



A new approach for database “Column-oriented DBMS”

Neeraj Kumar Sirohi¹, Divakar Yadav²

¹Department of Computer Science and Engineering

IMS Engineering College

Ghaziabad (U.P) INDIA

²Faculty of Computer Science and Engineering

Jaypee Institute of Information Technology

Gautam Budh Nagar, Noida (U.P) INDIA

Abstract: As we all know that many different kind of organization need database for required information from different area like analytic, different types of reports and other many reasons. As the organization grows the database also grow into millions of records spared over many tables and business intelligent/science. And if any organization want any specific report strategic analysis, and other business intelligence activities for their work need data for many data warehouses where no need to access all the attributes just needed few attributes of a table. This is not possible for traditional DBMS to read and fetch only those attribute which is needed by query. Because as we all know that traditional row oriented DBMS access the data row by row in this process we read that column also those are not needed by the user which cause the time delay for query.

In this paper we going to present the need of column-oriented DBMS, architecture design and the working structure of a new type of DBMS which is called Column oriented DBMS also known as read-optimized relational DBMS .in contrasts with most currently used systems, which are write-optimized. Among the differences in its design are: The data can be stored in column rather than row, carefully coding of object into storage including main memory during query processing, storing an overlapping projections, rather than the current fare of tables and indexes, a non-traditional implementation of transaction which include easy availability for read only transaction and the extensive use of bitmap indexes to complement B-tree structure here we present the performance on a data and show that the system we are created called column store is much faster than traditional and popular commercial product is called row oriented system.

Keywords: WWW; Column oriented database; Tuple Mover; RS;WS; join index; sid

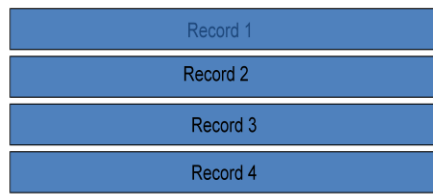
I. INTRODUCTION

Generally most DBMS implemented record-oriented storage system, where the attributes of a record (or tuple) are placed contiguously in memory. With this *row oriented* architecture, a Single disk write suffices to push all of the fields of a single record out to disk, that's why high Performance writes are achieved, and we said that DBMS with a row oriented structure is a *write-optimized* system. In contrast, systems oriented toward ad-hoc ("Ad Hoc" is a Latin phrase which means "for this purpose" and in today's parlance generally means "on the fly," or "spontaneously.") Query. An ad-hoc query is a query that is run at the spur of the moment, and generally we not saved the Ad-hic query for run again. The Ad- hoc queries are generally run using a SQL statement created by a tool or an administrator. So therefore, such a query is one that might suit a situation which is only there for the moment and later on will become irrelevant.) Querying of large amount of data should be *read-optimized* .data warehouse represent one class of read-optimized system. In such environments, column store architecture should be more efficient. A row-oriented implementation of a DBMS would store every attribute of a given row in sequence, with the last entry of one row followed by first entry of the next. On the other hand a column-oriented implementation of a DBMS would store every attribute of a given column in sequence, with the column values for the same column stored in sequence, with the end of one column followed by the beginning of the next (as shown in figure 1).

In this paper, we discuss the design of a column store called C-store that includes a number of novel features relative to existing systems. Commercial relational DBMSs store complete tuples of tabular data along with auxiliary B-tree indexes on attributes in table. Such indexes can be primary or secondary. In primary the row of the table are stored in close order on the specify attribute as possible. And in secondary indexes no attempts is made to keep the underlying records in order on the indexed attribute. Such indexes are not performing well in red-optimized world C-store physically stores a collection of columns, each stored on some attribute(s). Same types of column are stored on the same attribute which is called a “projection”; the same columns are present in multiple projections, possibly also stored on a different attribute in each. Now we can say that collection of

“Grid” computers will be the cheapest hardware architecture for computing and storage intensive application such as DBMSs [DEW192].

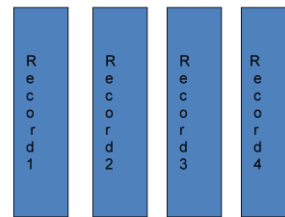
Terminology -- “Row Store”



Are represented as



Column Stores



Are represented as



Figure 1: Row store and Column store architectures

In C-Store approaches we combine in a single piece of system software, both a read-optimized column store and an update/insert-oriented writeable store, connected by a *tuple mover*, as noted in figure 2. At the top level, there is a small Writeable Store (WS) component, which is architected to support high performance inserts and updates. There is also a much larger component called the Read-optimized Store (RS), which is capable of supporting very large amounts of information. RS, as the name implies, is optimized for read and supports only a very restricted form of insert, namely the batch movement of records from WS to RS, a task that is performed by the tuple mover of figure 2.

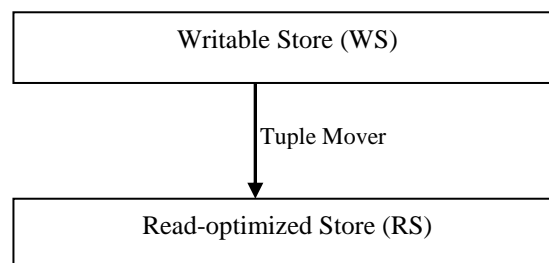


Figure 2: Representation of Writeable Store (WS) and Read-optimized Store (RS) components

Insertions are sent direct to WS, while deletion must be marked to RS for later purging by the tuple mover. to support of high-speed tuple mover , we are used a variant of the LSM-tree concept [ONE196], which support a merge out process that moves tuples from Ws to RS. The architecture of figure 2 must support transaction in an environment of many large ad-hoc, smaller update transaction and perhaps continuous insert. Finally the most commercial optimizer and executors are row-oriented. So both RS and WS are column-oriented, it makes sense to build a column-oriented optimizer and executor.

The architecture of column-oriented DBMS is reduced the number of disk per query. The important features of column-oriented DBMS are:

1. Redundant storage of element of a table in several overlapping Projection in different orders, and the query can be solved using the most advantageous projection
2. A column-oriented optimizer and executor , with different primitives than in a row-oriented system
3. A hybrid architecture with a WS component optimizer for frequent insert and update and an RS component optimizer for query performance
4. Heavily compressed column using one of several coding schemes.
5. High availability and improved performance through K-safety using a sufficient number of overlapping projections.
6. The use of snapshot isolation to avoid 2PC and locking for query

The rest of the paper we organized as follows. In Section 2 presents the related work done in this field. In Section discuss about Need and benefits of Column-Oriented DBMS. In Section 4 we present the data model implemented by column-oriented system, followed by the WS portion in section 5. In section 6 we consider the allocation of column-oriented data structure to node in a grid followed by column-oriented recovery by updates and transaction in section 7. Section 8 deals with tuple mover component and in section 9 present the query optimizer and executor and finally in section 10 we conclude this paper.

II. Related Work

One of the thrusts in the warehouse market is in maintained is so-called “data cubes”. This work done by Arbor software in early 1990’s which was effective at “slicing and dicing” large data set [GRAY97]. Efficiently building and maintained specific aggregates on store data set has been widely used if the workload cannot be calculated in advance, it is very difficult to decide what to precompute. The column stores aim at the latter problem. Storing data in column has been implemented in several systems that is Sybase IQ, Addamark, Bubba [COPE88], Monet[BONC04], and KBD, of these, Monet is probably closest to column store in design style. Similarly, storing tables using an inverted organization is well known, here every column stored using some type of indexing, and record identifiers are used to find corresponding column in other columns. Column store uses this sort of organization in WS but executed the architecture with RS and a tuple mover.

Roth and Van [ROTH93] provide the excellent summary of the techniques which have been developed. Our coding scheme is like to the same technique all of which are derived from a long history of work on the topic in the broader field of computer science [WITT87]. Our observation that it is possible to compute directly on compressed data has been made before [GRAE91, WESM00].

Finally, materialized views, snapshot isolation, transaction management, and high availability have also be extensively studied. The contribution of column store is an innovation combination of these techniques that simultaneously provides performance.

III. Need of Column-Oriented DBMS

According to Wintercorp’s 2005 Top Ten Program Summary, during the five year period between 1998 and 2003, the size of the largest data warehouses grew at an exponential rate, from 5TB to 30TB. But in the four year period between 2001 and 2005, that exponential rate increased, with the largest data warehouses growing from 10TB to 100TB. At the same time, the average hourly workload for data warehouses approached 2,000,000 SQL statements per hour, and in cases the number of SQL statements per hour reaching nearly 30,000,000.1

(A) Engineered for analytic performance

Because of the limitations described previously, there are limits to the performance which row-oriented systems can deliver when tasked with a large number of simultaneous, diverse queries. The usual approach of adding increasing numbers of space-consuming indexes to accelerate queries becomes untenable with diverse query loads, because of storage and CPU time required to load and maintain those indexes. With column-oriented systems, indexes are of a fundamentally different design. They are engineered to store the data, rather than as a reference mechanism pointing to another storage area containing the row data. As a result, only the columns used in a query need be fetched from storage. This I/O is conducted in parallel, because large columns are automatically distributed across multiple storage RAID groups. Once retrieved, columns are maintained in cache using caching algorithms intended to optimize memory access behavior patterns, further reducing storage traffic.

(B) Rapid joins and aggregation

In addition, data access streaming along column-oriented data allows for incrementally computing the results of aggregate functions, which is critical for business intelligence applications. In addition, there is no requirement for different columns of data to be stored together; allocating columnar data across multiple processing units and storage allows for parallel accesses and aggregations as well, increasing the overall query performance. An underlying query analyzer can evaluate ways to break the query down in ways that not only exploits the availability of multiple CPUS for parallel operations, but also arranges the execution steps to leverage the ability to interleave multiple operations across the available computational and memory resources, creating additional opportunities for parallelization. For example, complex aggregations and groupings can be “pipelined” by streaming the results of parallelized data scans into joins that in turn feed groupings, all happening simultaneously.

(C) Smaller storage footprint

One of the dependent factors among row-based systems to accommodate the aforementioned “data explosion” is the need for additional structures beyond the row-based data storage. These include the addition of indexes, tables for pre-aggregation and materialized views to the already burdensome storage requirements. A more efficient design is used in column-oriented systems, where data is stored within the indexes, eliminating the storage penalty of the types of indexes used in row-based systems. Some column-oriented systems not only store data by column, but store the data values comprising a column within the index itself; efficient bit-mapped indexes are selected to optimize storage, movement and query efficiency for each individual column’s data type.

(D) Suitability for compression

The columnar storage of data not only eliminates storage of multiple indexes, views and aggregations, but also facilitates vast improvements in compression, which can result in an additional reduction in storage while maintaining high performance. One example, typically employed where a column contains often-repeating values, involves tokenizing the commonly used data values, mapping those values to tokens that are stored for each row in the column. Mapping the original form of the data to a “token”, and storing the column as a list of

these tokens requires much less storage space and yet, to the application, is totally transparent. For example, instead of maintaining copies of city name values in address records, each city name (“Phoenix, AZ”) can be mapped to an integer value (“3449”) which requires 2 bytes to store, rather than 10-12 bytes. The resulting compression is 4-6 times in this example.

(E) *Optimized for query efficiency*

The bitmap-based data structures used in some column-oriented analytic stores provide superior query performance by using sophisticated aggregation and bitmap operations within and across columns. If we use a bitmap to describe data, we can make use of the mapping to also count unique occurrences at loading time, and therefore we can provide that pre-aggregated count rather than analyzing the actual column data. The resulting systems may provide orders of magnitude improvements in query processing performance. As an example, consider the use of a tokenized column for listing cities. Some queries, such as those that count occurrences, never need to access the data itself, because some amount of metadata (such as “counts”) that is an inherent result of preparing the column for storage can be captured at loading time.

IV. Data Model

Column-oriented system supports the standard relational logical data model, It represents relational tables using vertical fragmentation, by storing each column in a separate $\{(surrogate,value)\}$ table, called a BAT (Binary Association Table). The left column, also denoted as *head*, contains object-identifiers (OID), and the right column called *tail*, contains attribute values. The *head* is always a dense and sorted list. Structure of BAT table’s is show with Table 1 in figure 3. C-store use a low-level relational algebra which is known as the *BAT algebra*, it take two values BATs and scalar as parameter for the input. The result of an SQL query is a collection of BATs and the BATs (intermediate) always stored the complete result in it, and. BAT storage maintains two simple memory arrays, one store the values for the head and the second used to store the value for the tail column. For variable-width the tail column are divided into two parts. The first part contains the concatenated data values and the second part with offsets into the former. In column-oriented tables we are consists a unique *primary key* or be a *foreign key* that references a primary key of another table. The column-oriented query language is assumed to be SQL.

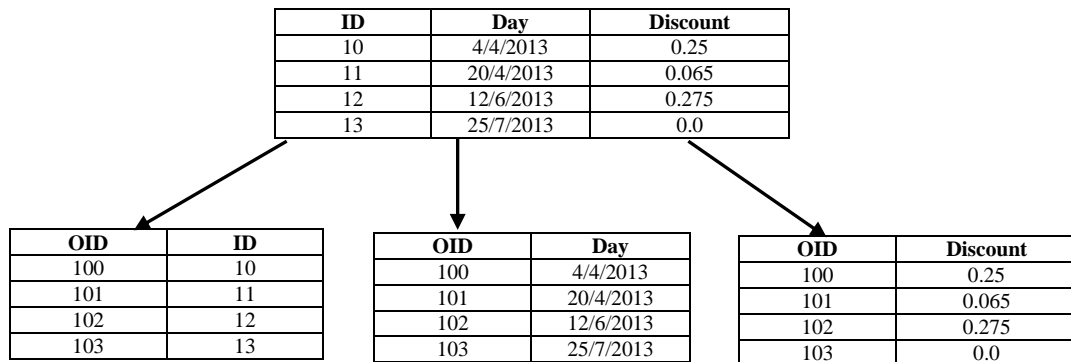


Table 1: simple Discount table (Figure 3)

We denote the *i*th projection over table *X* is *X_i*, followed by the name of the field in a projection. Attributes from the other table are represented by the name of table in which they come from. For example here we have a table EMP(name, age, salary, dept) DEPT(dname, floor). Following are the possible sets of projections of these tables are

- EMP1(name, age)
- EMP2(dept, age, DEPT.floor)
- EMP3(name, salary, DEPT.dname)
- EMP4(dname, floor, EMP.name)

Example 1: possible projections for EMP and DEPT

The Tuples in a projection are stored in column-wise. So if there is an *N* attributes in a projection, then there will be *N* data structure, each storing a single column, each of which is stored on the same *sort key*. The sort key can be any column or columns in a projection. Tuples in a projection are sorted on the key(s) in left to right manner.

A possible ordering for the above projection would be:

- EMP1(name, age | age)
- EMP2(dept, age, DEPT.floor | DEPT.floor)
- EMP3(name, salary, Dept.dname | salary)
- EMP4(dname, floor, EMP.name | floor)

Example 2: Projection in Example 1 with sort orders

Finally every projection is horizontally partitioned into 1 or more segments, and allot the segment identifier, Sid to every segment where Sid > 0. Column-oriented DBMS supports only value-based partitioning on the bases of sort key of a projection, every segment of a given projection is associated with a *key range* of the sort key. There in projection every column of every table is stored in at least one projection. Column-oriented DBMS must able to reconstruct the whole row of a table by the collection of stored segment. This is done by the joining of segment from different projections, which we accomplish using *storage keys* and *join indexes*.

Storage Keys: The storage key (SK) is associated with every data of every column with each segment. Values from different column in the same segment with machine store key belong to the same logical row. Storage keys are numbered 1, 2, 3, in RS and are not physically stored, but are inferred from a tuple's. storage key are physically stored in WS and are represented as integer, larger than the largest integer storage key for any segment in RS.

Join Indexes: To reconstruct all of the records in a table X from its numbers of projection Column-oriented DBMS uses join indexes. If X1 and X2 are two projections that cover a table X, a join index from the S segments in S1 to the Y segments in S2 is logically a collection of S tables, one per segment, S, of X1 consisting of rows of the form: (s: SID in X2, k: storage key in segment s). If we want to reconstruct X from the segment of X1, ..., Xk it must be possible to find a path through a set of join indices that maps each attribute of X into some sort order O*. A path is a collection of join indices originating with a sort order specified by some projection, X1, that passes through zero or more intermediate join indices and ends with a projection in sorted order. For Example, to reconstruct the EMP from its projection which is shown in example 2, we need at least two join indexes. Now we choose age as a common sort order, we could build two indices that map EMP2 and EMP3 to the ordering of EMP1. Alternatively, we could create a join index that maps EMP2 to EMP3 and one that maps EMP3 to EMP1. Figure 4 shows a simple example of a join index that maps EMP3 to EMP1, with single segment (SID =1) for each projection. For example, the first entry In EMP3(ram, 25,000), corresponding the second entry of EMP1, and first entry of the join index has store key 2.

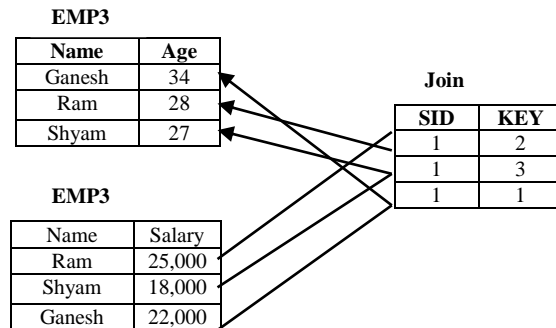


Figure 4: A join index from EMP3 to EMP1

The segment of the projections in a database and there connection join indexes must be allocated to the various node in a column-oriented system. The column-oriented administrator can optimally specify that the table in a database that must be K-safe. In this case, the loss of K-nodes in the grid will still allow all tables in a database to be reconstructed. When a failure occurs, column-oriented simply continues with K-1 safety until the failure is repaired and the node is brought back up to speed. Join indexes is used at the time when we want to use connect the various projection which is base at the same table. As noted earlier, a join index is a collection of (SID, storage_key) pairs. Each of these two fields can be stored as normal columns.

V. Writeable Store (WS)

WS is a column store and implements the physical DBMS design as RS. The same projection and join indexes are present in WS. The storage key, SK, for each record is explicitly stored in each WS segment. When a new logical tuple is inserted into table then a new unique SK is given to each tuple. This SK is an integer larger than the number of records in the largest segment in the database. There is a 1:1 mapping between RS segment and WS segment. A(SID, storage_key) pair identified a record In either of these contains. We assume that WS is trivial in size relative to RS, each projection uses B-tree indexing to maintain a logical sort-key order. Every column in a WS projection is represented as pair of collection, (v, sk), such that v is a value in the column and sk is its corresponding storage key. The sort key(s) of each projection is additionally represented by pairs (s, sk) such that s is a sort key and sk is the storage key describe where s first appears. Every projection is represented as a collection of pairs of segments, one in Ws and one in RS. For each record in the “sender”, must store the sid and storage key of a corresponding record in the “receiver”. it will be useful to horizontally partition the join index in the same way as the “sending” projection and then to co-locate join index.

VI. Storage Management

The storage management issue is the allocation of segment to nodes in a grid system; column-oriented will perform this operation automatically using storage allocator. As we now join indexes should be co-locator with their “sender” segment. Also, each WS segment will be co-located with the RS segment that contains the same key range. Everything is a column; storage is simply the persistence of a collection of columns. Our analysis show that a raw device offers little benefit relative to today’s file system. Hence big column (megabytes) is stored in individual files in the underlying operating system.

VII. Recovery

A crashed side recovers by running a query (copying state) from other projections. Column oriented maintained K-safety i.e sufficient projections and join indexes are maintained, so the K sites can fails with in t. the time to recover, and the system will be able to maintain transactional consistency. There are the cases to consideration first, if the failed side has no data loss then bring it up to date by executing updates that will be queued for it anywhere in network. Hence recovery from the most common type of crashes I straightforward. Second case is consider for catastrophic failure in this both the RS and WS are destroyed .in this case we have no option but only reconstruct both segment. And in third case occurs if WS is destroyed but RS not since RS is written only by the tuple mover. Hence, we discuss this common case in detail below.

VIII. Efficiently Recovering the WS

Suppose we have a WS segment, Sr, is a projection with a sort key SK and a key range R on a recovering site r along with a collection C of other projection, M1,...,Mb which contain the sort key of Sr. The tuple mover guarantees that each WS segment S, contains all tuples with an insertion timestamp later then some time $t_{lastmove}(S)$, Which represent the most recent insertion time of any record in S’s corresponding RS segment. For recovering site first inspect every projection in C for a collection of columns that covers the range of key K with each segment having $t_{lastmove}(S) \leq t_{lastmove}(Sr)$. it can run a collection of queries of the form

```

SELECT desired_fields
      insertion_epoch
      deletion_epoch
FROM   recovery_segment
WHERE  insertion_epoch > tlastmove(Sr)
      AND insertion_epoch <= HWM
      AND deletion_epoch >= 0
      OR deletion_epoch >= LWM
      AND sort_key in K

```

As long as the above query return a storage key, other fields in the segment can be found by following appropriate join index. As long as there is a collection of segments that cover the key ranges of Sr, this technique will restore Sr to the current HWM (High water mark is the maximum amount of database blocks used so far by a segment.) Executing queued updates will then complete the task.

IX. Tuple Mover

Tuple mover is important part of column-oriented database its move the blocks of tuples in a WS segment to the corresponding RS segment, it operates as a background task looking for related segment pairs. When it found any one, it perform a merge-our-process, MOP on this (RS, WS) segment pair. MOP find all records in the WS segment with an insertion time at or before the LWM, and then divides it into two parts

- At the time of deletion or insert the value before LWM are discarded, because the use cannot run queries because it is the time of execution
- If not deleted or delete after LWM these are moved to RS MOP will create a new RS segment which is identified by RS. It reads the block of columns of the RS segment, deletes the RS items which have values in the DRV less than or equal to the LWM, and merge these values with columns from WS. The merged data then written in the new RS’ segment. The time of most recent inserted is become a new $t_{lastmove}$ of RS’ segment and it always less than or equal to the LWM. This old-master/new-master approach is more efficient that the update-in-place approach, since essentially all data objects will move, maintenance of the DRV is also mandatory. Once RS’ contains all the WS data and join indexes are modified on RS’, the system cuts over from RS to Rs’. The disk space used by the old Rs can now be freed.

X. Column-oriented Query Execution

The query optimizer will accept a SQL query and construct plan of execution nodes. In this section we describe the nodes that can appear in a plan and then the architecture of the optimizer itself.

(A) *Query Operators and Plan Format*

There are 8 different types and each accepts operands or produce results of type projection (Proj) column (Col), or bitwise (Bits). A projection is simple with the same ordering. A bitstring is a list of zeros and ones indicating that the associated values are present in the record subset being described. Column oriented DBMS also accept predicates (pred), join indexes (JI), attribute name (Att) and expression (Exp) as argument. Here we summarize each operator as follows:

1. *Decompress*: convert a compressed column to an uncompressed (Type 4) representation.
2. *Select*: is equivalent to the selection operator of the relation algebra it produces a bitstring representation of the result
3. *Project*: equivalent to the projection operators of the relational algebra (Π).
4. *Sort*: Sort all columns in a projection by some subset of those column (the sort column)
5. *Aggregate operators*: Like SQL aggregate over a name column, and for each group identified by the values in a projection.
6. *Permute*: permutes a projection according to ordering defined by a join index.
7. *Join*: joins two projections according to a predicate that correlates them.
8. *Bitstring Operators*: Band produces the bitwise AND of two bitstring. Bor produces a bitwise OR. BNot produces the complement of a bitstring.

XI. Conclusions

This paper is presented the design of column store, a fundamental departure from the architecture of current DBMSs. Unlike current system, it is aimed at the read “read-mostly” DBMS market. The innovation contributions embodied in column include:

- A column store representation, with an associated query execution engine.
- A hybrid architecture that allow transaction on a column store.
- A focus on economizing the store representation on disk by coding data values dense-packing the data.
- A design optimized for a shared nothing machine environment.
- A distributed transaction without a redo log or two phase commit.
- Efficient snapshot isolation.

XII. References

- [1]. <http://www.monetdb.org/Documentation/Manuals/MonetDB/Architecture>
- [2]. http://www.sybase.in/files/White_Papers/Performance-Column-DBM-041509-WP.pdf
- [3]. Peter Boncz et.al .MonetDB/x100: Hyper-pipelining Query Execution. In *proceedings CIDR* 2004.
- [4]. Gray et al. data Cubes: A Relational Aggregation Operator Generalization Group- By, Cross-Tab, and Sub- Totals.Data Mining and Knowledge Discovery, 1(1) 1997.
- [5]. P. O’Neil and D. Quass, Improved Query Performance with variant indexes, In *Proceedings of SIGMOD, 1997*
- [6]. Oracle Corporation. *Oracle 9i database for data warehousing and Business Intelligence*. White Paper <http://www.oracle.com/solutions/>
- [7]. Paule Westerman. *Data Warehousing: Using the Wal-Mart Model*. Morgan-Kaufmann Publishers, 2000
- [8]. Harizopoulos, S., Liang, V., Abadi, D.J., and Madden, S.: Performance tradeoffs in read- optimized databases. In *Proc. VLDB, 2006*.
- [9]. Tsirogiannis, D., Harizopoulos, S., Shah, M.A., Wiener, J.L.,and Graefe, G.: Query processing techniques for solid state drives. In *Proc. SIGMOD, 2009*.
- [10]. Abadi, D.J., Myers, D.S., DeWitt, D.J., and Madden, S.R.: Materialization strategies in a column-oriented DBMS. In *Proc. ICDE, 2007*
- [11]. Zukowski, M., Heman, S., Nes, N., and Boncz, P.A.: Superscalar ram-cpu cache compression. In *Proc. ICDE, 2006*
- [12]. A. Ailamaki. “Database Architecture for New Hardware.”Tutorial. In *Proc. VLDB, 2004*.
- [13]. S. Agrawal, V. R. Narasayya, B. Yang. “Integrating Vertical and Horizontal Partitioning Into Automated Physical Database Design.” In *Proc. SIGMOD, 2004*.