

# CMSC724: Literature Survey

## Improving MapReduce Performance Using Data Compression

Greg Benjamin, Samet Ayhan, Kishan Sudusinghe  
University of Maryland, College Park

### 1 Introduction

Although MapReduce provides programmers with tools and utilities to accomplish their parallel computing needs with large amounts of data, performance has reportedly been a bottleneck. Much recent work comparing the performance of MapReduce with parallel databases has concluded that many problems encountered by MapReduce users may be overcome by applying techniques learned from over three decades of research on parallel DBMSs. However, translating these techniques to a MapReduce implementation such as Hadoop presents a unique challenge that can lead to new design choices.

Inspired by Floratoru et al. [2], we propose to explore the area of column-oriented storage techniques for MapReduce which will investigate implementation of compression techniques. The paper only experiments with simple LZO compression, citing a belief that more advanced techniques which provide higher compression ratios induce too much overhead and actually harm performance. We plan to test this claim on a variety of compression techniques, in particular arithmetic coding, which we believe is well-suited to this application.

### 2 Inspiration: Floratou et al.

Our starting point for this project was a 2011 paper by Floratou et al., entitled *Column-oriented storage techniques for MapReduce*. This work seeks to address many of the issues related to performance in the Hadoop implementation of MapReduce by incorporating techniques used in column-oriented and parallel database systems. In particular, contributions of the paper include:

- A novel column-oriented data storage format called `ColumnInputFormat`, which improves on existing column-storage formats like `RCFile` [28]

by allowing each column of a relation to be stored in separate disk blocks, increasing efficiency pay-offs in caching and compression ratio.

- A *lazy deserialization* scheme based on the SkipList [26] data structure. The scheme is designed so that column values need only be read in from disk if they are actually used, and blocks of unused values may be skipped over using pointer information embedded in the column file. This avoids the heavy cost of deserializing every value in the column, whether or not that value is actually needed by the MapReduce task.
- Two techniques for compressing column files in a chunked manner, such that the SkipList deserialization scheme also avoids having to decompress a chunk if the values it contains are not needed. These techniques make use of the lightweight, non-optimal LZO algorithm [23] to do the compression of individual chunks.
- An experimental analysis of the above techniques, in comparison with other commonly-used storage formats and compression techniques. In general, `ColumnInputFormat` and the compression schemes presented here are found to provide an order-of-magnitude speedup over the other configurations.

Of particular interest to us is the authors' claim that using heavier compression algorithms with better compression ratios does not lead to any speedup. Apparently this is because such algorithms incur too much CPU overhead during decompression, and this outweighs any benefit gained by having to read in less data. However, such claims were *only* tested using the `zlib` compression library, which is an implementation of the Lempel-Ziv '77 algorithm [15].

## 3 Background

We now present necessary background information and additional work relevant to the concepts presented in Floratou et al. In this section, we focus on the MapReduce framework, data processing alternatives such as column-stores and parallel databases, and general data compression techniques.

### 3.1 MapReduce

MapReduce (MR) is a new programming framework and implementation model for processing large data sets. The model uses a map function to process key/value pairs and generate a set of intermediate key/value pairs. The MapReduce library groups together all intermediate values associated with the same intermediate key and passes them to the reduce function. The reduce function accepts an intermediate key and a set of values for that key and merges these intermediate values to form a smaller set of values.

These programs can be parallelized and executed on large cluster of commodity hardware resulting in highly scalable systems.

Dean et al. present a performance experiment in 2008 [17] where they test the system with a number of tasks (sort, grep, etc.). They notice that the input rate is higher than the shuffle rate and the output rate. They say that the reason is because most data is read from local disk. Also, they note that the shuffle rate is higher than the output rate because the output phase writes two copies of the sorted data. The reason why two copies are written is because of reliability and availability requirements.

### 3.2 Hadoop

Due to fact that our implementation will be built upon Hadoop, an open source framework implementing the MapReduce algorithm, we made an effort to learn more about it. We came across a book, *Hadoop In Action* [7], which provides a deep insight as to what Hadoop is, how it is installed and configured and how the MR libraries are used inside of user programs. In the book, Chuck Lam provides an answer to a simple question, we all may have asked us at some point: *what makes Hadoop better than standard relational databases?*

- One reason is that SQL (Structured Query Language) is by design targeted at structured

data. Many of Hadoops initial applications deal with unstructured data such as text. From this perspective Hadoop provides a more general paradigm than SQL. For working only with structured data, the comparison is more nuanced. In principle, SQL and Hadoop can be complementary, as SQL is a query language which can be implemented on top of Hadoop as the execution engine. But in practice, SQL databases tend to refer to a whole set of legacy technologies.

- Scaling commercial relational databases is expensive. Their design is more friendly to scaling up. To run a bigger database you need to buy a bigger machine. Whereas Hadoop is designed to be a scale-out architecture operating on a cluster of commodity PC machines. Adding more resources means adding more machines to the Hadoop cluster.
- A fundamental tenet of relational databases is that data resides in tables having relational structure defined by a schema. Although the relational model has great formal properties, many modern applications deal with data types that dont fit well into this model. Text documents, images, and XML files are popular examples. Also, large data sets are often unstructured or semistructured. Hadoop uses key/value pairs as its basic data unit, which is flexible enough to work with the less-structured data types. In Hadoop, data can originate in any form, but it eventually transforms into (key/value) pairs for the processing functions to work on.
- SQL is fundamentally a high-level declarative language. You query data by stating the result you want and let the database engine figure out how to derive it. Under MapReduce you specify the actual steps in processing the data, which is more analogous to an execution plan for a SQL engine. Under SQL you have query statements; under MapReduce you have scripts and codes. MapReduce allows you to process data in a more general fashion than SQL queries.
- Hadoop is designed for offline processing and analysis of large-scale data. It doesnt work for random reading and writing of a few records, which is the type of load for online transaction processing.

We also explored Hadoops architecture and found out that on a fully configured cluster, running Hadoop means running a set of daemons, or resident programs, on the different servers in the network. These daemons have specific roles; some exist only on one server, some exist across multiple servers. The daemons include:

- NameNode
- DataNode
- Secondary NameNode
- JobTracker
- TaskTracker

Hadoop Distributed File System (HDFS) is a filesystem designed for large-scale distributed data processing under frameworks such as MapReduce. In fact, it is considered that HDFS has its roots from the Google File System (GFS), which was introduced by Ghemawat et al. in 2003 [24]. The GFS is scalable and implemented for large, distributed, data-intensive applications.

According to Zaharia et al. [20], Hadoop’s performance is tied to its task scheduler, which assumes that cluster nodes are homogeneous and tasks make progress linearly. Based on this assumption, task scheduler decides when to re-execute tasks that appear to be stragglers. In practice, this assumption doesn’t hold. Zaharia et al. shows that Hadoop’s scheduler can cause severe performance degradation in heterogeneous environments. They experiment this Amazon’s EC2 and come up with a new scheduling algorithm: Longest Approximate Time to End (LATE). They claim that this algorithm can improve Hadoop response time by a factor of 2.

They identify a flaw in Hadoop with regard to threshold-based scheduling algorithm. So, they end up implementing a new scheduling algorithm, LATE, that uses estimated finish times to speculatively execute the tasks that hurt the response time the most. According to the authors, LATE performs significantly better than Hadoops default speculative execution algorithm in real workloads on Amazons Elastic Compute Cloud.

### 3.3 Column-stores

Choice of data structure has a great role in compression. Turning list of key/value pairs as HDFS uses into column-store would obviously allow us obtain

the best benefit out of compression. Hence we have carefully been exploring column-stores.

Column-stores store each database table column separately, with attribute values belonging to the same column stored contiguously, compressed, and densely packed, as opposed to traditional database systems that store entire records one after the other. Reading a subset of a tables columns becomes faster, at the potential expense of excessive disk-head seeking from column to column for scattered reads or updates. In 2009 paper, Abadi et al [8] indicate that one of the most-often cited advantages of column-stores is data compression. The intuitive argument for why column-stores compress data well is that compression algorithms perform better on data with low information entropy (high data value locality). However, column-stores come with a price; two of the most-often cited disadvantages of column-stores are write operations and tuple construction. Write operations are generally considered problematic for two reasons:

- Inserted tuples have to be broken up into their component attributes and each attribute must be written separately
- The dense packed data layout makes moving tuples within a page nearly impossible

Stonebraker et al. make a reference to Abadi et al. paper and implement a column-store system, C-Store. This implementation reflects all column-store properties as outlined in their 2005 paper [21].

From another perspective Dao Vo et al. in their 2007 paper [5] formalize the notion of column dependency as a way to capture information redundancy across columns and discuss how to automatically compute and use it to substantially improve table compression. They elaborate on a compression technique for table data based on discovering dependency relations among the columns of a table and using them to transform the data to enhance compressibility. Assuming some measure to assess the strength of column dependency, the notion of a k-transform is introduced to represent a data transformation plan based on column dependency.

According to the authors, experiments based on a variety of table data show that the presented algorithm performs well. They experiment the performance and put the system in a benchmark with **pzip**. The results show that their system outperforms **pzip**. According to the authors, **pzip** represents a significant improvement over conventional

techniques for compressing table data. However, its method of grouping columns based on some external compressor does not directly address the information redundancy arising from column dependencies. Further, the plan computation in **pzip**, i.e., its column grouping, is slow and must be done off-line for an entire class of tables. As such, there is no recovery mechanism in **pzip** if a computed plan no longer matches the data characteristics in a particular table. In contrast, the presented algorithm delivers dependency transform computation which is fast and can be used on-line for each table. In fact, for large tables where the table characteristics may change between sufficiently long sequences of rows, a new dependency transform could be recomputed whenever compression degradation is detected.

### 3.4 Compression Techniques

While surveying the relevant literature, we discovered several compression techniques which we consider interesting in the context of the problem described by Floratou et al.

In a 2007 study, Samara et al. [12] examined and compared six lossless compression methods in the context of the storage and transmission of serialized Java objects. In particular, the study explored the tradeoff between high compression ratio and fast compression/decompression performance. According to Samara et al., LZMA and **bzip2** strike the best compromise between these two metrics, which should make them particularly useful for improving performance in MapReduce according to the Floratou paper.

Nicolae 2010 [6] performed a similar analysis of various compression techniques, here in the context of transparent compression for cloud-based storage services. Performance of the algorithms was considered, but the primary metric here was compression ratio, as the paper was motivated by attempting to reduce storage fees in the commercial cloud. The implementation of a sampling algorithm to dynamically decide whether to compress data or not (based on the entropy of what was being stored) was also discussed. The paper focused particularly on LZO and **bzip**, and concluded that both were appropriate in such a context due to their relatively light CPU overheads, which again should make both viable candidates for our consideration.

The study conducted by Balakrishnan et al., in 2006 [4], focused heavily on lossless compression techniques targeting scientific applications deployed

in extreme-scale parallel machines (i.e., IBM Blue Gene/L). However, compression was primarily on large amount of log data. The study found that lossless compression facilitates storage, archiving, and network transfer of data saving storage and network bandwidth. Experimental results showed 35.5% improvement in compression ratio and 76.6% improvement in compression time over **bzip2** and 21% improvement over compression ratio and 72.5% compression time over LZMA level 9 compression by 7zip. The approach defines a compression pipeline with different compression algorithms used in each stage. Best compression technique for each stage is chosen from a set of compression utilities.

Entropy compression is another technique widely looked at in the literature to alleviate data movement bottleneck (i.e., IBM's DB2 - individual records are compressed using dictionary scheme). Study by Vijayshankar et. al [25] takes an approach based on a mix of column and row (tuple) coding. Columns are coded using Huffman coding to exploit skew in value frequencies that will result in variable length field codes. Field codes are then concatenated to form tuplecodes and then sorted and delta coded to take advantage of the lack of order within a relation. However, there exist many issues with the proposed approach. Row or page compression was simpler to implement but it worsened memory and cache behavior due to decompression, even though it reduces I/O costs. Experimental studies also suggest CPU cost is high for decompression.

The technique we consider to have the most potential, however, is arithmetic coding. As explained in a 1987 paper by Witten et al. [14], this technique is a lightweight dictionary encoding scheme with better performance than Huffman coding [13] in both CPU overhead and compression ratio. Like all compression algorithms, the goal is to transmit more common symbols with fewer bits than less probable ones. Arithmetic coding accomplishes this by mapping symbols to a subset of the interval  $[0, 1)$ . The encoded value of a message starts as the entire interval, but is then reduced to a subset thereof by each successive symbol in the message. Decoding of an interval works by repeating the process, reading off (in order) the symbols that reduce the starting range to one which contains the encoded interval, and then reducing the starting range accordingly, until the two converge. This scheme is lossless and encodes an entire message as a single interval, unlike other dictionary coding schemes, which work on a symbol-by-

symbol basis. We think this is especially applicable to the work done by Floratou et al., both because of the lightweights of the scheme, but also because such a scheme should operate well on chunks of data at a time, allowing us to maximize the gains of using a SkipList representation to store data.

## 4 MapReduce Performance

By far, the most frequent complaint we’ve encountered regarding the MapReduce framework (and the Hadoop implementation, particularly) is that the performance of the framework is unacceptably slow [22]. Many papers have been written detailing MapReduce efficiency concerns and how to fix them; after surveying the relevant literature, we believe that most approaches fit into one of four topics: improved performance reading data from HDFS into memory, better ways to lay the data out within files on disk, accessing data using an index rather than sequential scan, or reducing the amount of data read in using data compression. We now consider each of these in detail.

### 4.1 Serialization and Data Format

In the 2010 paper *MapReduce: A Flexible Data Processing Tool* [18], authors Dean and Ghemawat responded to many of the criticisms of MapReduce inefficiency presented in [22]. The authors immediately pointed out that all experiments in the Stonebraker et al. paper were conducted using raw text files as the underlying storage format. This is inherently slow, as each file must be deserialized during parsing; the binary data read in will automatically be converted into an ASCII string in memory before anything can be done with it. By contrast, most well-practiced users of MapReduce store their data in a binary format such as Google’s Protocol Buffers. This allows mapper tasks to skip the stage of deserialization entirely, significantly improving the performance of the map phase.

The Jiang et al. 2010 paper *The Performance of MapReduce: An In-depth Study* [11] conducted experiments to compare the load times of textual data with those of binary files in a variety of formats. The study found that binary file formats were indeed faster to read in, but that the improvement was less than would be expected. The authors reasoned that most input formats parse input using immutable Java objects (such as Strings and Integers) for each object read in, and that the high CPU overhead of

creating all these objects was to blame for such poor performance. Further experimentation using mutable objects during parsing was able to improve the performance of the load phase by a factor of 10, compared to a meager 2x improvement from switching from text to binary data storage.

### 4.2 Data Layout

Many papers have examined the improvements that may be gained using a different data layout within input files. These papers draw on experience gained from the sphere of traditional databases, where it has been shown that column-oriented storage can significantly outperform row-oriented storage for certain types of queries [10]. However, there are additional concerns that arise when this experience is applied to MapReduce; the penalty in a row-oriented store for having to access every attribute in a record is higher because these attributes must all be parsed and deserialized, even if they aren’t used. On the other hand, partitioning data across columns is troublesome because columns may not be shuffled onto the same map nodes, and so reconstructing a given record may require communication across the network to access attributes for all the relevant columns.

Most attempts to resolve this issue draw on the PAX (Partition Attributes Across) file format [1]. PAX was a scheme put forward for column-stores to improve cache performance by striking a compromise between a pure row-store and a pure column-store. Essentially, relations are carved up into blocks, where each block is a contiguous group of records that will fit within a single page in memory. Then, blocks are vertically partitioned and written to disk as a mini-column-store. In this way, one gets the benefits of performing fewer I/Os to query a relation when only touching a few attributes, but cache performance is better when it comes time to reconstruct tuples, since the entire block (and therefore all necessary attributes) is already in cache. This was experimentally shown to lead to speedups between 10-50%, depending on the type of query.

Facebook’s RCFfile and the Trojan layout proposed in [3] both make use of a similar technique to improve performance in MapReduce. RCFfile applies the PAX scheme almost verbatim, but uses “row groups” (each the size of an HDFS block) instead of memory pages as the unit of column storage. Since each HDFS block is stored atomically on a single map node, this resolves the problem of having to communicate over

the network to reconstruct tuples, and successfully increases the performance of map tasks by about 15%.

The Trojan layout paper takes this a step further, suggesting that we can improve performance even more by vertically partitioning row groups into “column groups” of various sizes. For a relation  $R$  with 5 attributes ( $R = (a, b, c, d, e)$ ), we could store  $R$  as 5 distinct columns, or as a pair of column groups  $(a, b)$  and  $(c, d, e)$ , or as a set of 3 groups, or as anything in between. Since HDFS replicates data blocks anyway, we can precompute a couple of these column groupings which are optimal for the expected query workload, and then store a different one in each replica. By dynamically re-routing queries to use the column grouping most efficient for the query type, the authors observe a 3-5x speedup over PAX and other layouts. Of course, this pre-supposes that we know our query workload *a priori*, which may be an unrealistic assumption in most cases.

The Floratou et al. paper also uses column-oriented storage, but deviates from prior work in that columns are stored contiguously across all relations, rather than only contiguously within row groups. The `ColumnInputFormat` modifies HDFS’s block placement policy to ensure that column data for the same records does indeed end up on the same map node, so that no network communication is necessary to reconstruct tuples.

### 4.3 Indexing

The Dean and Ghemawat 2010 paper also responded to a previous criticism that MapReduce cannot make use of any indexing structures by pointing out that “natural indicies” in the structure of the format of input files can be exploited for improved performance. For example, if a MapReduce task operates on log data and the logs files are named with the dates of their rollover, map tasks can discard data from old logs without having to read it in first by simply examining filenames. This proposal was also examined experimentally in [11], and it was confirmed that such natural indices could contribute a performance speedup of 3-5x.

Hadoop++ [19] implemented a “lightweight, non-invasive” technique for generating actual data indices during data load and embedding them in input files. These indices could then be utilized by customized `InputFormat` classes for faster access during the map phase. It was shown that this technique could improve performance 5-20x, depending on the type of query. Unfortunately, such techniques also

incurred a 10x performance hit during data loading to actually build the index.

While embedded indexing techniques remain an open area of research, we think that the quick-and-dirty nature of most MapReduce tasks makes it unlikely that the extra work to build an index will be seen as worthwhile in the common case. However, exploiting natural indices in the input data is certainly a beneficial technique and can really help performance.

### 4.4 Data Compression

The final category for performance-boosting techniques involves compressing input data so that fewer I/Os are required to load everything into memory. Compression of table data is a standard and well-known technique for improving performance in both row-stores and column-stores [9]. However, in the sphere of MapReduce this seems to be a largely unexplored direction of research. Aside from the Floratou et al. paper, we were only able to find two papers which considered the implications of data compression in the MapReduce framework.

The `RCFile` paper briefly mentions the use of a compression algorithm such as `gzip` to reduce the size of columns within a block, but does not go into substantial detail, nor consider the results of using other compression algorithms at all.

Although not much of work is performed within the Hadoop MapReduce framework, a study by Jacopo et. al [16] introduce a MapReduce programming model based compression and decompression technique that have linearly scaled performance gains. MapReduce framework handled compression/decompression to reduce the amount of work done but the dictionary scheme is used to store data. However, dictionary scheme introduced had many failures as centralized approach is not feasible due to need of large memory to store the dictionary. If dictionary did not fit in memory, accessing it through storage is very costly due to seeks and I/O costs. Though caching can help to increase the performance, only few terms will benefit from such approach. Majority of the terms occur very rarely hence it hard to cache and gain performance. Compression is order of magnitude faster compared to decompression in MapReduce framework promising improvements with possibly better technique than dictionary.

Another paper we were able to find focusing on compression in MapReduce was a workshop submission from 2010 by Chen et al. [27] on energy-efficiency tradeoffs for applying compression to input data. Es-

sentially, the authors ran a series of experiments demonstrating that enabling Hadoop’s built-in `gzip` and `LZO` compression options on input and output data files can shift a portion of the workload from I/O to CPU time. Since CPU time is less costly in terms of energy consumption, the authors concluded that compression can be used to make MapReduce “greener”. However, the compression schemes built in to Hadoop simply compress entire files before being written to HDFS, rather than doing anything particularly clever or insightful.

## 5 Next Steps

Because of the sparsity of work involving applying compression techniques to improve the performance of the MapReduce framework, we think this continues to be an interesting area of work with some potentially low-hanging fruit. We plan to continue along this line of inquiry by contacting the authors of Floratou et al. and trying to get ahold of their `ColumnInputFormat` source code. If successful, we plan to repeat their experiments, and then continue further by hooking in a variety of different compression schemes to see how performance is affected. Of particular interest are Huffman and arithmetic coding, `LZO`, `LZMA`, and `bzip`. We would also like to consider ways in which the benefits of compression could be increased through the use of appropriate chunking at layout schemes. Hopefully this line of inquiry will be fruitful in terms of improving the performance of the MapReduce framework.

## References

- [1] A. Ailamaki, D. DeWitt, M. D. Hill, and M. Skounakis. Weaving relations for cache performance. In *VLDB*, pages 169180, 2001.
- [2] Avriila Floratou, Jignesh M. Patel, Eugene J. Shekita, and Sandeep Tata. 2011. Column-oriented storage techniques for MapReduce. *Proc. VLDB Endow.* 4, 7 (April 2011), 419-429.
- [3] Alekh Jindal, Jorge-Arnulfo Quian-Ruiz, and Jens Dittrich. 2011. Trojan data layouts: right shoes for a running elephant. In *Proceedings of the 2nd ACM Symposium on Cloud Computing (SOCC '11)*. ACM, New York, NY, USA, Article 21, 14 pages. DOI=10.1145/2038916.2038937 <http://doi.acm.org/10.1145/2038916.2038937>
- [4] Balakrishnan, R.; Sahoo, R.K., "Lossless compression for large scale cluster logs," *Parallel and Distributed Processing Symposium*, 2006. *IPDPS 2006. 20th International* , vol., no., pp.7 pp., 25-29 April 2006 doi: 10.1109/IPDPS.2006.1639692 URL: <http://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=1639692&isnumber=34366>
- [5] Binh Dao Vo and Kiem-Phong Vo. 2007. Compressing table data with column dependency. *Theor. Comput. Sci.* 387, 3 (November 2007), 273-283. DOI=10.1016/j.tcs.2007.07.016 <http://dx.doi.org/10.1016/j.tcs.2007.07.016>
- [6] Bogdan Nicolae. 2010. High throughput data-compression for cloud storage. In *Proceedings of the Third international conference on Data management in grid and peer-to-peer systems (Globe'10)*, Abdelkader Hameurlain, Franck Morvan, and A. Min Tjoa (Eds.). Springer-Verlag, Berlin, Heidelberg, 1-12.
- [7] Chuck Lam. *Hadoop in Action*. 2011.
- [8] Daniel J. Abadi, Peter A. Boncz, and Stavros Harizopoulos. 2009. Column-oriented database systems. *Proc. VLDB Endow.* 2, 2 (August 2009), 1664-1665.
- [9] D. J. Abadi, S. Madden, and M. Ferreira. Integrating Compression and Execution in Column-Oriented Database Systems. In *SIGMOD*, pages 671682, 2006.
- [10] D. Abadi, S. R. Madden, and N. Hachem. Column-Stores vs. Row-Stores: How Different Are They Really? In *SIGMOD*, pages 967980, 2008.
- [11] D. Jiang, B. C. Ooi, L. Shi, and S. Wu. The Performance of MapReduce: An In-depth Study. *PVLDB*, 3(1):472483, 2010.
- [12] Ghassan Samara, Ahmad El-Halabi, and Jalal Kawash. 2007. Compressing serialized java objects: a comparative analysis of six compression methods. In *Proceedings of the third conference on IASTED International Conference: Advances in Computer Science and Technology (ACST'07)*. ACTA Press, Anaheim, CA, USA, 427-431.
- [13] Huffman, D. (1952). A Method for the Construction of Minimum-Redundancy Codes. *Proceedings of the IRE*, 40(9), 1098-1101. Retrieved from <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=4051119>
- [14] Ian H. Witten, Radford M. Neal, and John G. Cleary. 1987. Arithmetic coding for data compression. *Commun. ACM* 30, 6 (June 1987), 520-540. DOI=10.1145/214762.214771
- [15] Jacob Ziv and Abraham Lempel; A Universal Algorithm for Sequential Data Compression, *IEEE Transactions on Information Theory*, 23(3), pp. 337343, May 1977.
- [16] Jacopo Urbani, Jason Maassen, and Henri Bal. 2010. Massive Semantic Web data compression with

- MapReduce. In Proceedings of the 19th ACM International Symposium on High Performance Distributed Computing (HPDC '10). ACM, New York, NY, USA, 795-802. DOI=10.1145/1851476.1851591 <http://doi.acm.org/10.1145/1851476.1851591>
- [17] J. Dean and S. Ghemawat. MapReduce: Simplified Data Processing on Large Clusters. CACM, 51(1):107113, 2008. <http://doi.acm.org/10.1145/214762.214771>
  - [18] J. Dean and S. Ghemawat. MapReduce: A Flexible Data Processing Tool. CACM, 53:7277, January 2010.
  - [19] J. Dittrich, J.-A. Quian-Ruiz, A. Jindal, Y. Kargin, V. Setty, and J. Schad. Hadoop++: Making a Yellow Elephant Run Like a Cheetah. PVLDB, 3(1):518529, 2010.
  - [20] Matei Zaharia, Andy Konwinski, Anthony D. Joseph, Randy Katz, and Ion Stoica. 2008. Improving MapReduce performance in heterogeneous environments. In Proceedings of the 8th USENIX conference on Operating systems design and implementation (OSDI'08). USENIX Association, Berkeley, CA, USA, 29-42.
  - [21] Mike Stonebraker, Daniel J. Abadi, Adam Batkin, Xuedong Chen, Mitch Cherniack, Miguel Ferreira, Edmond Lau, Amerson Lin, Sam Madden, Elizabeth O'Neil, Pat O'Neil, Alex Rasin, Nga Tran, and Stan Zdonik. 2005. C-store: a column-oriented DBMS. In Proceedings of the 31st international conference on Very large data bases (VLDB '05). VLDB Endowment 553-564.
  - [22] Michael Stonebraker, Daniel Abadi, David J. DeWitt, Sam Madden, Erik Paulson, Andrew Pavlo, and Alexander Rasin. 2010. MapReduce and parallel DBMSs: friends or foes?. Commun. ACM 53, 1 (January 2010), 64-71. DOI=10.1145/1629175.1629197 <http://doi.acm.org/10.1145/1629175.1629197>
  - [23] Oberhumer, M.F.X.J.: Lempel-ziv-oerhumer. <http://www.oberhumer.com/opensource/lzo> (2009)
  - [24] Sanjay Ghemawat, Howard Gobioff, and Shuntak Leung. 2003. The Google file system. SIGOPS Oper. Syst. Rev. 37, 5 (October 2003), 29-43. DOI=10.1145/1165389.945450 <http://doi.acm.org/10.1145/1165389.945450>
  - [25] Vijayshankar Raman and Garret Swart. 2006. How to writing a table dry: entropy compression of relations and querying of compressed relations. In Proceedings of the 32nd international conference on Very large data bases (VLDB '06), Umeshwar Dayal, Khuyong Whang, David Lomet, Gustavo Alonso, Guy Lohman, Martin Kersten, Sang K. Cha, and Young-Kuk Kim (Eds.). VLDB Endowment 858-869.
  - [26] W. Pugh. Skip Lists: A Probabilistic Alternative to Balanced Trees. CACM, 33(6):668676, 1990.
  - [27] Yanpei Chen, Archana Ganapathi, and Randy H. Katz. 2010. To compress or not to compress - compute vs. IO tradeoffs for mapreduce energy efficiency. In Proceedings of the first ACM SIGCOMM workshop on Green networking (Green Networking '10). ACM, New York, NY, USA, 23-28. DOI=10.1145/1851290.1851296 <http://doi.acm.org/10.1145/1851290.1851296>
  - [28] Y. He, R. Lee, Y. Huai, Z. Shao, N. Jain, X. Zhang, and Z. Xu. RCFile: A Fast and Space-efficient Data Placement Structure in MapReduce-based Warehouse Systems. In ICDE, 2011.