

Simple and Efficient Coupling of Hadoop With a Database Engine

Jiamin Lu, Ralf Hartmut Güting

Faculty of Mathematics and Computer Science, FernUniversität Hagen, Germany

{jiamin.lu, rhg}@fernuni-hagen.de

Abstract

The growing need of processing massive amounts of data leads database researchers to explore the possibility of combining their existing single-computer database systems with the popular parallel processing platform Hadoop. These hybrid systems not only can keep the efficiency of database processing, but also achieve a remarkable scalability.

This poster intends to propose such a system named Parallel *SECONDO*. It combines Hadoop with a number of extensible *SECONDO* database engines, in order to scale up the capability of processing extensible data models in *SECONDO* to a cluster of computers. It is also evaluated with the join operation on standard, spatial and spatio-temporal data upon different sizes of clusters.

1 Motivation & Approach

Most existing Hadoop extensions like HadoopDB [3] rely on Hadoop to shuffle intermediate data, in order to get a balanced workload assignment on cluster nodes. During the process, intermediate data have to be transformed to key-value pairs and delivered to HDFS. These overheads cause a performance degradation of Hadoop hybrid systems.

In this poster, we present a novel method to couple Hadoop and our extensible database system *SECONDO* [4, 2] at the engine level rather than the SQL

level [5, 1], by performing the task of shuffling intermediate data with distributed databases instead of Hadoop. Therefore, all data are exchanged between database engines directly, hence to reduce the unnecessary transform and transfer overhead. A prototype system Parallel *SECONDO* is developed based on this method and a customized distributed file system PSFS (Parallel *SECONDO* File System) is proposed accordingly.

Moreover, Parallel *SECONDO* also develops a parallel data model to indicate distributed data and state parallel queries. It supports all existing *SECONDO* data types and their related operators, including spatial and spatio-temporal data (or moving objects). Thereby, the user can describe parallel queries as usual in *SECONDO* executable language, like using a normal single-computer system.

2 Evaluations

Parallel *SECONDO* is evaluated with two parallel join methods, HDJ and SDJ, which respectively rely on Hadoop and *SECONDO* to shuffle intermediate data. The evaluation is made in both our own small-scale cluster and large-scale clusters consisting of hundreds of Amazon Web Service (AWS) instances, revealing that HDJ well inherits the scalability from Hadoop, while SDJ performs more efficiently in limited scale clusters of up to 100 computers. As a result, we obtain for the first time a highly scalable generic system for moving objects management.

Acknowledgments

We are grateful for the research grant provided by AWS in Education, which supports our evaluation in EC2 and the first author is also thankful for the financial support from Chinese Scholarship Council (CSC).

Copyright © 2013 by the Association for Computing Machinery, Inc. (ACM). Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the Owner/Author.

SoCC'13, 1–3 Oct. 2013, Santa Clara, California, USA.
ACM 978-1-4503-2428-1. <http://dx.doi.org/10.1145/2523616.2525940>

References

- [1] Parallel Secondo. <http://dna.fernuni-hagen.de/secondo/ParallelSecondo>.
- [2] Secondo. <http://dna.fernuni-hagen.de/secondo/>.
- [3] A. Abouzeid, K. Bajda-Pawlikowski, D. Abadi, A. Silberschatz, and A. Rasin. HadoopDB: An Architectural Hybrid of MapReduce and DBMS Technologies for Analytical Workloads. *Proc. VLDB Endowment*, 2(1):922–933, 2009.
- [4] R. Güting, T. Behr, and C. Düntgen. SECOND0: A Platform for Moving Objects Database Research and for Publishing and Integrating Research Implementations. *IEEE Data Eng. Bull.*, 33(2):56–63, 2010.
- [5] J. Lu and R. Güting. Simple and Efficient Coupling of Hadoop With a Database Engine. *Fernuniversitat in Hagen, Informatik-Report*, 366, 2012. <http://dna.fernuni-hagen.de/papers/CouplingHadoop366.pdf>.