

# Basic Building Blocks for Column-Stores

Andreas Schmidt<sup>\*†</sup>, Daniel Kimmig<sup>†</sup>, and Reimar Hofmann<sup>\*</sup>

<sup>\*</sup> Department of Computer Science and Business Information Systems,  
Karlsruhe University of Applied Sciences  
Karlsruhe, Germany

Email: {andreas.schmidt, reimar.hofmann}@hs-karlsruhe.de

<sup>†</sup> Institute for Applied Computer Science  
Karlsruhe Institute of Technology  
Karlsruhe, Germany

Email: {andreas.schmidt, daniel.kimmig}@kit.edu

**Abstract**—A constantly increasing CPU-memory gap as well as steady growth of main memory capacities have increased interest in column store systems due to potential performance gains within the realm of database solutions. In the past, several monolithic systems have reached maturity in the commercial and academic spaces. However, a framework of low-level and modular components for rapidly building column store based applications has yet to emerge. A possible field of application is the rapid development of high-performance components in various data-intensive areas such as text-retrieval systems and recommendation systems. The main contribution of this paper is a column-store-tool-kit, a basic building block of low-level components for constructing applications based on column store principles. We present a minimal amount of necessary structural elements and associated operations required for building applications based on our column-store-kit. The eligibility of our toolkit is demonstrated subsequently in using the components of our toolkit for building different query execution plans. This part of work is a first step in our effort for the construction of a pure column-store based query optimizer.

**Keywords**—Column store; basic components; framework; rapid prototyping; TPC-H benchmark; query-optimizer; query-execution plan.

## I. INTRODUCTION

Within database systems, values of a dataset are usually stored in a physically connected manner. A *row store* stores all column values of each single row consecutively (see Figure 1, bottom left). In contrast to that, within a *column store*, all values of each single column are stored one after another (see Figure 1, bottom right). In column stores, the relationship between individual column values and their originating datasets are established via Tuple IDentifiers (TID). The main advantage of column stores during query processing is the fact that only data from columns which are of relevance to a query have to be loaded. To answer the same query in a row store, all columns of a dataset have to be loaded, despite the fact, that only a small portion of them are actually of interest to the processing. On the other side, the column store architecture is disadvantageous for frequent changes (in particular insertions) to datasets. As the values are stored by column, they are distributed at various locations, which leads to a higher number of required write operations exceeding those within a row store to perform the same changes. This characteristic makes this type of storage interesting especially for applications with very

high data volume and few sporadic changes only (preferably as bulk upload), as it is the case in, e.g., data warehouses, business intelligence systems or text retrieval systems.

In our previous work [1], we identified the basic building blocks of our Column Store ToolKit (CSTK) and its interfaces with respect to providing a toolkit for building column store-based applications. In this paper, we extend our formulated ideas with a number of experiments which demonstrate the suitability of our toolkit with regard to building a query optimizer for column stores and the general suitability for scientific questions in the field of column store research.

Interest in column store systems has recently been reinforced by steady growth of main memory capacities that meanwhile allow for main memory-based database solutions and, additionally, by the constantly increasing CPU-memory gap [2]. Today's processors can process data much quicker than it can be loaded from main memory into the processor cache. Consequently, modern processors for database applications spend a major part of their time waiting for the required data. Column stores and special cache-conscious [3] algorithms are attempts to avoid this "waste of time". A number of commercial and academic column store systems have been developed in the past. In the research area, MonetDB [4] and C-Store [5] are widely known. Open Source and commercial systems include Sybase IQ, Infobright, Vertica, LucidDB, and Ingres. All these systems are more or less complete database systems with an SQL interface and a query optimizer.

As column stores are a young field of research, numerous aspects remain to be examined. For example, separation of datasets into individual columns result in a series of additional degrees of freedom when processing a query. Abadi et al. [6] developed several strategies as to when a result is to be "materialized", i.e., at which point in time result tuples shall be composed. Depending on the type of query and selectivity of predicates, an early or late materialization may be reasonable. Interesting studies were published about compression methods [7], various index types as well as the execution of join operations, e.g., Radix-Join [2], Invisible Join [8] or LZ-Join [9]. In addition to that, there are attempts at creating hybrid approaches that try to combine the advantages of column and row stores. The main objective of this paper is to present a number of low-level building blocks for constructing applications based on column store systems. Instead of copy-

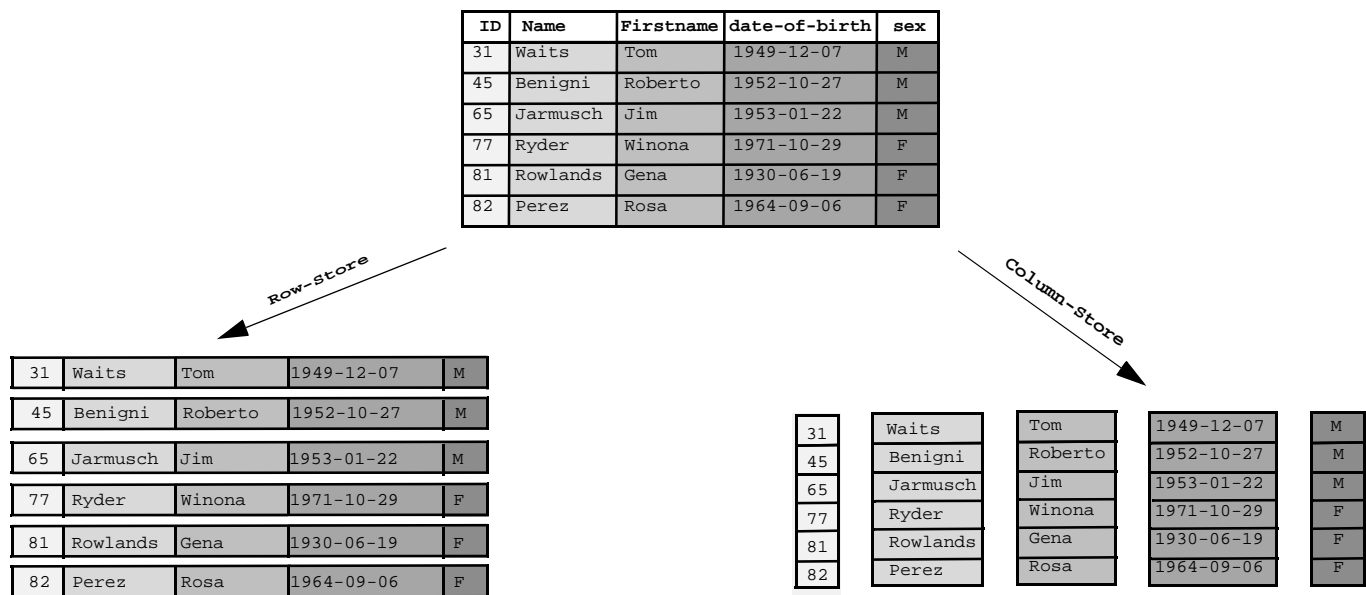


Fig. 1. Comparison of the layouts of a row store and a column store

ing the low-level constructs of existing sophisticated column stores, our research work is focused on identifying components and operations that allow for building specialized column store based applications in a rapid prototyping fashion. As our components can be composed in a “plug and compute”-style, our contribution is a column-store-tool-kit, which is a building block for experimental and prototypical setup of applications within the field of column stores. A possible field of application is the rapid development of high-performance components in various data-intensive areas such as text-retrieval systems.

This paper is structured as follows. In the next section, related work is mentioned. Then, in Section III, some basic considerations about column stores will be outlined. Afterwards, in Section IV and V the identified components and corresponding operations will be explained on a logical level. On this basis, various implementations of logical components and operations will be presented in Section VI. Finally, results will be summarized and an outlook will be given on future research activities.

## II. RELATED WORK

In the field of database systems, there are a number of related approaches. For example the C++-fastbit-library [10] provides a number of searching functions based on compressed bitmap indexes. Beside the low-level bitmap components, also a SQL interface exists in this library. The approach is comparable to the bitmap index in some relational database systems (i.e., Oracle, PostgreSQL). In contrast to these indexes, the fastbit bitmaps are compressed and therefore also usable for high cardinality attributes. The CSTK described in this paper can benefit from the compressed bitmap classes when implementing the *PositionLists* (see Section IV). Whether this implementation variant is advantageous depends on a number of factors. For details see [11]. In the field of query optimization there are a number of different tools, i.e., the Volcano project [12], developed by Goetz Graefe. Volcano is an optimizer generator, which means, that the source code

of the optimizer is generated, based on a model specification which consists of algebraic expressions. The library itself contains modules for a file-system, buffer management, sorting, duplicate elimination, *B+*-trees, aggregation, different join implementations, set operations, and aggregation functions. Based on the experiences gained with Volcano, the Cascades framework [13] was started, which later forms the base for the SQL Server 7.0 query optimizer [14].

## III. COLUMN STORE PRINCIPLES

Nowadays, modern processors utilize one or more cache hierarchies to accelerate access to main memory. A cache is a small and fast memory that resides between main memory and the CPU. In case the CPU requests data from main memory, it is first checked, whether it already resides within the cache. In this case, the item is sent directly from the cache to the CPU, without accessing the much slower main memory. If the item is not yet in the cache, it is first copied from the main memory to the cache and then further sent to the CPU. However, not only the requested data item, but a whole cache line, which is between 8 and 128 bytes long, is copied into the cache. This prefetching of data has the advantage, that requests to subsequent items are much faster, because they already reside within the cache. Meanwhile, the speed gain when accessing a dataset in the first-level cache is up to two orders of magnitude compared to regular main memory access [15]. Column stores take advantage of this prefetching behavior, because values of individual columns are physically connected together and, therefore, often already reside in the cache when requested, as the execution of complex queries is processed column by column rather than dataset by dataset. This also means that the decision whether a dataset fulfills a complex condition is generally delayed until the last column is processed. Consequently, additional data structures are required to administrate the status of a dataset in a query. These data structures are referred to as *Position Lists*. A *PositionList* stores the TIDs of matching datasets. Execution of a complex query generates

a *PositionList* with entries of the qualified datasets for every simple predicate. Then, the *PositionLists* are linked by *and/or* semantics. As an example, Figure 2 shows a possible execution plan for the following query:

```
select birthday, name
  from person
 where birthdate < '1960-01-01'
    and sex='F'
```

First, the predicates *birthdate < '1960-01-01'* and *sex = 'F'* must be evaluated against the corresponding columns (*birthdate* and *sex*), which results in the *PositionLists* *PL1* and *PL2*. These two evaluations could also be done in parallel. Next, an *and*-operation must be performed on these two *PositionLists*, resulting in the *PositionList* *PL3*. As we are interested in the *birthdate* and name of the persons that fulfil the query conditions, we have to perform another two operations (extract), which finally returns the entries for the TIDs, specified by the *PositionList* *PL3*.

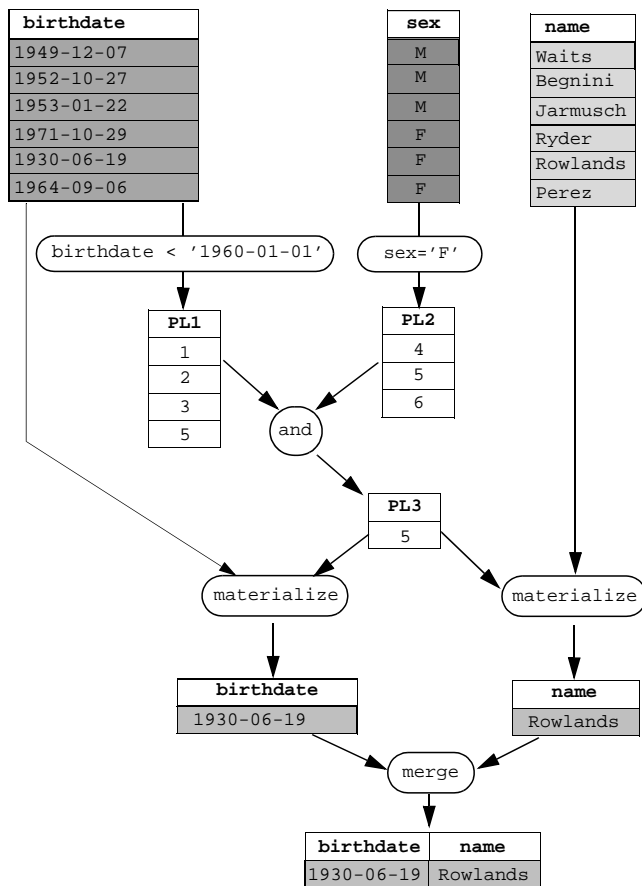


Fig. 2. Processing of a query with PositionLists

#### IV. CONCEPT

The main focus of our components is to model the individual columns, which can occur both in the secondary store as well as main memory. Their types of representation may vary. To store all values of a column, for example, it is not necessary to explicitly store the TID for each value, because it can be determined by its position (dense storage). To handle the results

of a filter operation however, the TIDs must be stored explicitly with the value (sparse storage). Another important component is the *PositionList* already mentioned in Section III. Just like columns, two different representation forms are available for main and secondary storage. To generate results or to handle intermediate results consisting of attributes of several columns, data structures are required for storing several values (so-called multi columns). These may also be used for the development of hybrid systems as well as for comparing the performance of row and column store systems. The operations mainly focus on writing, reading, merging, splitting, sorting, projecting, and filtering data. Predicates and/or *PositionLists* are applied as filtering arguments. Figure 3 illustrates a high level overview of the most important operations and transformations between the components. In Section V, they will be described in detail. Moreover, the components are to be developed for use on both secondary store and main memory as well as designed for maximum performance. This particularly implies the use of cache-conscious algorithms and structures.

#### V. PRESENTATION OF LOGICAL COMPONENTS

In the following sections, the aforementioned components will be presented together with their structure and their corresponding operations. Section VI will then outline potential implementations to reach highest possible performance.

##### A. Structure

1) *ColumnFile*: The *ColumnFile* serves to represent a column on the secondary storage. Supported primitive data types are: uint, int, char, date und float. Moreover, the composite type *SimpleStruct* (see V-A2) is supported, which may consist of a runtime definable list of the previously mentioned primitive data types. As a standard, the TID of a value in the *ColumnFile* is given implicitly by the position of the value in the file. If this is not the case, a *SimpleStruct* is used, which explicitly contains the TID in the first column.

2) *SimpleStruct*: *SimpleStruct* is a dynamic, runtime definable data structure. It is used within *ColumnFile* as well as within *ColumnArrays* (see below). The *SimpleStruct* plays a role in the following cases:

- Result of a filter query, in which the TIDs of the original datasets are also given.
- Combination of results consisting of several columns.
- Setup of hybrid systems having characteristics of both column and row stores. For example, it may be advantageous to store several attributes in a *SimpleStruct* that are frequently requested together.
- Representation of sorted columns, where TIDs are required. This is particularly reasonable for Join operators or a run-length-encoded compression on their basis.

3) *ColumnArray* and *MultiColumnArray*: A *ColumnArray* represents a column in main memory, which consists of a flexible number of lines. The data types correspond to those of the previously defined *ColumnFile*. If the data type is a *SimpleStruct*, it is referred to as *MultiColumnArray*. In addition to the actual column values, the TIDs of the first and last dataset

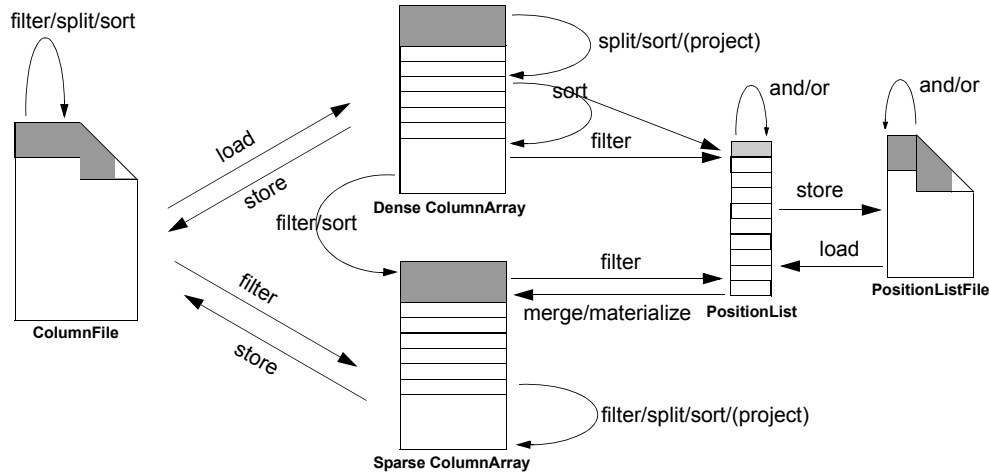


Fig. 3. Components and Operations

and the number of datasets stored are given in the header of the *(Multi)ColumnArray*. Two types of representations are distinguished:

- **Dense:** The type of representation is dense, if no gaps can be found in the datasets, i.e., if the TIDs are consecutive. In this case, the TID is given virtually by the TID of the first data set and the position in the array and does not have to be stored explicitly (see Figure 4, left side). This type of representation is particularly suited for main memory-based applications, in which all datasets (or a continuous section of them) are located in main memory.
- **Sparse:** This type of representation explicitly stores the TIDs of the datasets (see Figure 4, right). The primary purpose of a sparse *ColumnArray* is the storage of (intermediate) results. As will be outlined in more detail in Section V, it may be chosen between two physical implementations depending on the concrete purpose.

4) *ColumnBlocks* and *MultiColumnBlocks*: Apart from the *(Multi)ColumnArrays* of flexible size, *(Multi)ColumnBlocks* exist, which possess an arbitrary, but fixed size. They are mainly used to implement *ColumnArrays* with their flexible size. In addition, they may be applied in the implementation of a custom buffer management as a transfer unit between secondary and main memory and as a unit that can be indexed.

5) *PositionList*: A *PositionList* is nothing else than a *ColumnArray* with the data type *uint(4)* as far as structure is concerned. However, it has a different semantics. The *PositionList* stores TIDs. A *PositionList* is the result of a query via predicate(s) on a *ColumnFile* or a *(Multi)ColumnArray*, where the actual values are of no interest, but rather the information about the qualified data sets. *Position Lists* store the TIDs in ascending order without duplicates. This makes the typical and/or operations very fast, as the cost for both operations is  $O(|Pl_1| + |Pl_2|)$ . As will be outlined in Section VI, various types of implementations may be applied. Analogously to the *(Multi)ColumnArray*, there is a representation of the *PositionList* for the secondary store, which is called *PositionFile*.

Dense		Sparse	
StartPos:1024 EndPos :2047	StartPos:1024 EndPos :2047	StartPos : 1024 EndPos : 2047 Entries : 351	StartPos : 1024 EndPos : 2047 Entries : 351
	name sex	11	name sex
		21	
		45	
		51	
		89	
		93	
		..	
ColumnArray	MultiColumnArray	ColumnArray	MultiColumnArray

Fig. 4. Types of ColumnArrays

### B. Operations

1) *Transformations on ColumnFiles*: Several operations are defined on *ColumnFiles*. A filter operation (via predicate and/or *PositionList*) can be performed on a *ColumnFile* and the result can be written to another *ColumnFile* (with or without explicit TIDs). Other operations are the splitting of a *ColumnFile* as well as sorting among different criterias (see Section V-B6) with and without explicitly storing the TID.

2) *Transformations between ColumnFile and (Multi)-Column-Array*: *ColumnFiles* and *(Multi)ColumnArrays* are different types of representation of one or more logical columns. Physically, *ColumnFiles* are located in the secondary storage, while *ColumnArrays* are located in main memory. Consequently, both types of representations can also be transformed into each other using the corresponding operators.

A *ColumnFile* can be transformed completely or partially into a dense *(Multi)ColumnArray*. If not all, but only certain datasets that match special predicates or *PositionLists* are to be loaded into a *(Multi)ColumnArray*, this can be achieved using filter operations that generate a sparse *(Multi)ColumnArray*. A sparse *(Multi)ColumnArray* may also be transformed into a *ColumnFile*. In this case, the TIDs are stored explicitly in combination with the values. Other operations refer to the insertion of new values and the deletion of values. An outline of the most important operations of *ColumnFiles* is given in Table I.

TABLE I. OUTLINE OF OPERATIONS ON COLUMNFILES

Operation	Result type
read(ColumnFile)	ColumnArray (dense)
read(ColumnFile, start, length)	ColumnArray (dense)
filter(ColumnFile, predicate)	ColumnArray (sparse)
filter(ColumnFile, predicate-list)	ColumnArray (sparse)
filter(ColumnFile, positionlist)	ColumnArray (sparse)
filter(ColumnFile, positionlist-list)	ColumnArray (sparse)
filter(ColumnFile, predicate-list, positionlist-list)	ColumnArray (sparse)
fileFilter(ColumnFile, predicate)	ColumnFile (explicit TIDs)
fileFilter(ColumnFile, predicate-list)	ColumnFile (explicit TIDs)
fileFilter(ColumnFile, positionlist)	ColumnFile (explicit TIDs)
fileFilter(ColumnFile, positionlist-list)	ColumnFile (explicit TIDs)
fileFilter(ColumnFile, predicate-list, positionlist-list)	ColumnFile (explicit TIDs)
split(ColumnFile, predicate)	ColumnFile, ColumnFile
split(ColumnFile, position)	ColumnFile, ColumnFile
merge(ColumnFile-list, predicate-list)	MultiColumnArray (sparse),
sort(ColumnFile, column(s), direction)	ColumnFile
sort(ColumnFile, Orderlist)	ColumnFile
mapSort(ColumnFile)	ColumnFile, Orderlist
mapSort(ColumnFile)	ColumnArray, Orderlist
insert(ColumnFile, value)	Tupel-ID
insert(MultiColumnFile, value <sub>1</sub> , ...)	Tupel-ID
delete(ColumnFile, Tupel-ID)	boolean
delete(ColumnFile, Positionlist)	integer
delete(ColumnFile, predicate)	integer
delete(ColumnFile, predicate-list)	integer

3) *Operations on ColumnArrays*: Filter operations can be executed on (*Multi*)*ColumnArrays* using predicates and/or *PositionLists*. This may result in a sparse (*Multi*)*ColumnArray* or a *PositionList*. Furthermore, *ColumnArrays* may also be linked with each other by and/or semantics. If the (*Multi*)*ColumnArrays* have the same structure, the result also possesses this structure. The results correspond to the intersection or union of the original datasets. The result is a sparse (*Multi*)*ColumnArray*. If (*Multi*)*ColumnArrays* of differing structure are to be combined, only the and operation is defined. The result is a (*Multi*)*ColumnArray* that contains a union of all columns of the involved (*Multi*)*ColumnArrays* and returns the values for the datasets having identical TIDs. If the (*Multi*)*ColumnArrays* used as input are dense and if they have the same TID interval, the resulting *MultiColumnArray* is also dense. An outline of the most important operations of *ColumnArrays* is given in Table II. *ColumnArray* may also refer to a *MultiColumnArray*. A *MultiColumnArray*, however, only refers to the version having several columns.

4) *Transformation from PositionList to ColumnArray*: If the column values of the stored TIDs inside a *PositionList* are needed, an *extract* operation must be performed. Input to this operation is a *PositionList* as well as a *dense (multi) ColumnArray*. The result is a *sparse (Multi) ColumnArray*.

5) *Operations between PositionLists*: Several *PositionLists* may be combined by *and*, *or* semantics, with the result being a *PositionList*. The result list is sorted in ascending order corresponding to the TIDs. In addition, operations exist to load and store *PositionLists*. An outline of operations of *PositionLists* can be found in Table III.

6) *Sorting*: One basic operation on (*Multi*) *ColumnArrays* as well as *ColumnFiles* is sorting. Beside the obvious task to bring the result of a query in a specific order, sorting also plays an important role regarding performance considerations. For the elimination of duplicates, for join operations and for compression using run-length encoding, previous sorting can

TABLE II. OUTLINE OF OPERATIONS ON ColumnArrays

Operation	Result type
filter(ColumnArray, predicate)	ColumnArray (sparse)
filter(ColumnArray, predicate-list)	ColumnArray (sparse)
filter(ColumnArray, positionlist)	ColumnArray (sparse)
filter(ColumnArray, positionlist-list)	ColumnArray (sparse)
filter(ColumnArray, predicate-list, positionlist-list)	ColumnArray (sparse)
filter(ColumnArray, predicate)	PositionList
filter(ColumnArray, predicate-list)	PositionList
filter(ColumnArray, positionlist)	PositionList
filter(ColumnArray, positionlist-list)	PositionList
filter(ColumnArray, predicate-list, positionlist-list)	PositionList
and(ColumnArray, ColumnArray)	ColumnArray
or(ColumnArray, ColumnArray)	ColumnArray
and(ColumnArray, ColumnArray)	PositionList
or(ColumnArray, ColumnArray)	PositionList
project(MultiColumnArray, columns)	(Multi)ColumnArray
asPositionList(ColumnArray, column)	PositionList
split(ColumnArray, predicate)	ColumnArray (sparse) ColumnArray (sparse)
sort(ColumnArray)	ColumnArray
sort(ColumnArray, Orderlist)	ColumnArray
mapSort(ColumnArray)	ColumnArray, Orderlist
split(ColumnArray(dense), position) ColumnArray (dense)	ColumnArray (dense)
split(ColumnArray (sparse), position) ColumnArray (sparse)	ColumnArray (sparse)
merge(ColumnArray-list (sparse), predicate-list)	MultiColumnArray (sparse)
store(ColumnArray (dense))	ColumnFile
store(ColumnArray (sparse))	ColumnFile (explicit TIDs)

TABLE III. OUTLINE OF OPERATIONS ON PositionLists

Operation	Result type
load(ColumnFile)	PositionList
store(PositionList)	ColumnFile
and(PositionList, PositionList)	PositionList
or(PositionList, PositionList)	PositionList
materialize(PositionList, ColumnArray, ...)	ColumnArray
materialize(PositionList, ColumnFile, ...)	ColumnArray
read(PositionListFile)	PositionList
store(PositionList)	PositionListFile

dramatically improve performance. As a consequence of sorting, the natural order is lost. This is critical for dense columns with implicit TIDs, because the relation to the other column values is lost. The problem can be solved by an additional data structure, similar to a *PositionList* that contains the mapping information to the original order of the datasets. Figure 5 gives an example of this situation. The *Multi ColumnArray* on the left side is to be sorted according to the column “name”. Additionally to the sorting of the *MultiColumn* (top right), a list is generated which holds the information about the original positions (down right). The list can then be reused by applying it as a sorting criterion to other columns later, as shown in Figure 6.

7) *Compression*: Compression plays an important role in column stores [7], as it reduces the data volume that needs to be loaded. Nevertheless, we decided not to include compression in the first prototype and to concentrate on the interfaces of the components. To a certain extent, this constraint can be compensated by the use of dictionary-based compression [16], which will be implemented above the basic components. In later versions, various compression methods will be integrated, so first of all run-length encoding (RLE) [17].

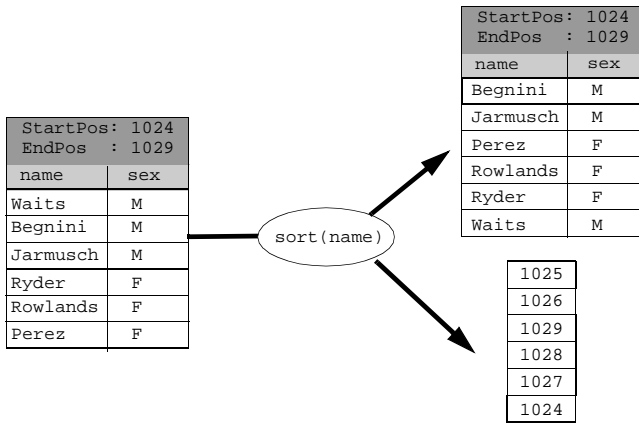


Fig. 5. Sorting with explicit generation of an additional mapping list

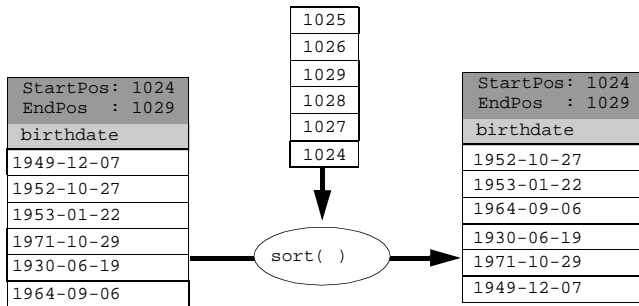


Fig. 6. Sorting with explicit given sort-order

VI. IMPLEMENTATION-SPECIFIC CONSIDERATIONS

After presenting the logical structure and the required operations, this section will now focus on considerations for achieving a performance-oriented implementation. Due to the constantly increasing CPU-memory gap, cache-conscious programming is indispensable. For this reason, the implementation was made in C/C++. All time-critical parts were implemented in pure C using pointer arithmetics. The uncritical parts were implemented using C++ classes. The *ColumnBlock* was established as a basic component of the implementation. It is the basic unit for data storage. Its size is defined at creation time and it contains the actual data as well as information on its structure and the number of datasets. The structurization options correspond to those of the *(Multi)ColumnArray*. The *ColumnBlock* also handles all queries by predicates and/or *PositionLists*. A *(Multi)ColumnArray* consists of *n* *ColumnBlock* instances. All operations on a *(Multi)ColumnArray* are transferred to the underlying *ColumnBlocks*.

*PositionLists* play a central role in column store applications. One important point is the size of a *PositionList*. If the *PositionLists* are short (i.e., if they contain a few TIDs only), representation as *ColumnArray* is ideal. Four bytes are required per selected entry. If the lists are very large, however, memory of 400 MB is required for ten million entries, for instance. In this case, a bit vector is recommended for representation. This bit vector uses for each dataset a bit at a fixed position to indicate whether a dataset belongs to the set of results or not. If, for example, 10 million data sets exist for a table, only 1.25 MB are required to represent the *PositionList* for

certain selectivities. Moreover, the two important operations *and* and *or* can be mapped on the respective primitive processor commands, which makes the operations extremely fast. If *PositionLists* are sparse, bit vectors can be compressed very well using run-length encoding (RLE) (e.g., to a few KB in case of 0.1% selectivity). The necessary operations can be performed very efficiently on the compressed lists, which further increases the performance. An implementation based on the word-aligned hybrid algorithm [18] with satisfactory compression for medium-sparse representations was developed within the framework of the activities reported here [19], [20].

Figure 7 gives an overview of the memory consumption for different implementations of a *PositionList*. Here, we compare the behavior of a dynamic array containing 4-byte TIDs with a plain uncompressed bitvector and different implementations (32, 64 bit) of the Word Aligned Hybrid (WAH) algorithm [18], both compressed and uncompressed. As we can see in the figure, the behavior of the dynamic array implementation is quite good for very small selectivities, but changes for the worse for medium and high densities. Uncompressed bitmaps (plain bitvector or WAH uncompressed) behave independently for all densities. Their size is determined by the number of tuples in a table only. Compressed bitmaps show a very good behavior for all densities. If selectivities become low, they behave like uncompressed bitmaps (compared to a pure uncompressed implementation of a bitvector, there will be a slight overhead of 1/32. resp. 1/64.). From a selectivity of about 3%, the array has a higher memory consumption than the uncompressed bitvector. Beside the memory consumption, also the runtime behavior of the different implementation variants plays a very important role. In [21], an elaborate analysis of the memory consumption and runtime behavior of different implementation variants (array, bitvector, compressed bitvector) for positionlists can be found. The bottom line of this paper is that the choice of the right implementation variant is not a trivial task and depends heavily on the selectivity of the predicates. The differences in the runtime behavior are over two orders of magnitude for typical *PositionList* operations.

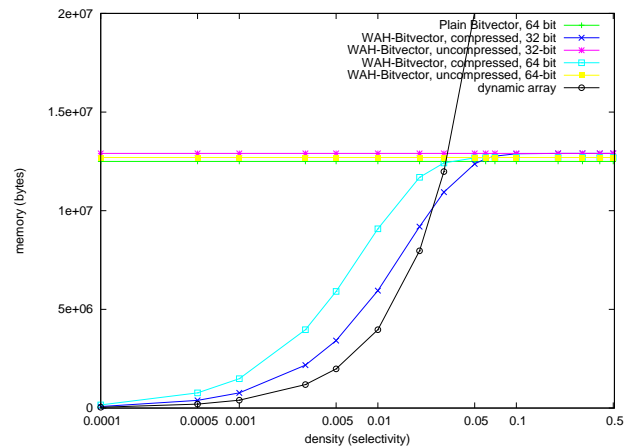


Fig. 7. Comparison of the memory consumption for different implementation variants of *PositionLists*

*MultiColumnArrays* may exist in two different physical layouts. In the first version, the *n* values are written in a physically successive manner and correspond to the classical *n*-ary storage model (NSM). This type of representation is

particularly suited, if further queries are to be performed on this *MultiColumnArray* with predicates on the respective attributes. The individual values of a dataset are stored together in the cache and all attribute values are checked simultaneously rather than successively with the help of additional *PositionLists* (see Figure 8, left). The second type of representation corresponds to the PAX format [22]. Here, every column is stored in a separate *ColumnArray*. In addition, a *PositionList* is stored, which identifies the datasets (see Figure 8, right). This type of representation is recommended, for instance, for collecting values for subsequent aggregation functions. Several (*Multi*)*ColumnArrays* may share a single *PositionList*.

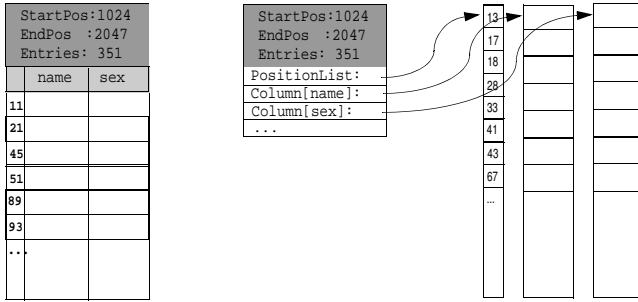


Fig. 8. Comparison of storage formats for ColumnArrays

## VII. USE CASES

In this section, we want to deal with the usage of the CSTK-components. We will present the general mechanism for building complex queries from the components, demonstrate the suitability of our components for scientific questions in the field of column store research and present an execution plan and the runtime behavior for a typical data warehouse query from the TPC-H [23] benchmark. The aim of this experiment is to gain further insight into the costs of the different operations and to derive rules for a query optimizer for column stores [24].

### A. Usability of the Components

1) *Materialization*: In [6], Abadi et al. propose different strategies to construct the final result sets from the intermediate *PositionLists*. This step is called “materialization”. One strategy is to keep the *PositionList* values as long as possible and to only materialize the attribute values in a very last step. This is called “late materialization”. On the other hand, “early materialization” means that the values should already be extracted in every selection step. The quintessence of the paper is that the superiority of any strategy depends on the characteristic of the query.

In the paper, Abadi et al. identified four different datasource operators (DS1, .., DS4) from which data could be read from disk or main memory. Additionally, they identified the *AND* operator for *PositionLists* and two more tuple construction operators, *MERGE* and *SPC* (Scan, Predicate, and Construct) for the construction of result tuples.

Based on these operators, they formed different query plans to implement early and late materialization strategies. Figure 9 shows the different execution plans for the following query, implementing an early materialization strategy (a, b) or a late materialization strategy (c, d).

```
SELECT l_shipdate, l_linenum
FROM lineitem
WHERE l_shipdate < C1
AND l_linenum < C2
```

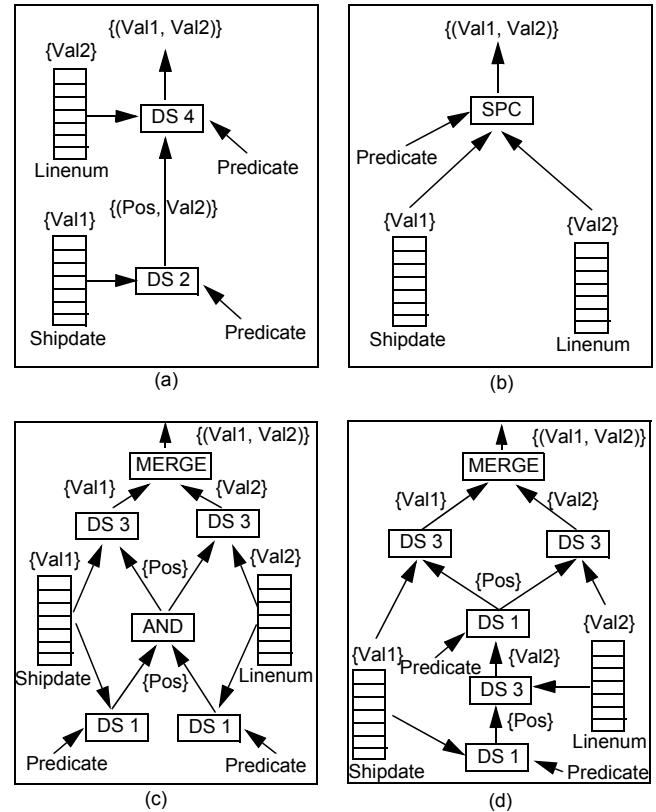


Fig. 9. Different query-plans from [6]

Using the components from the CSTK, these query plans can easily be rebuilt, using the operations from Tables I, II, and III. This is shown in Figure 10. In contrast to the original execution plans, which do not distinguish between file and main memory representation in each case, this is done with the execution plans built with the CSTK.

2) *Complex Queries*: In the following, a step-by-step explanation of a join operation is performed based on an example. The underlying dataset is from TPC-H benchmark (*lineitem* and *partkey* table).

The SQL query is the following:

```
SELECT p_name, l_quantity
FROM part
JOIN lineitem
ON p_partkey = l_partkey
WHERE l_orderkey = 34
```

Figure 11 shows the corresponding operations on the required columns. First, the *WHERE*-clause on the *l\_orderkey* column is executed (1) to get the corresponding TIDs (*l\_TID*) from the *lineitem* table. The extracted TIDs (5,6,7) are then used to read the corresponding values (883, 894, 169) from the *l\_partkey* column of the *lineitem* table (2). Next, the (*l\_TID*, *l\_partkey*) tuples are sorted based on their *l\_partkey*

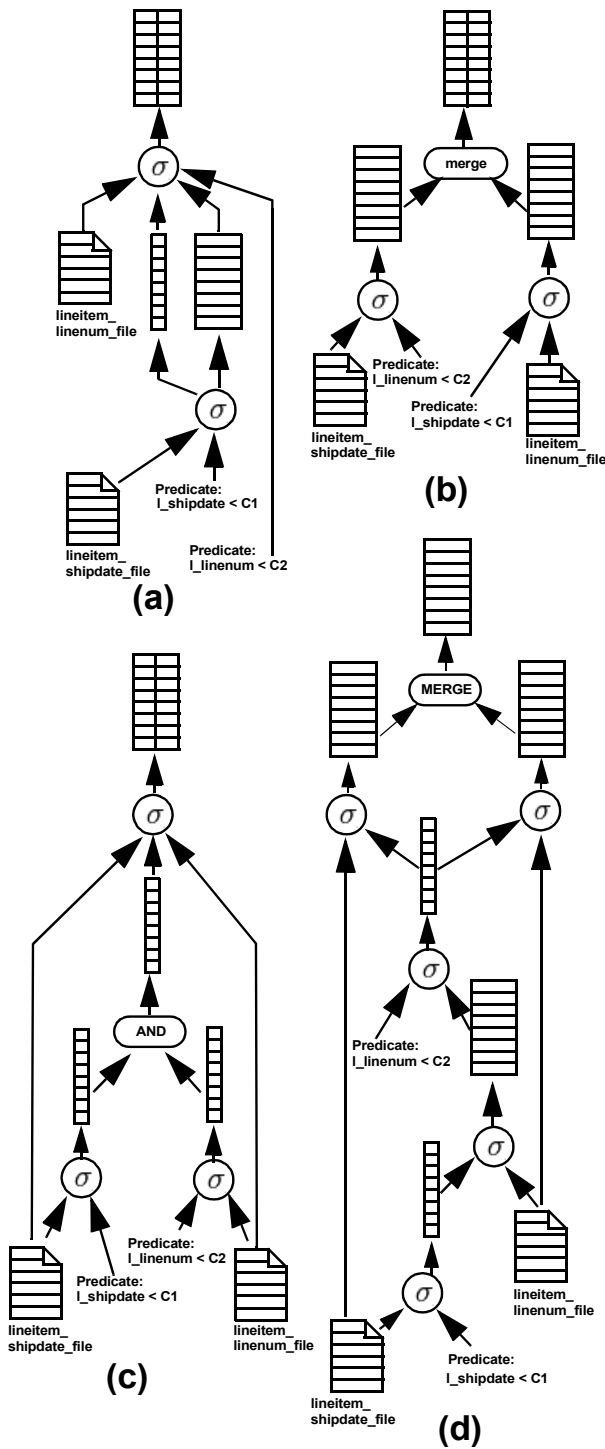


Fig. 10. Different materialization strategies from [6] using the CSTK components

values (3). The resorted tuples can then be merged with the sorted  $p\_partkey$  column of the  $partkey$  table (5), which has to be sorted priorly (4) and enriched with the  $p\_TID$  column, which was implicitly given by the position of the values in the unsorted  $p\_partkey$  column.

The result of the merge operation are tuples of the form  $(l\_TID, p\_TID)$ . They represent the result of the join operation

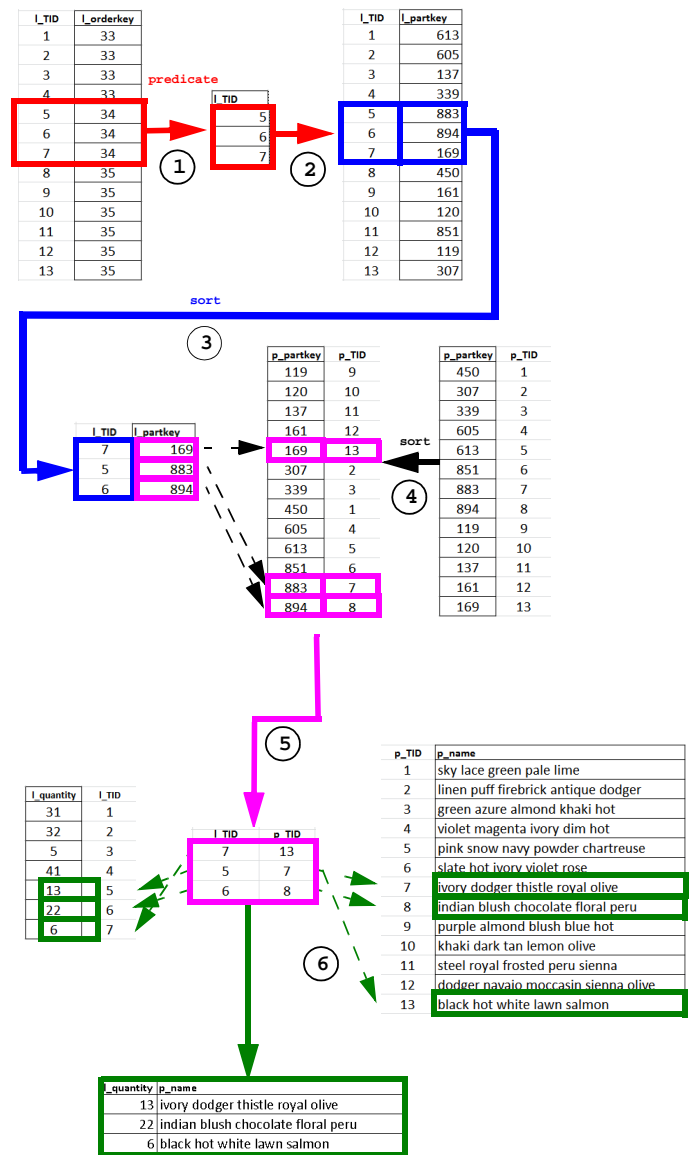


Fig. 11. Join-Operation with the Column-Store-Tool-Kit

between the  $lineitem$  and  $partkey$  table on the  $partkey$  column. In the last step, the materialization (6) takes place. The  $l\_TID$  and  $p\_TID$  values are replaced by their corresponding values from the  $p\_name$  and  $p\_quantity$  columns.

After demonstration of a CSTK-Join on a concrete example, the principle data flows, based on the operations on Tables I, II, and III are shown. Figure 12 shows an execution plan performing the following SQL query:

```
SELECT *
FROM orders o
JOIN lineitem l
ON l_orderkey=o_orderkey
```

In the current execution plan, a sort-merge join is performed. As a first step, the entries in the two column files  $orders\_orderkey\_file$  and  $lineitem\_orderkey\_file$  must be sorted (remember: in the files, the TIDs are implicit given by the



position of the values in the file). This is done with the *mapSort*-operation. The *mapSort* operation sorts the column values and provides an additional data structure *pl\_o* and *pl\_l*, which contains the TIDs for each sorted value. The structure is similar to a *PositionList*, but the TIDs are no longer sorted.

After the preparatory sorting step, the values in the columns are compared position by position (operation *cmp*). For each matching value from the two columns, the corresponding entries in the previously generated *PositionList* *pl\_o* and *pl\_l* are taken and written into the joined *PositionList* (*pl\_o',pl\_l'*). In a final step (not shown in Figure 12), the joined *PositionList* is materialized.

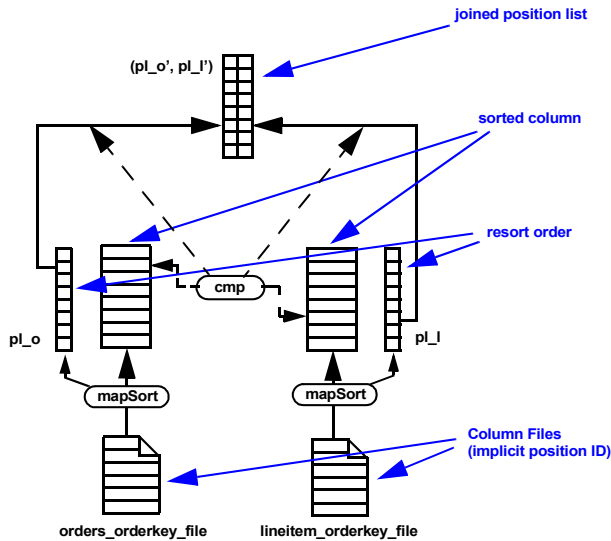


Fig. 12. Join-Operation with the Column-Store-Tool-Kit

An additional *WHERE* clause (see below) leads to the execution plan in Figure 13.

```
SELECT *
FROM orders o
JOIN lineitem l
ON l_orderkey=o_orderkey
WHERE o_orderdate= '1992-01-13'
```

The evaluation of the condition on the *orders\_orderdate\_file* generates a *PositionList* (*pl\_o*), which acts as a filter criterion for the *orders\_orderkey\_file*. After filtering, the *PositionList* also represents the TIDs for the *orders\_orderkey* column. In the subsequent *mapSort* operation, the *orders\_orderkey* column is resorted along its values and the corresponding TIDs in the *PositionList* *pl\_o* get resorted, respectively (*pl\_o'*). The rest of the join operation is similar to that already described in Figure 12.

### B. Performance Tests

To complete our case study concerning our toolkit, we present a more complex query from the TPC-H repository (Query 3). We model an execution plan using our components and run some performance tests, which we compare with MySQL and Infobright.

The SQL query we use is the following:

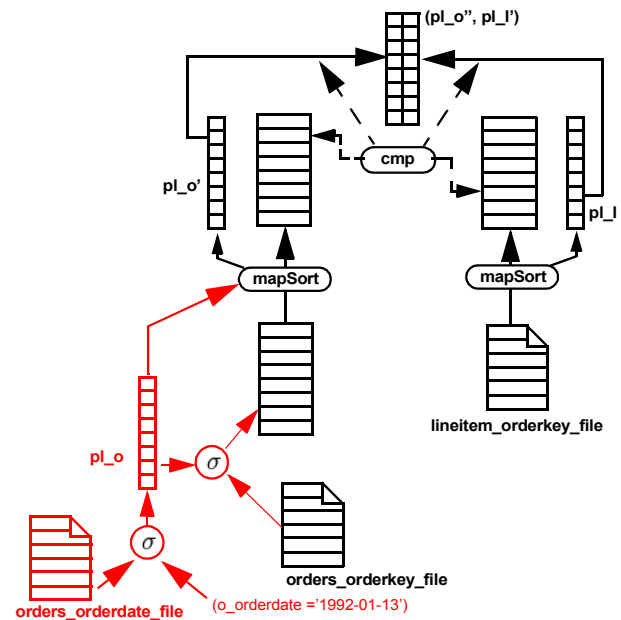


Fig. 13. Join-Operation with the Column-Store-Tool-Kit

```
select l_orderkey,
       sum(l_extendedprice*(1-l_discount))
       as revenue,
       o_orderdate,
       o_shippriority
from customer,
orders,
lineitem
where c_mktsegment = 'BUILDING'
and c_custkey = o_custkey
and l_orderkey = o_orderkey
and o_orderdate < date '1995-03-15'
and l_shipdate > date '1995-03-15'
group by l_orderkey,
o_orderdate,
o_shippriority
order by revenue desc,
o_orderdate
```

A possible corresponding execution plan for this query using late materialization is shown in Figure 14. Beside the used operations and the intermediate results, shown. The input consists of about 6 million lineitem tuples, 727 thousand orders and over 30 thousand customers from the TPC-H benchmark dataset. The machine settings are the following: Intel® Core™ i7-3520M CPU, 2.9 GHz processor with 2 physical cores, 8 GB main memory, running Windows 7 Enterprise, 64 bit. The cache sizes are: First level cache: 128KB, second-level cache: 512KB, third-level cache: 4MB.

The operation mainly consists of a join over the three tables and a subsequent grouping according to three columns. The overall execution time is about 1.107 sec. About 20% of the overall time is spent reading the needed columns from file and performing the selections based on predicates or *PositionLists*. The most expensive operations are the *mapSort*-operations, which take about 25% of the execution time. The subsequent sorting of the corresponding *PositionLists* takes another 15%.

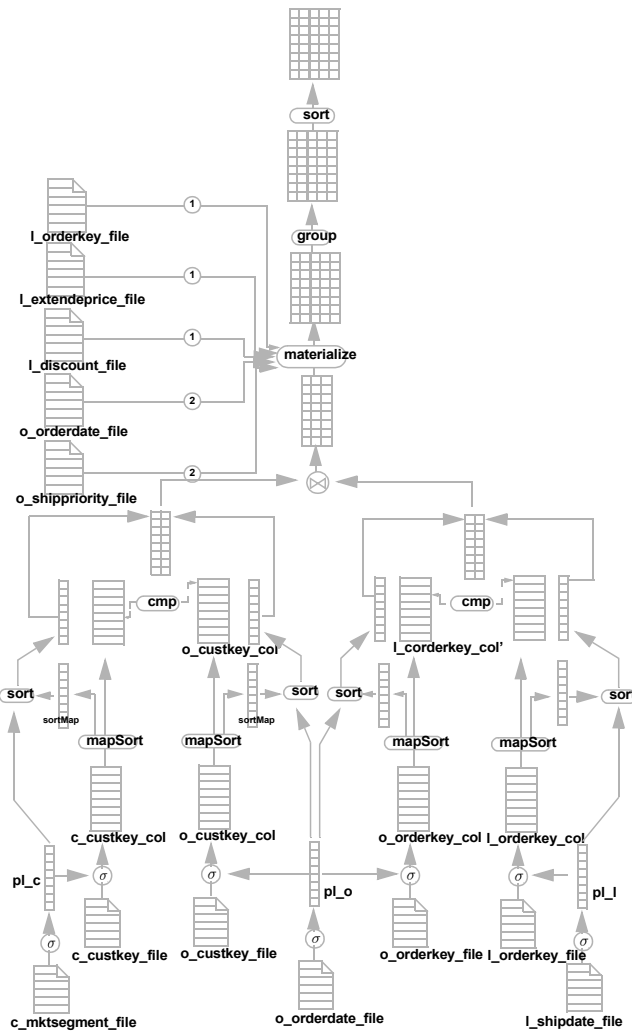


Fig. 14. TPC-H Query 3, execution plan and time behavior with CSTK components

Currently, we use a standard quicksort implementation without any optimizations. By exchanging the sorting algorithm with a more sophisticated version, we expect a further improvement of the runtime behavior. After sorting of the columns, we can use a merge-join implementation, which performs its task in about 0.03 seconds for an input cardinality of over 3.2 million tuples (lineitem datasets) and 727 thousand tuples (order datasets).

About one third of the complete execution time is spent accessing files on disk. Using a main-memory implementation could further reduce the overall execution time significantly. In comparison, the execution time of the same query using MySQL (with indexes on all foreign keys as well as on the columns which are predicated) takes about 116 seconds (cold start) with empty cache and about 13 seconds for repeated executions. Infobright [25], a column store-based version of MySQL, takes about 3 seconds to execute the query.

### VIII. CONCLUSION

This paper presented a collection of basic components to build column store applications. The components are semanti-

cally located below those of the existing column store database implementations and are suited for building experimental (distributed) systems in the field of column store databases.

As a proof of concept, we used these components to retrace the materialization experiments carried out by Abadi et al. [6]. Additionally, we show that typical operations like joining tables and grouping results can be carried out. Finally, we construct an execution plan from the TPC-H benchmark and point out that the performance is quite good, compared to existing column store databases. It is planned to use these components to obtain further scientific findings in the area of column stores and to develop data-intensive applications.

### IX. FUTURE WORK

A first version of the column store tool kit is available without support for compression. The next steps planned are the integration of compression and the use in concrete areas, such as text retrieval systems. A future activity will be the implementation of a scripting language interface for the components. With the help of this interface, it will be possible to assemble the developed components more easily without losing the performance of the underlying C/C++ implementation. In this case, the scripting language act as glue between the components, allowing the developer to build up complex high performance applications with very little effort [26]. As an alternative, a custom domain-specific language (DSL) [27] may be used for building column store applications. A bachelor's thesis [28] focused on the extent to which various degrees of flexibility regarding the structure of *MultiColumnArrays* and expression of the predicates affect the performance. According to the thesis, the structural definition at compilation time is of significant advantage compared to the structural definition at runtime. If the implemented flexibility of the *SimpleStruct* is not required at runtime, an alternative implementation may be used. It may be realized by defining a language extension for C/C++, for example. Thus, the respective structures and operations can be defined using a simple syntax. With a number of macros of the C++ preprocessor or a separate inline code expander [29], these could then be transformed into valid C/C++ code.

### REFERENCES

- [1] A. Schmidt and D. Kimmig, "Basic components for building column store-based applications," in *DBKDA'12: Proceedings of the The Forth International Conference on Advances in Databases, Knowledge, and Data Applications*. iaria, 2012, pp. 140–146.
- [2] S. Manegold, P. A. Boncz, and M. L. Kersten, "Optimizing database architecture for the new bottleneck: memory access," *The VLDB Journal*, vol. 9, no. 3, pp. 231–246, 2000.
- [3] T. M. Chilimbi, B. Davidson, and J. R. Larus, "Cache-conscious structure definition," in *PLDI '99: Proceedings of the ACM SIGPLAN 1999 conference on Programming language design and implementation*. New York, NY, USA: ACM, 1999, pp. 13–24.
- [4] P. A. Boncz, M. L. Kersten, and S. Manegold, "Breaking the memory wall in monetdb," *Commun. ACM*, vol. 51, no. 12, pp. 77–85, 2008.
- [5] M. Stonebraker, D. J. Abadi, A. Batkin, X. Chen, M. Cherniack, M. Ferreira, E. Lau, A. Lin, S. Madden, E. O'Neil, P. O'Neil, A. Rasin, N. Tran, and S. Zdonik, "C-store: a column-oriented dbms," in *VLDB '05: Proceedings of the 31st international conference on Very large data bases*. VLDB Endowment, 2005, pp. 553–564.
- [6] D. J. Abadi, D. S. Myers, D. J. Dewitt, and S. R. Madden, "Materialization strategies in a column-oriented dbms," in *In Proc. of ICDE*, 2007.

- [7] D. J. Abadi, S. R. Madden, and M. Ferreira, "Integrating compression and execution in column-oriented database systems," in *SIGMOD*, Chicago, IL, USA, 2006, pp. 671–682.
- [8] D. J. Abadi, S. R. Madden, and N. Hachem, "Column-stores vs. row-stores: How different are they really," in *In SIGMOD*, 2008.
- [9] L. Gan, R. Li, Y. Jia, and X. Jin, "Join directly on heavy-weight compressed data in column-oriented database," in *WAIM*, 2010, pp. 357–362.
- [10] K. Wu, "Fastbit reference manual," Scientific Data Management Lawrence Berkeley National Lab, Tech. Rep. LBNL/PUB-3192, august 2007. [Online]. Available: <http://lbl.gov/%7Ekwu/ps/PUB-3192.pdf>
- [11] A. Schmidt and D. Kimmig, "Considerations about implementation variants for position lists," in *DBKDA'13: Proceedings of the The Fifth International Conference on Advances in Databases, Knowledge, and Data Applications*. iaria, 2013, pp. 108–115.
- [12] G. Graefe and W. J. McKenna, "The volcano optimizer generator: Extensibility and efficient search," in *Proceedings of the Ninth International Conference on Data Engineering, April 19-23, 1993, Vienna, Austria*. IEEE Computer Society, 1993, pp. 209–218.
- [13] G. Graefe, "The cascades framework for query optimization," *IEEE Data Eng. Bull.*, vol. 18, no. 3, pp. 19–29, 1995.
- [14] B. Nevarez, *Inside the SQL Server Query Optimizer*. United Kingdom: Red gate books, 2011.
- [15] P. A. Boncz, M. Zukowski, and N. Nes, "Monetdb/x100: Hyper-pipelining query execution," in *CIDR*, 2005, pp. 225–237.
- [16] C. Binnig, S. Hildenbrand, and F. Färber, "Dictionary-based order-preserving string compression for main memory column stores," in *SIGMOD '09: Proceedings of the 35th SIGMOD international conference on Management of data*. New York, NY, USA: ACM, 2009, pp. 283–296.
- [17] S. Smith, *The scientist and engineer's guide to digital signal processing*. California Technical Publishing, 1997.
- [18] K. Wu, E. J. Otoo, and A. Shoshani, "Optimizing bitmap indices with efficient compression," *ACM Trans. Database Syst.*, vol. 31, no. 1, pp. 1–38, 2006.
- [19] A. Schmidt and M. Beine, "A concept for a compression scheme of medium-sparse bitmaps," in *DBKDA'11: Proceedings of the The Third International Conference on Advances in Databases, Knowledge, and Data Applications*. iaria, 2011, pp. 192–195.
- [20] M. Beine, "Implementation and Evaluation of an Efficient Compression Method for Medium-Sparse Bitmap Indexes," Bachelor Thesis, Department of Informatics and Business Information Systems, Karlsruhe University of Applied Sciences, Karlsruhe, Germany, 2011.
- [21] A. Schmidt and D. Kimmig, "Considerations about implementation variants for position-lists," in *DBKDA'13: Proceedings of the Fifth International Conference on Advances in Databases, Knowledge, and Data Applications*, 2013.
- [22] A. Ailamaki, D. J. DeWitt, and M. D. Hill, "Data page layouts for relational databases on deep memory hierarchies," *The VLDB Journal*, vol. 11, no. 3, pp. 198–215, 2002.
- [23] "TPC Benchmark H Standard Specification, Revision 2.1.0," Transaction Processing Performance Council, Tech. Rep., 2002.
- [24] A. Schmidt, D. Kimmig, and R. Hofmann, "A first step towards a query optimizer for column-stores," Poster presentation at the Forth International Conference on Advances in Databases, Knowledge, and Data Applications, DBKDA'12, Saint Gilles, Reunion, 2012.
- [25] D. Ślezak and V. Eastwood, "Data warehouse technology by infobright," in *Proceedings of the 2009 ACM SIGMOD International Conference on Management of data*, ser. SIGMOD '09. New York, NY, USA: ACM, 2009, pp. 841–846.
- [26] J. K. Ousterhout, "Scripting: Higher-Level Programming for the 21st Century," *IEEE Computer*, vol. 31, no. 3, pp. 23–30, 1998.
- [27] M. Mernik, J. Heering, and A. M. Sloane, "When and how to develop domain-specific languages," *ACM Comput. Surv.*, vol. 37, no. 4, pp. 316–344, 2005.
- [28] M. Herda, "Entwicklung eines Baukastens zur Erstellung von Column-Store basierten Anwendungen Bachelor's thesis, Department of Informatics, Heilbronn University of Applied Sciences, Germany," Jun. 2011.
- [29] J. Herrington, *Code Generation in Action*. Greenwich, CT, USA: Manning Publications Co., 2003.