Flávio R. C. Sousa, Gustavo A. C. Santos, José Gilvan Rodrigues Maia, and Javam C. Machado. This paper focuses on the multi-tenant approaches to improving the use of resources, by reducing the operation cost of services, and it proposes an approach to improving quality of service for multi-tenant RDBMS, by employing the migration techniques of tenants, system monitoring, allocation strategy, forecast approach, and benefits of cloud infrastructure to improve performance and reduce provider cost.

(2) Paper "PolarDBMS: Towards a Cost-Effective and Policy-Based Data Management in the Cloud" by Ilir Fetai, Filip M. Brinkmann and Heiko Schuldt. This paper reports the work in progress PolarDBMS towards a flexible and dynamically adaptable system for managing data in the Cloud. PolarDBMS derives policies from application and service objectives, from which it automatically deploys the most efficient and cost-optimized set of modules and protocols, and monitors their compliance. Further, the modules and their customization are changed at running time if necessary.

(3) Paper "SLA-driven Workload Management for Cloud Databases" by Dimokritos Stamatakis and Olga Papaemmanouil. This paper focuses on the challenges related to Service-Level-Agreements (SLAs) specification and management, and argues that SLA management for cloud databases should itself be offered as a cloud-based automated service. For this, it talks about the design of a framework that (a) enables the specification of custom application level performance SLAs and (b) offers workload management mechanisms that can automatically customize their functionality towards meeting these application-specific SLAs.

### 3.2 Query Processing

The next three papers investigate query processing in cloud based systems for different types of data, from relational data to graphs to data streams.

(4) Paper "Parallel Join Executions in RAMCloud" by Christian Tinnefeld, Donald Kossmann, Joos-Hendrik Boese and Hasso Plattner. This paper studies the utilization of the processing power of large-scale storage systems for supporting

query execution, and it evaluates the parallel execution of join operations in Stanford's RAMCloud, a DRAM-based storage system connected via RDMA-enabled network adapters. To do this, a system model is proposed to derive the execution costs for the Grace Join, the Distributed Block Nested Loop Join, and the Cyclo Join algorithm and their corresponding implementations in RAMCloud, together with a set of heuristics for parameterizing the execution of multiple join operations in parallel for maximizing the throughput.

(5) Paper "Data Stream Partitioning Re-Optimization Based on Runtime Dependency Mining" by Emeric Viel and Haruyasu Ueda. This paper focuses on the optimization of communication cost in distributed data stream processing systems. As programs made of multiple queries can be parallelized by partitioning input streams according to partitioning keys, different partitioning keys for different queries often require intermediary re-partitions, which, as a result, causes extra communication cost and reduces the throughput. It is known that re-partitionings could be avoided by detecting dependencies among the partitioning keys applicable to each query. This paper extends existing (compile-time) partitioning optimization methods by adding a runtime re-optimization module, based on the usage of temporal approximate dependencies among partitioning keys, a type of dependency approximately valid over a sliding window.

(6) Paper "Neighbor-base Similarity Matching for Graphs" by Hang Zhang, Hongzhi Wang, Jianzhong Li and Hong Gao. This paper investigates approximate graph pattern matching which has various cloud data management related applications.

## 3.3 System Architecture

The following two papers aim at novel cloud systems that could provide better services to deal with the requirements of various applications.

(7) Paper "Curracurrong Cloud: Stream Processing in the Cloud" by Vasvi Kakkad, Akon Dey, Alan Fekete and Bernhard Scholz. This paper focuses on the stream processing systems in a cloud environment, and describes a novel system *Curracurrong Cloud* that allows the computation and data origins to share a cloud-hosted cluster, offers a lightweight algebraic-style description of the processing pipeline, and

supports automated placement of computation among computing resources.

(8) Paper "ORESTES: a Scalable Database-as-a-Service Architecture for Low Latency" by Felix Gessert, Florian Bcklers and Norbert Ritter. This paper first describes three major problems that hinder the applicability of database systems in cloud environments: (a) high network latencies for remote/mobile clients, (b) lack of elastic horizontal scalability mechanisms, and (c) missing abstraction of storage and data models. It then introduces an architecture, a REST/HTTP protocol and a set of algorithms towards solving these problems with a Database-as-a-Service middleware called ORESTES, which exposes cloud-hosted NoSQL database systems through a scalable tier of REST servers. These together provide database-independent and object-oriented schema design, a client-independent REST-API for database operations, globally distributed caching, cache consistency mechanisms and optimized database ACID transactions.

## 3.4 Benchmarking

Database system benchmarks like TPC-C and TPC-E are very useful on evaluating the performance of DBMS systems. The last two papers study the benchmarking of Web-scale transactional databases and cloud-based tagging services, respectively.

(9) Paper "YCSB+T: Benchmarking Web-scale Transactional Databases" by Akon Dey, Alan Fekete, Raghunath Nambiar, Uwe Röhm. Cloud service benchmark frameworks like YCSB are designed for performance evaluation of distributed NoSQL key-value stores, which initially did not support transactions. Recent implementations of Web-scale distributed NoSQL systems, such as Spanner and Percolator, offer transaction features to cater to new Web-scale applications. To fix this gap in standard benchmarks, this paper first identifies the issues to be addressed when evaluating transaction support in NoSQL databases. It then introduces YCSB+T, an extension of YCSB, that wraps database operations within transactions, incorporates a validation stage to detect and quantify database anomalies resulting from any workload, and gathers metrics that measure the transactional overhead.

(10) Paper "Benchmarking Cloud-based Tagging Services" by Tanu Malik, Kyle Chard and Ian Foster. Tagging services have emerged as a

useful and popular way to organize data resources. However, an efficient implementation of tagging services is a challenge since highly dynamic schemas and sparse, heterogeneous attributes must be supported within a shared, openly writable database. This case-study paper describes a benchmark for tagging services, and proposes benchmarking modules that can be used to evaluate the suitability of a database for workloads generated from tagging services. The modules have been incorporated as part of OLTP-Bench, a cloud-based benchmarking infrastructure, to understand performance characteristics of tagging systems on several relational DBMSs and cloud-based database-as-a-service (DBaaS) offerings.

We would like to encourage the interested readers to look into our workshop proceedings for the details of the ten accepted research papers.

## 4. CONCLUSIONS

CloudDB was held successfully six times associated with CIKM from 2009 to 2013 and with ICDE in 2014. During these six years, cloud computing has undergone significant development and attracted major interests from both industry and academia. Topics of CloudDB cover all sorts of aspects on cloud data management, such as cloud computing infrastructure for big data storage and computing, cloud privacy and security, query processing and indexing access control in cloud computing systems, service-level agreements, business models and pricing policies, novel data-intensive computing applications, massive parallel query execution, data intensive scalable computing, and large-scale analytical methodology and algorithm.

Further, all the participants agreed that many open challenges remain open for cloud data management, such as big data management in the cloud and cloud data security and privacy.

## 6. REFERENCES

[1] D. W.-L. Cheung, I.-Y. Song, W. W. Chu, X. Hu, and J. J. Lin, editors. *The 18th ACM Conference on Information and Knowledge Management, Hong Kong, China, November 2-6*. ACM, 2009.

[2] I. Cruz, E. Ferrari, and Y. Tao, editors. *The 30th IEEE International Conference on Data Engineering, Chicago, IL, USA, March 31-April 4*. IEEE Computer Society, 2014.

[3] Q. He, A. Iyengar, W. Nejdl, J. Pei, and R. Rastogi, editors. *The 22nd ACM International Conference on Information and Knowledge Management, San Francisco, CA, USA, October 27 - November 1*. ACM, 2013.

[4] J. Huang, N. Koudas, G. J. F. Jones, X. Wu, K. Collins-Thompson, and A. An, editors. *The 19th ACM Conference on Information and Knowledge Management, Toronto, Ontario, Canada, October 26-30*. ACM, 2010.

[5] C. Macdonald, I. Ounis, and I. Ruthven, editors. *The 20th ACM Conference on Information and Knowledge Management, Glasgow, United Kingdom, October 24-28*. ACM, 2011.

[6] X. wen Chen, G. Lebanon, H. Wang, and M. J. Zaki, editors. *The 21st ACM International Conference on Information and Knowledge Management, Maui, HI, USA, October 29 - November 02*. ACM, 2012.