

Chapter 14

Big Data Processing and Big Analytics

Jaroslav Pokorný

Charles University, Czech Republic

Bela Stantic

Griffith University, Australia

ABSTRACT

Development and wide acceptance of data-driven applications in many aspects of our daily lives is generating waste volume of diverse data, which can be collected and analyzed to support various valuable decisions. Management and processing of this big data is a challenge. The development and extensive use of highly distributed and scalable systems to process big data have been widely considered. New data management architectures (e.g., distributed file systems and NoSQL databases) are used in this context. However, features of big data like their complexity and data analytics demands indicate that these concepts solve big data problems only partially. A development of so called NewSQL databases is highly relevant and even special category of big data management systems is considered. In this chapter, the authors discuss these trends and evaluate some current approaches to big data processing and analytics, identify the current challenges, and suggest possible research directions.

INTRODUCTION

One of the main characteristics of Big Data is its volume, which exceeds the normal range of databases in practice. For example, web clicks, social media, scientific experiments, and datacentre monitoring belong among data sources that generate vast amounts of raw data every day. To enable management of such volume of data there is need to have appropriate processing, i.e. Big Data computing, therefore, Big Data processing is an issue of the highest importance, generally referred in literature as *Big Analytics*. Big Analytics is another buzzword denoting a combination of Big Data and Advanced Analytics. J. L. Leidner (R&D at Thomson Reuters) in the interview with R. V. Zicari (ODMS.org, 2013) emphasizes that buzzwords like “Big Data” do not by themselves solve any problem – they are not magic bullets. He gives an advice how to tackle and solve any problem. There is need to look at the input data, specify

DOI: 10.4018/978-1-5225-8446-9.ch014

the desired output, and think hard about whether and how you can compute the desired result, which is basically nothing but “good old” computer science.

The recent advances in new hardware platforms, methods, algorithms as well as new software systems enables efficient Big Data processing and Big Analytics. Big Data and Big Analysis are decision making drivers in a 21st century world driven by web, mobile and IoT technologies.

Effective use of systems incorporating Big Data in many application scenarios requires adequate tools for storage and processing such data at low-level and analytical tools on higher levels. Moreover, applications working with Big Data are both transactional and analytical. However, they require usually different architectures.

Big Analytics is the most important aspect of Big Data computing mainly from a user’s point of view. Unfortunately, large datasets are expressed in different formats, e.g., relational, XML, textual, multimedia or RDF, which may cause difficulties in their processing by data mining algorithms. Also, increasing either data volume in a repository or the number of users of this repository requires more feasible solution of scaling in such dynamic environments than it is offered by traditional database architectures.

Clearly, Big Analytics is done also on big amounts of transaction data as extension of methods used usually in technology of *data warehouses* (DW). Generally DW technology is focused on structured data in comparison to much richer variability of Big Data as it is understood today. Therefore, analytical processing of Big Data Analytics requires not only new database architectures but also new methods for integrating and analyzing heterogeneous data.

Big Data storage and processing are essential for cloud services. This reinforces requirements on the availability and scalability of computational resources offered by cloud services.

Users have a number of options associated with above mentioned issues. For storing and processing large datasets they can use:

- Traditional DBMS - relational (SQL), OO, OR,
- Traditional relational parallel database systems (shared nothing architectures),
- Distributed file systems and Hadoop technologies,
- Key-value datastores (so called NoSQL databases),
- New database architectures (e.g., NewSQL databases).

In particular, three last categories are not mutually exclusive and can and they should co-exist in many enterprises. Another adept to coexistence is, of course, a relational database management system (RDBMS) that ensures transactional data processing in the enterprise.

The NoSQL and NewSQL databases present themselves as data processing alternatives that can handle huge volumes of data and provide the required scalability. NoSQL databases are a type of databases which were initiated by Web companies in early 2000s. NewSQL databases are aiming to provide the scale-out advantages of NoSQL databases often on commodity hardware and maintain the transactional data consistency guarantees of traditional relational DBMS. They are also compatible with SQL. Especially, *massively parallel analytic databases* play an important role here. Algorithms supporting Big Analytics are presented on the top of these systems or they are a native part of their implementation.

The chapter is an attempt to cover principles and core features of these systems and to associate them to main application areas of Big Data processing and management in practice, particularly in relation to Big Analytics. We also focus in more extent on challenges and opportunities associated with Big Data. The text follows the work of Pokorny & Stantic (2016).

29 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage:
www.igi-global.com/chapter/big-data-processing-and-big-analytics/230694

Related Content

XML Data Integration: Merging, Query Processing and Conflict Resolution

Yan Qi, Huiping Cao, K. Selçuk Candan and Maria Luisa Sapino (2010). *Advanced Applications and Structures in XML Processing: Label Streams, Semantics Utilization and Data Query Technologies* (pp. 333-360).

www.irma-international.org/chapter/xml-data-integration/41511

A Logic Programming Perspective on Rules

Leon Sterling and Kuldar Taveter (2009). *Handbook of Research on Emerging Rule-Based Languages and Technologies: Open Solutions and Approaches* (pp. 195-213).

www.irma-international.org/chapter/logic-programming-perspective-rules/35860

Using Device Detection Techniques in M-Learning Scenarios

Ricardo Queirós and Mário Pinto (2013). *Innovations in XML Applications and Metadata Management: Advancing Technologies* (pp. 118-134).

www.irma-international.org/chapter/using-device-detection-techniques-learning/73176

A Fuzzy RDF Graph-Matching Method Based on Neighborhood Similarity

Guanfeng Li and Zongmin Ma (2019). *Emerging Technologies and Applications in Data Processing and Management* (pp. 184-198).

www.irma-international.org/chapter/a-fuzzy-rdf-graph-matching-method-based-on-neighborhood-similarity/230689

Modeling Temporal Information With JSON

Zhangbing Huang and Li Yan (2019). *Emerging Technologies and Applications in Data Processing and Management* (pp. 134-153).

www.irma-international.org/chapter/modeling-temporal-information-with-json/230687