

Holistic Indexing in Main-memory Column-stores

Eleni Petraki
CWI Amsterdam
petraki@cwi.nl

Stratos Idreos
Harvard University
stratos@seas.harvard.edu

Stefan Manegold
CWI Amsterdam
manegold@cwi.nl

ABSTRACT

Great database systems performance relies heavily on index tuning, i.e., creating and utilizing the best indices depending on the workload. However, the complexity of the index tuning process has dramatically increased in recent years due to ad-hoc workloads and shortage of time and system resources to invest in tuning.

This paper introduces *holistic indexing*, a new approach to automated index tuning in dynamic environments. Holistic indexing requires zero set-up and tuning effort, relying on adaptive index creation as a side-effect of query processing. Indices are created incrementally and partially; they are continuously refined as we process more and more queries. Holistic indexing takes the state-of-the-art adaptive indexing ideas a big step further by introducing the notion of a system which never stops refining the index space, taking educated decisions about which index we should incrementally refine next based on continuous knowledge acquisition about the running workload and resource utilization. When the system detects idle CPU cycles, it utilizes those extra cycles by refining the adaptive indices which are most likely to bring a benefit for future queries. Such idle CPU cycles occur when the system cannot exploit all available cores up to 100%, i.e., either because the workload is not enough to saturate the CPUs or because the current tasks performed for query processing are not easy to parallelize to the point where all available CPU power is exploited.

In this paper, we present the design of holistic indexing for column-oriented database architectures and we discuss a detailed analysis against parallel versions of state-of-the-art indexing and adaptive indexing approaches. Holistic indexing is implemented in an open-source column-store DBMS. Our detailed experiments on both synthetic and standard benchmarks (TPC-H) and workloads (SkyServer) demonstrate that holistic indexing brings significant performance gains by being able to continuously refine the physical design in parallel to query processing, exploiting any idle CPU resources.

Categories and Subject Descriptors

H.2.2 [Database Management]: Physical Design

Keywords

Holistic Indexing, Self-organization

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
SIGMOD'15, May 31–June 4, 2015, Melbourne, Victoria, Australia.
Copyright 2015 ACM 978-1-4503-2758-9/15/05\$15.00.
<http://dx.doi.org/10.1145/2723372.2723719>.

1. INTRODUCTION

The big data era is causing the research community to rethink fundamental issues in the design of database systems towards more usable systems [32] that can access data better and faster [36, 37, 39, 40], that can better exploit modern hardware and opportunities for massive parallelization [15] that can support efficient processing of OLTP and/or OLAP queries [6, 34, 35, 41].

The Physical Design Problem. Physical design, and in particular proper index selection, is a predominant factor for the performance of database systems and has only become more crucial in the big data era. With new dynamic and exploratory environments, physical design becomes especially hard given the instability of workloads and the continuous stream of big data; a single physical design choice is not necessarily correct or useful for long stretches of time, while at the same time workload knowledge is scarce given the exploratory user behavior.

State-of-the-Art. In database applications, where “the future is known”, e.g., in scientific databases, social networks, web logs, etc. In particular, the query processing patterns follow an exploratory behavior, which changes so arbitrarily that it cannot be predicted. Such environments cannot be handled by offline indexing. Online indexing [10, 47] and adaptive indexing [27] are two approaches to automatic physical design in such dynamic environments, but none of them in isolation handles the problem sufficiently. Online indexing periodically refines the physical design but it may negatively affect running queries every time it needs to use resources for reconsidering the physical design and it may not be quick to follow the workload changes as it reacts only periodically. Adaptive indexing does not have this problem as it introduces continuous, incremental and partial index refinement but it adjusts the physical design only during query processing based on queries.

Unfortunately, for many modern applications “the future is unknown”, e.g., in scientific databases, social networks, web logs, etc. In particular, the query processing patterns follow an exploratory behavior, which changes so arbitrarily that it cannot be predicted. Such environments cannot be handled by offline indexing. Online indexing [10, 47] and adaptive indexing [27] are two approaches to automatic physical design in such dynamic environments, but none of them in isolation handles the problem sufficiently. Online indexing periodically refines the physical design but it may negatively affect running queries every time it needs to use resources for reconsidering the physical design and it may not be quick to follow the workload changes as it reacts only periodically. Adaptive indexing does not have this problem as it introduces continuous, incremental and partial index refinement but it adjusts the physical design only during query processing based on queries.

Always Indexing. In this paper, we make the observation that in real systems there are plenty of resources that remain under-utilized and we propose to exploit those resources to be able to better address dynamic and ad-hoc environments. In particular, we focus on exploiting CPU cycles to the maximum by continuously detecting idle CPU cycles and using them to refine the physical design (in parallel with query processing). Such idle CPU cycles occur when the system does not exploit all available cores up to 100%. We distinguish between “idle time” as in “there is no user-driven work-

Indexing	Statistical analysis	Exploitation of idle resources <i>before</i> query processing	Exploitation of idle resources <i>during</i> query processing	Index materialization	Updates, projection cost	Workload
Offline	✓	✓	×	full	high	static
Online	✓	×	✓	full	high	dynamic
Adaptive	×	×	×	partial	low	dynamic
Holistic	✓	✓	✓	partial	low	dynamic

Table 1: Qualitative difference among offline, online, adaptive and holistic indexing.

load at all and the entire CPU (all its hardware contexts) is idle (except possible occasional operating system background activity)” and “idle CPU resources” as in “the active user-driven workload does / can not use all physically available CPU hardware contexts.” Intuitively, there are two options when resources are under-utilized but still there are active queries in the system. The first option is to introduce more parallel query processing algorithms to maximize utilization for the existing workload. The second one is to exploit the extra resources towards a different goal (extra indexing actions in our case). We investigate and compare both directions.

Holistic Indexing. In this paper, we introduce a new indexing approach, called *holistic indexing*. Holistic indexing addresses the automatic physical design problem in modern applications with dynamic and exploratory workloads. It continuously monitors the workload and the CPU utilization; when idle CPU cycles are detected, holistic indexing exploits them in order to partially and incrementally adjust the physical design based on the collected statistical information. Each index refinement step may take only a few microseconds to complete and the system will typically perform several such steps in one go depending on available system resources. Everything happens in parallel to query processing but without disturbing running queries. The net effect is that holistic indexing refines the physical design, improving performance and robustness by enabling better data access patterns for future queries. Table 1 summarizes the qualitative difference between holistic indexing and current state-of-the-art indexing approaches. Compared to past approaches, holistic indexing manages to minimize both initialization and maintenance costs, as it relies on partial indexing, and to exploit all possible CPU resources in order to provide a more complete physical design.

Contributions. Our contributions are as follows:

- We introduce the idea of exploiting idle CPU resources towards continuously adapting the physical design to ad-hoc and dynamic workloads.
- We discuss in detail the design of holistic indexing on top of modern column-store architectures, i.e., how to detect and exploit idle CPU resources during query processing.
- We implemented holistic indexing in an open-source column-store, MonetDB [26, 52]. Through a detailed experimental analysis both with microbenchmarks and with TPC-H we demonstrate that we can exploit idle CPU resources to prepare the physical design better, leading to significant improvements over past indexing approaches in dynamic environments.

Paper Structure. The rest of the paper is structured as follows: Section 2 provides an overview of related work. In Section 3, we shortly recap the basics of column-store architectures and the basics of adaptive indexing. Then, Section 4 introduces holistic indexing, while Section 5 presents a detailed experimental analysis. Finally, Section 6 discusses future work and concludes the paper.

2. RELATED WORK

In this section we briefly discuss previous approaches to automated physical design, i.e., offline, online and adaptive indexing.

Offline Indexing. Offline indexing is the earliest approach on self-tuning database systems. Nowadays, all major database products offer auto-tuning tools [5, 14, 50] to automate the database physical design. Auto-tuning tools mainly rely on a “what-if analysis” [12] and close interaction with the optimizer [11] to decide which indices are potentially more useful for a given workload.

Offline indexing requires heavy involvement of a database administrator (DBA). Specifically, a DBA invokes the tool and provides its input, i.e., a representative workload. The tool analyzes the given workload and recommends an appropriate physical design. However, the DBA is the one that decides which of the changes in the physical design should be applied. The main limitation of offline indexing appears when the workload cannot be predicted and/or there is not enough idle time to invest in the offline analysis and the physical design implementation.

Online Indexing. Online indexing addresses the limitation of offline indexing. Instead of making all decisions a priori, the system continuously monitors the workload and the physical design is periodically reevaluated. System COLT [47] was one of the first online indexing approaches. COLT continuously monitors the workload and periodically in specific epochs, i.e., every N queries, it reconsiders the physical design. The recommended physical design might demand creation of new indices or dropping of old ones. COLT requires many calls to the optimizer to obtain cost estimations. A “lighter” approach, i.e., requiring less calls to the optimizer, was proposed later [10]. Soft indices [43] extended the previous online approaches by building full indices on-the-fly concurrently with queries on the same data, sharing the scan operator.

The main limitation of online indexing is that reorganization of the physical design can be a costly action that a) requires a significant amount of time to complete and b) requires a lot of resources. This means that online indexing is appropriate mainly for moderately dynamic workloads where the query patterns do not change very frequently. Otherwise, it may be that by the time we finish adapting the physical design, the workload has changed again leading to a suboptimal performance.

Adaptive Indexing. Adaptive indexing is the latest and the most lightweight approach in self-tuning databases. Adaptive indexing addresses the limitations of offline and online indexing for dynamic workloads; it instantly adjusts to workload changes by building or refining indices partially and incrementally as part of query processing. By reacting to every single query with lightweight actions, adaptive indexing manages to instantly adapt to a changing workload. As more queries arrive, the more the indices are refined and the more performance improves. Adaptive indexing has been studied in the context of main-memory column-stores [27, 48], Hadoop [46] as well as for improving more traditional disk-based settings [20]. It has been shown to work for many core database architecture issues such as updates [28], multi-attribute queries [29], concurrency control [8, 16, 17], partition-merge-like logic [20, 31]. In addition, [18] shows how to benchmark adaptive indexing techniques, while stochastic database cracking [21] shows how to be robust on various workloads and [19] shows how adaptive indexing can apply to key columns. Finally, recent work on parallel adaptive indexing studies CPU-efficient implementations and proposes algorithms to

exploit multi-cores [8, 44]. Recently, adaptive indexing concepts have been extended to provide adaptive indexes for time series data [51] as well as using incoming queries for more broad storage layout decisions, i.e., reorganizing base data (columns/rows) according to incoming query requests [7], or even about which data should be loaded [25]. In addition, adaptive indexing ideas have been used to design new generation data exploration tools such as touch-based data systems [30, 42].

The main limitation of adaptive indexing is that it works only during query processing. In this way, the only opportunity to improve the physical design is only when queries arrive.

Database Systems for the Multi-core Era. Modern hardware offers opportunities for high parallelism; a single machine may be equipped with chip multiprocessors, which contain multiple cores with support for multiple context threads. Recent research focuses on exploiting parallelism opportunities by a) processing multiple queries concurrently, and b) by parallelizing tasks in the critical path during query processing [22, 23, 24, 33]. Sorting is one of the most important database tasks (and a core component of adaptive indexing in column-stores) that can be highly-parallelized using modern hardware advances [9, 38, 45, 49].

Holistic Indexing. Contrary to past indexing approaches, holistic indexing results in an always-on self-tuning database system. Holistic indexing is inspired by all past approaches and it maintains the design points which are useful for dynamic workloads. For example, it uses adaptive indices as in adaptive indexing and it monitors the workload as in online indexing. Contrary to other approaches though, holistic indexing is always active, always trying to improve the physical design with every opportunity that occurs when it detects under-utilized CPU cores. Compared to recent research that tries to adapt modern database systems to the multi-core era by parallelizing core database operators, holistic indexing provides an additional way to exploit those resources when either the workload is not enough to saturate the CPUs or when simply we cannot fully parallelize all database actions.

3. BACKGROUND

Holistic indexing is designed for column-oriented database architectures and it exploits the main design points of adaptive indexing. In this section, we provide the necessary technical background about column-store database architectures [4] and we discuss the basics of adaptive indexing.

3.1 Column-oriented DBMS

Column-oriented DBMS are inspired by the Decomposition Storage Model [13]. Data is stored one column at a time instead of one row at a time as in traditional row-oriented DBMS. Every relational table is vertically fragmented into a set of columns (one for each attribute). Each value of a single tuple is stored in the same position across all columns. Full vertical fragmentation significantly reduces the I/O and memory footprint of queries that require only part of a table's attributes; only the attributes that are relevant to a query are loaded from disk to memory. Moreover, the alignment across all base columns allows for efficient late tuple reconstruction with tuple order-preserving operators. For instance, assume the following query.

select B from R where A > 10 and A < 20

A column-store DBMS performs two main steps to answer this query. First, a select operator searches column A for attribute values between 10 and 20. The intermediate result is a new column which contains the positions of the qualifying attribute values in the relation. Second, a project operator fetches the values residing

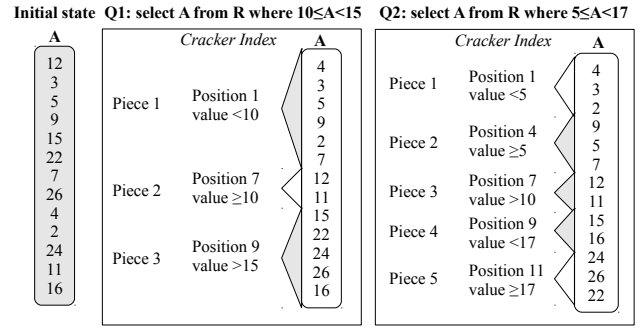


Figure 1: Adaptive indexing.

in attribute B at the positions specified by the intermediate result. Thus, column-oriented architectures allow for independent manipulation of the columns that are relevant to the query, which allows for operator implementations that exploit CPU and cache friendly patterns with tight for loops in an array-processing style.

With no indexing support, a column-store select operator has to completely scan a column in order to identify the qualifying tuples. The cost of scanning a column with N tuples is expressed in data accesses, i.e., how many of the column values are “touched”, and thus it is $O(N)$. The search for the qualifying tuples can be accelerated by orders of magnitude with a full indexing strategy. That is, we can first sort a column and then we may use binary search actions to find the qualifying tuples. The average cost of a binary search is $O(\log(N))$, while the average cost of sorting an entire column in memory, e.g., using quicksort, is $O(N * \log(N))$.

3.2 Adaptive Indexing

The cost of the sorting phase is a considerable overhead. Adaptive indexing is a significantly more lightweight approach that spreads the sorting cost across many queries in the workload. This helps when there is not enough idle time to create the full index upfront and the system has to answer queries. In addition, it is helpful when we do not know which columns are relevant for the workload a priori so we can sort them. This is especially important as today, as applications and schemas get more complex, we may have hundreds of columns in a database.

In adaptive indexing, the first time an attribute A is required by a query, a copy of the base column A is created, called the *cracker column* A_{CRK} of A. Each selection operator on A triggers the physical reorganization of A_{CRK} based on the requested range of the query. The query predicates are used as a hint on how the data should be stored, i.e., values that are less than the lower bound are moved before the lower bound while values that are greater than the upper bound are moved after the upper bound in the respective column. Thus, the column is partitioned on-the-fly, i.e., during query processing, based on the query predicates. The partitioning information for each cracker column is maintained in an AVL-tree, called *cracker index* of A. Index refinement is integrated with the query processing, since it is part of the select operator. Future queries on column A search the cracker index for the partition where the requested range falls. If the requested values already exist in the index, i.e., if past queries have cracked on exactly those ranges, then the select operator can return the result immediately. Otherwise, the select operator refines on-the-fly the column further, i.e., only the partitions/pieces of the column where the predicates fall will be reorganized.

For instance, Figure 1 shows an example of two queries cracking a single column, one after the other. The first query creates three partitions, while the second query refines only the first par-

tion, where its lower bound falls, and the third partition, where its upper bound falls. The more queries arrive, the more partitions/pieces are created. Thus, future queries have to refine smaller and smaller pieces, which results in a performance improvement as smaller pieces means that queries need to access less data.

The cost of the first cracking query in terms of data accesses is $O(N)$, since it has to analyze every tuple in the column. Assuming for simplicity that every query cracks a single piece in half, then the second query will touch $N/2$ tuples and so on. Thus, the cost of the i -th query becomes $N/2^{\text{floor}(\log_2(i))}$. The i -th query pays also the cost of searching the values in the cracker index. This cost is at most $\log(i)$, i.e., the depth of the AVL-tree.

The overall cost depends on the workload. However, the tendency is always the same; as more queries arrive, more partitions are created and, and each query has to touch less and less data during the select operator. In this way, adaptive indexing provides a lightweight alternative to full indexing; it can instantly adjust to a workload change and keep improving as the workload patterns persist, eventually reaching similar performance to full indexing but without requiring big initialization costs and workload knowledge.

Holistic Indexing vs. Adaptive Indexing. Holistic indexing maintains the design points of adaptive indexing, i.e., lightweight indexing with partial and incremental indices integrated with query processing. It extends adaptive indexing in that with holistic indexing the adaptive indices are continuously refined and improved; holistic indexing is always active; it continuously monitors the workload and the CPU utilization and reorganizes the physical design concurrently with query processing as soon as there are under-utilized CPU resources.

4. HOLISTIC INDEXING

In this section we discuss the fundamentals of holistic indexing. We designed holistic indexing on top of column-store architectures inspired by their flexibility on manipulating some attributes without affecting the rest. During query processing indices are built and optimized incrementally by adapting to query predicates, as in adaptive indexing. However, in contrast to adaptive indexing, index refinement actions are not triggered only as a side-effect of query processing; in holistic indexing incremental index optimization actions take place continuously in order to exploit under-utilized CPU cores. Thus, concurrently with user queries, system queries also refine the index space. Holistic indexing monitors the workload and CPU resources utilization and every time it detects that the system is under-utilized it exploits statistical information to decide which indices to refine and by how much.

Thus, with holistic indexing we achieve an always active self-organizing DBMS by continuously adjusting the physical design to workload demands.

Problem Definition: *Given a set of adaptive indices, statistical information about the past workload, storage constraints and the CPU utilization, continuously select indices from the index space and incrementally refine them, while the materialized index space size does not exceed the storage budget.*

In the rest of this section, we discuss in detail how we fit holistic indexing in a modern DBMS architecture.

4.1 Preliminary Definitions

First, we give a series of definitions.

Workload. A workload W consists of a sequence of user queries, inserts and deletes. Updates are translated into a deletion that is followed by an insertion.

CPU Utilization. CPU utilization in a time interval dt describes how much of the available CPU power is used in dt . Specifically, it expresses the percentage of total CPU time, i.e., the amount of time for which the CPU is used for processing user or kernel processes instead of being idle. CPU utilization is calculated using (operating system) kernel statistics.

Configuration. A configuration is defined as a set of adaptive indices that can be used in the physical design. There are three kinds of configurations. The *actual configuration*, C_{actual} , contains indices on attributes that have already been accessed by user queries in the workload. Indices are inserted in C_{actual} when they are created during query processing. For instance, assume a query Q enters the system and contains a selection on an attribute A . If the adaptive index on A does not exist, it is created on-the-fly and it is inserted in C_{actual} .

Besides C_{actual} , holistic indexing also maintains the *potential configuration*, $C_{potential}$, which contains indices on attributes that have not been queried yet. Indices are inserted in $C_{potential}$ either automatically by the system or manually by the user. Finally, the *optimal configuration*, $C_{optimal}$, contains indices that have reached the optimal status (the next paragraph describes when an index is considered optimal). The union of C_{actual} and $C_{potential}$ constitutes the index space IS , i.e., the indices which are candidates for incremental optimization when the system is under-utilized. Later, in Section 4.2, we describe how the system is educated to pick an index from IS . Indices from $C_{optimal}$ are not considered for further refinement during the physical design reorganization.

Optimal Index. Holistic indexing exploits adaptive indices. As seen in Section 3.2, an adaptive index is refined during query processing by physically reorganizing pieces of the cracker column based on query predicates. As more queries arrive, more pieces are created, and thus, the pieces become smaller. We have found that when the size of the pieces becomes equal to L_1 cache size ($|L_1|$), further refinements are not necessary; a smaller size increases administration costs to maintain the extra pieces and it does not bring any significant extra benefit as scanning inside L_1 is fast anyway (no cache misses). Pieces of size smaller than L_1 cache can either be sorted or queries simply need to scan them (a range select operator has to scan at most two L_1 pieces). An index I on an attribute A is considered optimal (I_{opt}), when the average size of pieces ($|p|$) in $ACRK$ is equal to the size of L_1 cache. Equation (1) describes the distance between I and I_{opt} .

$$d(I, I_{opt}) = |p| - |L_1| = \frac{N_A}{p_A} - |L_1| \quad (1)$$

N_A is the total number of tuples in $ACRK$ while p_A is the total number of pieces in $ACRK$. This information is readily available and thus we can easily calculate the average piece size in a cracker column and in turn we can calculate the distance of the respective cracker index from its optimal status.

Statistical Information. During query processing holistic indexing continuously monitors the workload and the CPU utilization. For each column in the schema it collects information regarding how many times it has been accessed by user queries, how many pieces the relevant cracker column contains, how many queries did not need to further refine the index because there was an exact hit. Besides the statistical information about the workload, kernel statistics are used in order to monitor the CPU utilization.

4.2 System Design

Holistic indexing is always active. It continuously monitors the workload and the CPU utilization. When under-utilized CPU cores are detected, holistic indexing exploits them in order to adjust the physical design based on the collected statistical information. The

system performs several index refinement steps simultaneously depending on available CPU resources. Everything happens in parallel to query processing, but without disturbing running queries.

We discuss in detail the continuous tuning process and how to exploit under-utilized CPU cycles. We also discuss how existing adaptive indexing solutions on core database architectures issues—such as updates and concurrency control can be directly adapted to work with holistic indexing.

Statistics per Column/Index. Statistics per column are collected during query processing. This is the job of the select operator as it is within the select operator that all (user query) adaptive indexing actions take place. Every time an attribute is accessed for a selection of a user query, the select operator updates a data structure, which contains all statistics for the respective index. Given that the select operator performs adaptive indexing actions anyway, it already has access to critical information such as how many new pieces were created during new cracking actions for this query, whether the select was an exact match, etc. All information is stored in a heap structure (one node per index) which allows us to easily put new indices in the configuration or drop old ones. The structure is protected with read/write latches as multiple queries or holistic workers (discussed later on) may be cracking in parallel.

Tuning Cycle. At all times there is an active *holistic indexing thread* which runs in parallel to user queries. The responsibility of the holistic indexing thread is to monitor the CPU utilization and to activate *holistic worker threads* to perform auxiliary index refinement actions whenever idle CPU cycles are detected. The tuning process is shown in Figure 2. The holistic indexing thread continuously monitors the CPU load at intervals of 1 second at a time. In case holistic worker threads are activated, the holistic indexing thread waits for all worker threads to finish and measures the CPU utilization within the next 1 second. In our analysis, we found that 1 second is the time limit that gives proper kernel statistics. When n idle CPU cores are detected, n holistic worker threads are activated. Each worker thread executes an instance of the *IdleFunction*, which picks an index from the Index Space IS and performs x partial index refinement actions on it. Every time an index is refined, the respective statistics, e.g., distance from the optimal index, are updated. When an index reaches the optimal status, it is moved into the optimal configuration $C_{optimal}$.

A side-effect of the tuning process is that some of the holistic worker threads might remain idle while the holistic indexing thread waits for all workers to finish. However, as we show later in Section 5.1 (Figure 6(d)), this happens only for very short periods of time and as the system adapts to the workload this phenomenon disappears (as the pieces queried in the adaptive indices become smaller and smaller the holistic indexing workers end up doing tasks of similar weight as none is going to touch a very big piece).

Index Refinement. Every time a worker thread wakes up, it performs x index refinements on a single column. x is a tuning parameter. In our analysis in Section 5.5 (Figure 15) we found that a good value for our hardware set-up is $x = 16$. The index refinements are performed by picking x random values in the domain of the respective attribute and cracking the column based on these values. In this way, each time a worker thread cracks a single piece of a column it splits this piece into two new pieces based on the pivot.

There are numerous choices on how to choose a pivot. We found that picking a random pivot is the most cost efficient choice. Other options include to crack the biggest piece of the column, i.e., with the rationale that this takes more work out of future queries. Another option is to crack the smallest piece, i.e., with the rationale that this piece is small because it is hot (because many queries access it for cracking). However, such options are hard to achieve

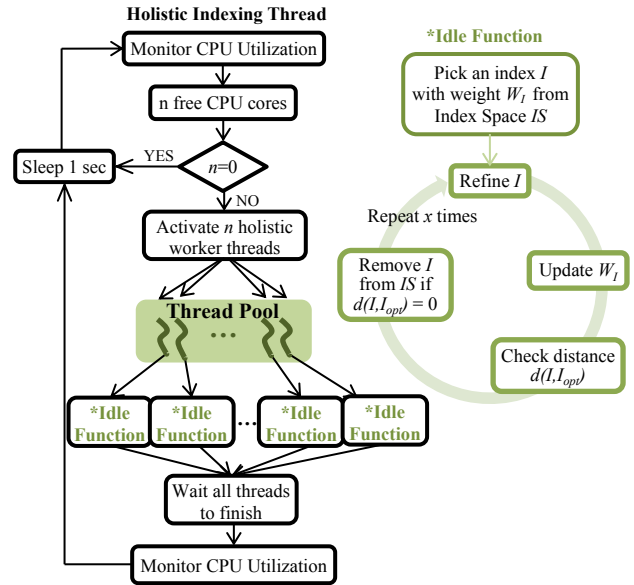


Figure 2: Tuning actions.

in a lightweight way as we need to maintain a structure such as a priority queue to know which piece is the biggest or smallest every time. Since every cracking action costs a few microseconds or milliseconds it is not worth the extra storage and CPU cost to maintain auxiliary structures. Random pivots converge quickly to cracking the whole domain, providing a column which is balanced in terms of which pieces are cracked and requiring no extra costs in deciding which pivot to choose.

Index Decision Strategies. Another decision we have to make is which index to refine out of the pool of candidate indices. Here, we describe four different strategies we can follow in order to pick an index from the index space. The notion behind the first three strategies is that, since the only information we have is the past workload, we can exploit this information in order to prepare the physical design for a similar future workload. On the contrary, the fourth strategy makes random choices.

For all strategies, a weight W_I is assigned to each index I in the index space. When an index I is added in the candidate indices, its weight is initialized to the distance between I and I_{opt} , which is given by Equation (1). For each index I , initially, there is only one partition ($p_I = 1$) in I , i.e., the entire column. Thus, the initial weight $W_{I_{init}}$ is equal to $N_I - L_{1_s}$, where N_I is the cardinality of the respective attribute (with type T) and L_{1_s} is the number of elements of type T that can fit into L_1 cache. The weight is used as a priority number in the first three strategies. The index with the highest priority, i.e., the maximum weight, in C_{actual} is refined first. When W_I becomes equal to zero, I is transferred from C_{actual} to $C_{optimal}$ and it is not considered for further refinement in the future. If C_{actual} is empty, an index is randomly picked from $C_{potential}$. The weight of each index is constantly updated after every index refinement regardless of whether it is caused by a user query or by holistic indexing. Below we describe the four strategies.

- **W1:** $W_I = d_I = d(I, I_{opt})$. Using this strategy, we give a priority to indices with large partitions.
- **W2:** $W_I = f_I * d_I$. Priority is given to indices that have large partitions and at the same time are accessed frequently in the workload. f_I is the number of user queries that access I .
- **W3:** $W_I = (f_I - f_{I_h}) * d_I$. In this strategy we try to identify indices that are accessed frequently in the workload and at the same time have large partitions, because they have a high

Example of latching actions in holistic indexing. It depicts latching actions during concurrent processing of two user queries and a holistic worker thread (HW). HW cracks a random piece. Red box indicates a write latch. Green box indicates a read latch. Grey boxes indicate pieces that are free to be refined by holistic worker threads.

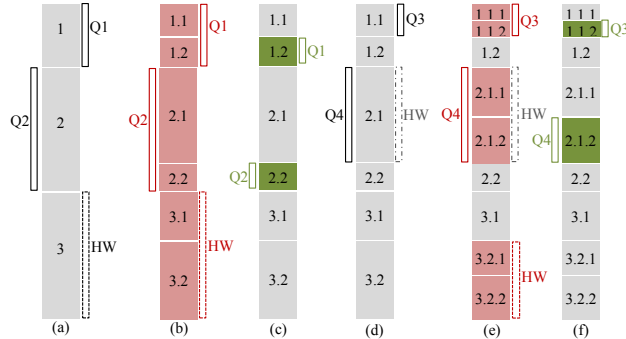


Figure 3: Concurrency control in holistic indexing.

hit rate. These indices have a smaller priority compared to indices with large partitions that are accessed less frequently. f_l is the number of user queries that access l , while f_{l_h} is the number of user queries that do not trigger a refinement of l because the requested value bound already exists in l .

- **W4:** Make a random choice.

Overall, our analysis, which is described later in Section 5.4 (Figure 13), with numerous workloads showed that even though small improvements can be achieved when picking the perfect strategy for each workload, the random strategy gives a good and robust overall solution that is always close to the best for all workloads.

Concurrency Control. An index refinement due to holistic indexing happens in parallel with user queries. Since user queries may also cause refinement of adaptive indices, we need to properly control these changes. In addition, as more than one holistic thread may be active at any time, they may be trying to refine the same index. The study of concurrency control for adaptive indexing [16, 17] showed that it is possible to allow multiple concurrent index refinements in adaptive indices via lightweight concurrency control, i.e., relying only on latches of individual pieces in an adaptive index. The point is that an index refinement only changes the structure of the index and not its contents (contrary to an actual update). In this case, an index refinement only rearranges values in a single piece of a column at a time. Thus it is sufficient to allow other queries to work on different pieces in parallel by taking read/write latches on individual pieces, called piece latches in [16, 17]. We exploit these techniques here in order to allow user queries and holistic workers to work concurrently over a single column, but we also identify extra opportunities to increase parallelism for holistic workers.

Figure 3 shows an example of an adaptive index where two queries are actively cracking it. Each query is interested in its own value range and needs to crack one piece, i.e., at the value of its selection. The idea is that *all* other pieces of the column are available for index refinement by holistic worker threads. One direction would be that each holistic worker decides which piece of an index to refine by picking from a list of pieces that currently have no locks. However, such information is expensive to maintain similarly to our discussion in the “Index Refinement” paragraph. Thus, holistic workers make random choices regarding which value to use as pivot and thus which piece to crack. However, when a holistic worker requests a write latch to crack a piece and it happens that the piece is locked at the moment, then if the latch is not given immediately, the worker picks another random pivot and repeats the procedure

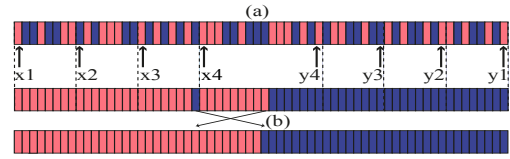


Figure 4: Refined Partition & Merge (multi-threaded) [44].

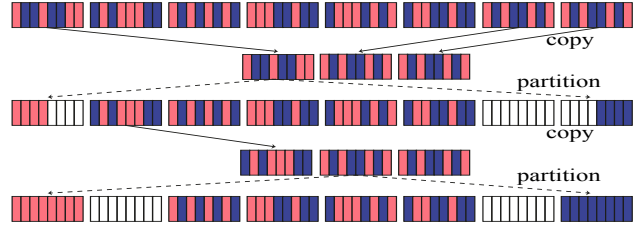


Figure 5: Vectorized Cracking [44].

until it finds a free piece to crack. In contrast, user queries need to always block in such cases and wait for the piece to be unlocked. For instance, in Figure 3(d) the holistic worker thread tries to lock piece 2.1, which is already locked by Q4. Instead of waiting for the lock to be released, the worker chooses another pivot. The new pivot falls in piece 3.2, which is not locked and it is reorganized finally by the worker (Figure 3(e)). As we process more queries and as we perform more holistic indexing, the number of pieces in an index grows; as a side-effect the waiting time for taking a latch decreases as there are more candidate pieces to pick from.

Updates. Updates for adaptive indexing have been studied in [29]. The design in [29] is that updates remain as pending updates and are merged during query processing, i.e., if a query requests a value range that contains one or more pending updates, then only those updates are merged on-the-fly and without destroying any of the information on the adaptive index. Each query needs to lock at most one column piece at a time for cracking and can update this piece at the same time if pending updates for this piece exist [16, 17]. Multiple queries may work in parallel updating and cracking separate pieces (value ranges) of the same column.

The difference here is that with holistic indexing, holistic workers not only perform auxiliary index refinement actions but also merge pending updates. That is, if a holistic worker picks a pivot which falls within a piece of the respective column and the value range, for which this piece holds values, has pending updates, then all those pending updates are merged by the holistic worker. In this way, holistic worker threads not only refine the adaptive indices in the background but also bring them more up to date which removes further load from future queries.

Storage Constraints. Holistic indexing works within a limited storage budget. Adaptive indices may be dropped or recreated at any time. They consist auxiliary information and thus dropping an index does not lead to any loss of data. In case the storage budget does not allow adding a new index triggered by a user query, then indices are removed with a least frequently used (LFU) policy from the index space at an index-level granularity or at a fine-grained granularity that allows for creating and dropping individual ranges dynamically, as partial cracking suggested in [29].

Multi-core Adaptive Indexing. The goal of holistic indexing is to improve the physical design by fully utilizing the available CPU resources. An alternative approach to achieve maximum CPU utilization is to parallelize the index refinement actions triggered by user queries. This problem was studied in [44], which introduced a multi-core, CPU efficient cracking algorithm shown in Figure 4. In this algorithm, the to-be-cracked piece is partitioned initially into as many slices as the number of threads, e.g., n (Figure 4(a)). The

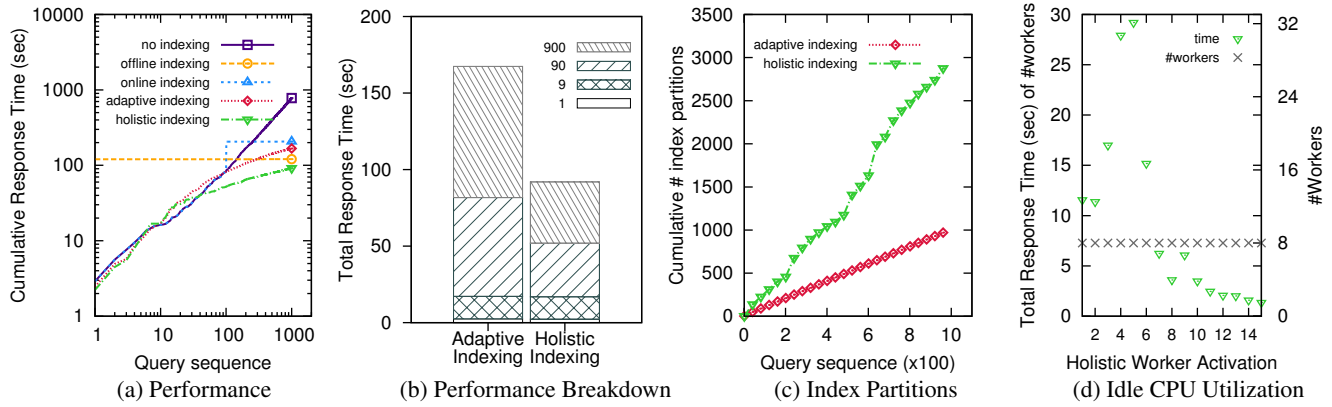


Figure 6: Improving performance with holistic indexing.

center slice is contiguous, while the remaining $n - 1$ slices consist of two disjoint halves, each, that are arranged concentrically around the center slice (x_i and y_i indicate the first and the last element of piece i respectively). n threads crack the n slices independently applying a vectorized, out-of-place cracking algorithm (Figure 5), which was proven in [44] to be the most CPU efficient single-threaded cracking implementation reported so far. Finally, the local data are merged into a big cracked piece (Figure 4(b)). We found that devoting all resources to perform adaptive indexing for user queries in parallel does not lead to the absolute best performance. Specifically, we found that we can improve performance even more by assigning part of the resources to holistic indexing. In this way, some of the CPU resources are assigned to parallel cracking for user queries but the rest of the CPU resources are distributed across several holistic workers for additional index refinements. In the experimental section we show why this approach is better than assigning all available resources to user queries.

5. EXPERIMENTAL ANALYSIS

In this section, we demonstrate that holistic indexing leads to a self-organizing always-on DBMS with substantial benefits in terms of response time; with zero administration or set-up effort holistic indexing improves performance adaptively by exploiting all available CPU resources to the maximum. We present a detailed experimental analysis using both standard benchmarks such as TPC-H and real-life workloads such as SkyServer as well as synthetic microbenchmarks for a fine-grained analysis.

We use a dual-socket machine equipped with two 2.00 GHz Intel(R) Xeon(R) CPU E5-2650 processors and with 256 GB RAM. Each processor has 8 hyper-threading cores resulting in 32 hardware threads in total. The operating system is Fedora 20 (kernel version 3.12.10). All experiments we report are based on an implementation of holistic indexing in MonetDB and assume a main-memory environment.

5.1 Improving over State-of-the-Art Indexing

In our first experiment we demonstrate that holistic indexing has the potential to bring substantial performance improvements over existing state-of-the-art indexing approaches. We test holistic indexing against parallel versions of adaptive indexing (database cracking), offline indexing, online indexing and plain scans.

For plain scans (no indexing), we use a parallel select operator implemented in MonetDB. For offline and online indexing we sort the columns using a highly parallel NUMA-aware sorting algorithm that was introduced in [9] (m-way, 16-byte keys) and is publicly available in [1]. Specifically, in offline indexing we pre-sort all the columns before query processing, while in online indexing

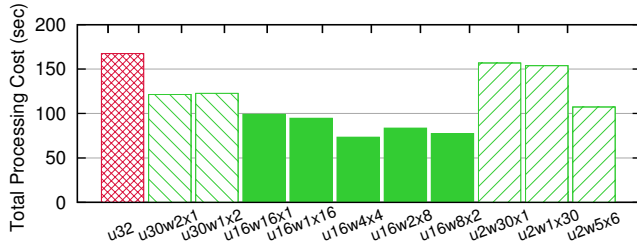
we assume that after processing a few queries we understand the workload patterns and then we sort the relevant columns. MonetDB automatically detects that a column is sorted and can use efficient binary search actions during select operations. For adaptive indexing we use the parallel vectorized database cracking algorithm that was introduced in [44] (see Section 4.2).

Here we use a synthetic benchmark. The query workload consists of 10^3 range select queries over a table of 10 attributes; each query touches a single attribute (we will see more complex queries later on). Each attribute consists of 2^{30} uniformly distributed integers, while the value range requested by each query (and thus the selectivity) is random. All queries are of the following form.

select A from R where A < v

The tested scenario assumes a dynamic and ad-hoc environment with zero workload knowledge and zero idle time to pre-index the data. Figure 6(a) shows the results. On the x-axis queries are ranked in execution order. The y-axis represents the cumulative response time as the query sequence evolves, i.e., each point (x, y) represents the sum of the execution time y for the first x queries. In this way, the graph shows how the response times evolve as we process more and more queries.

If there is no indexing support (plain scans), the entire column is scanned in parallel by 32 threads for every query. Because of this stable access pattern, the cumulative response time of the query sequence grows linearly as every query has similar cost. With offline indexing, on the other hand, it takes 12 seconds to completely sort each column, assuming a priori workload knowledge. This leads to a 120 seconds initialization overhead to sort all 10 columns. Since there is no idle time before the first query, the sorting cost is added to the execution time of the first query in Figure 6(a). After the first query, all queries are answered with fast binary search actions which results in a rather flat cumulative curve. With online indexing, the first 100 queries are answered without any index support and thus the cost grows linearly. After the 100th query, assuming enough workload knowledge has been obtained via monitoring, we proceed to sort all 10 columns. The sorting cost is added to the execution cost of the 101st query, since there is no idle time between the 100th and the 101st query. As with offline indexing, all subsequent queries are answered very fast with binary search over the sorted column. Thus, both in offline and in online indexing, the whole query sequence is affected by the sorting costs. On the other hand, adaptive indexing continuously improves performance requiring no workload knowledge and without penalizing individual queries. This improvement is due to the fact that adaptive indexing builds only partial indices which it incrementally refines as queries arrive. However, there is still room for big improvements.



Distribution of 32 Threads between Users and Workers

Figure 7: Performance improves if we distribute the threads equally between user queries and holistic workers.

Holistic indexing manages to further improve the performance of the workload by about 50%. Contrary to the other indexing approaches, MonetDB with holistic indexing enabled monitors the CPU utilization and constantly tries to maximize it. When holistic indexing detects idle CPU resources, it triggers index refinement actions on existing adaptive indices. In this experiment, an index is inserted in the index set, and specifically in C_{actual} , when a user query creates it. For holistic indexing the actual user queries behave in exactly the same way as in adaptive indexing, i.e., the first query will create an adaptive index and subsequent queries refine it using the very efficient and almost linearly scaling parallel vectorized cracking implementation from [44]. The difference is that with holistic indexing enabled idle CPU resources are exploited towards further refining the adaptive indices in a way which does not hurt running queries. Since parallel vectorized cracking [44] is designed to be CPU efficient, encountering only very little resource stalls, we generated kind of a worst-case scenario for holistic indexing: we limited the maximal number of hardware context assigned to user queries to 16 (equal to number of physical cores), leaving at least 16 (otherwise not effectively usable by the prime user query workload) hardware contexts (“hyper-threads”) available for holistic indexing. Constantly, our load-checker usually detects 16 idle hardware contexts, and consequently starts 8 holistic indexing workers (each using two threads) as shown in Figure 6(d). Figure 7 shows that the combination of using maximal 16 (out of 32) hardware contexts for user queries (performing parallel vectorized adaptive indexing [44]), while devoting any remaining idle hardware context to holistic indexing, improved the overall performance by a factor 2 over using all 32 hardware contexts for user queries (and thus none for holistic indexing).

Figure 6(b) is a breakdown of the performance of holistic indexing and adaptive indexing. The y-axis represents the total response time of the first query, the next 9 queries, etc. The total height of each bar represents the total response time to run the entire workload of 10^3 queries. The first few queries do not see any improvement because holistic indexing cannot concurrently refine a column if there are user queries cracking it. This is because initially columns have not been cracked at all and thus the first few user queries will lock big pieces for cracking. However, as each column is cracked into smaller pieces, holistic indexing may invoke actions to refine a column even if concurrent queries are cracking it. Essentially, each user query needs to lock at most one piece of an adaptive index at a time, i.e., the piece it is about to crack, and thus holistic indexing may choose any of the remaining pieces to perform further index refinements.

Holistic indexing outperforms adaptive indexing by a factor 2 by injecting additional index refinements on top of those that adaptive indexing does anyway. Figure 6(c) shows the amount of pieces which have been created in all 10 adaptive indices; holistic indexing creates more pieces than adaptive indexing. As a result, future

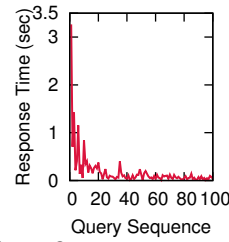


Figure 8: Adaptive Ind.

As discussed in Section 4.2, on some occasions the main holistic indexing thread waits for all workers to finish before assigning new tasks, leaving under-utilized CPU resources for some brief periods. Figure 6(d) shows how the total response time of all workers in every tuning cycle changes over time and as the query sequence evolves. The right y-axis shows the number of holistic worker threads the holistic indexing thread activates whenever it detects idle CPU resources. The maximum number of workers that holistic indexing can activate is 8. The left y-axis represents the total response time of all workers during a single tuning cycle. The x-axis represents the activations of holistic indexing. A single activation of holistic indexing triggers the activation of multiple holistic workers. The total response time of the workload is 90 seconds. Within this amount of time, holistic indexing is activated only 15 times, because of the waiting time (1 second) between two CPU load measurements and because of the waiting time until all workers finish in every tuning cycle. We observe that the response time of the workers is high only for the first few activations and reduces very fast as the index becomes fine-grained. In this way, the system adapts on its own and eventually no worker is a bottleneck.

Holistic indexing sees even bigger performance benefits when there is idle time before query processing. When there is idle time and no workload knowledge, holistic indexing chooses random indexes to insert in $C_{potential}$ and refines them until the first query arrives. Here using the same set-up as in the previous experiment, we manually induce idle time and holistic indexing adds 10 random indices in $C_{potential}$. Figure 9 shows the results. Compared to adaptive indexing, which does not exploit the a priori idle time

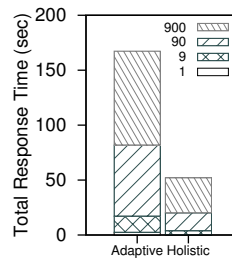


Figure 9: Indexing.

By being able to completely automatically utilize available CPU resources and direct them towards lightweight actions that may improve future requests, holistic indexing can bring significant performance gains on top of existing indexing approaches. It outperforms adaptive indexing by a factor 2 in terms of individual query performance. At the same time it outperforms offline and online indexing, especially in the beginning of the workload, when offline indexing penalizes the first queries with the index building cost and online indexing does not provide any indexing support until the 100th query. Besides the difference in performance, the qualitative difference described in Table 1, makes holistic indexing a very appealing indexing approach in modern dynamic environments.

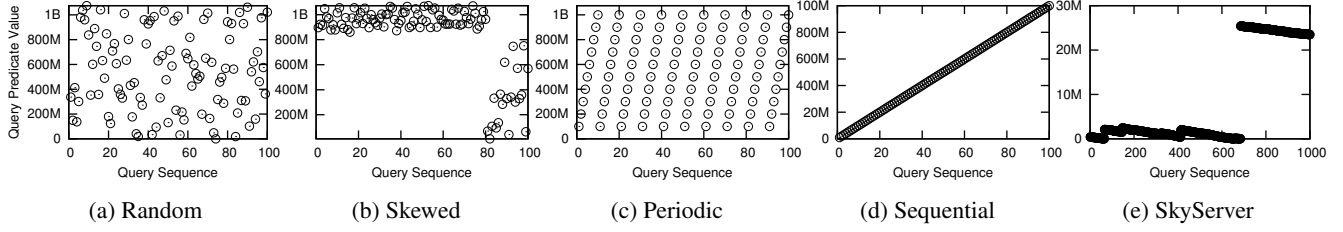


Figure 10: Various workload patterns.

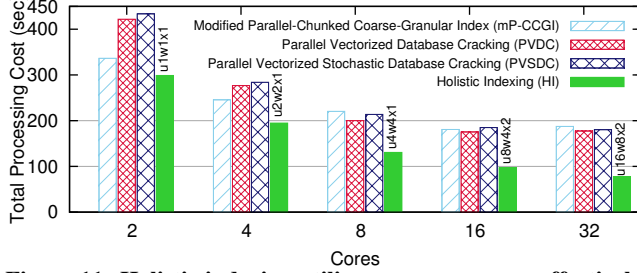


Figure 11: Holistic indexing utilizes resources more effectively than multi-core state-of-the-art adaptive indexing baselines.

5.2 Holistic Vs. Multi-core Adaptive Indexing

Assuming there are several CPU cores available in a modern system, in holistic indexing we utilize them by spreading auxiliary tuning actions across multiple indices. An alternative way to “keep the CPU busy”, is to parallelize the existing adaptive indexing algorithms. In this experiment we study how state-of-the-art adaptive indexing baselines compare to holistic indexing. In particular, we study parallel vectorized database cracking (PVDC), parallel vectorized stochastic database cracking (PVSDC) [44] and a modified version of parallel chunked coarse granular index (P-CCGI) [8] that we name mP-CCGI.

Stochastic cracking aims to provide robustness by performing auxiliary stochastic indexing actions. Although stochastic cracking studied the option of collecting statistics to target the proper value ranges to index, it was shown in [21] that reacting immediately to workload changes by auxiliary stochastic cracking actions has a better effect (i.e., more robust). This is because in stochastic cracking a running query performs auxiliary random cracking only within the piece that is already about to be cracked within a given column and as a result any action brings a benefit as it imposes more order. Holistic indexing considers a much more broad space of statistics keeping track of column-statistics to decide which columns to fine tune.

Both stochastic cracking and plain database cracking in this experiment utilize multi-cores as described in the last paragraph of Section 4. The original P-CCGI algorithm partitions the data into as many chunks as the available number of threads and the first query cracks each chunk in parallel having a separate cracking index for each chunk. Subsequent queries crack the chunks in parallel, while they benefit from the initial range partitioning. However, this way, data that belongs in a single value range is physically stored in separate chunks/arrays and feeding from there other relational operators is not compatible with a column-store such as MonetDB that relies on bulk processing; it does not allow to exploit tight for loops without intermediate if statements to detect when we should skip from one chunk to the next during an operator. To address this we extended the original P-CCGI algorithm [8] with the ability to consolidate selection results in a single array using the same techniques that were used for hybrid adaptive indexing which also operates on multiple chunks (but not in parallel) [31];

each query consolidates only the qualifying value ranges and each value range needs to be written to the contiguous array by a single query only, i.e., subsequent queries will only have to do consolidation if they need a new value range never consolidated before. In a vectorized column-store this could be done without consolidation, potentially improving performance further as has been indicated by partial sideways cracking [29]; vectorized processing for adaptive indexing is an open topic, though, orthogonal to this work.

The workload in this experiment consists of 10^3 select-project queries (as in the previous experiment) on 10 integer attributes. Each attribute consists of 2^{30} uniformly distributed integers. The value range requested by each query is random while we vary the number of available CPU cores from 2 to 32, i.e., the maximum number of cores in our system. For holistic indexing we give half of the cores to user queries and the rest of the cores are used by the workers (after testing all possible configurations, we found that this is the one that performs best in all cases). The labels on top of the bars that represent the performance of holistic indexing indicate the distribution of the threads between user queries and workers in every case (similar to Figure 7).

Figure 11 shows the results. In all cases, the performance improves as we invoke more cores into query processing. For multi-core database cracking and stochastic cracking the performance improves because many threads crack in parallel for one query at a time while for holistic indexing performance improves because many threads work in the background in parallel with query processing to further refine the various indices with auxiliary indexing actions. Holistic indexing sees a bigger improvement, because it is active all the time, i.e., maximizing CPU usage. On the contrary, multi-core vectorized database cracking and multi-core stochastic cracking improve the performance only during user queries and only during the cracking actions. Another subtle difference but one with a major performance impact is that while stochastic cracking and database cracking target all their adaptive indexing actions on specific pieces as they are driven by individual queries, holistic indexing spreads its actions across the whole range of the domain and thus across the whole range of an adaptive index (stochastic cracking does random actions but only within a single piece). This creates a nicely balanced index which has more potential to benefit future queries. The modified version of P-CCGI initially range partitions the data, which can be seen as a pre-index step. However, this is a cost that penalizes the first set of queries.

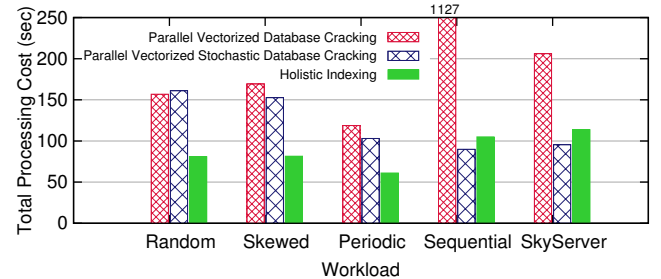


Figure 12: Holistic indexing is robust.

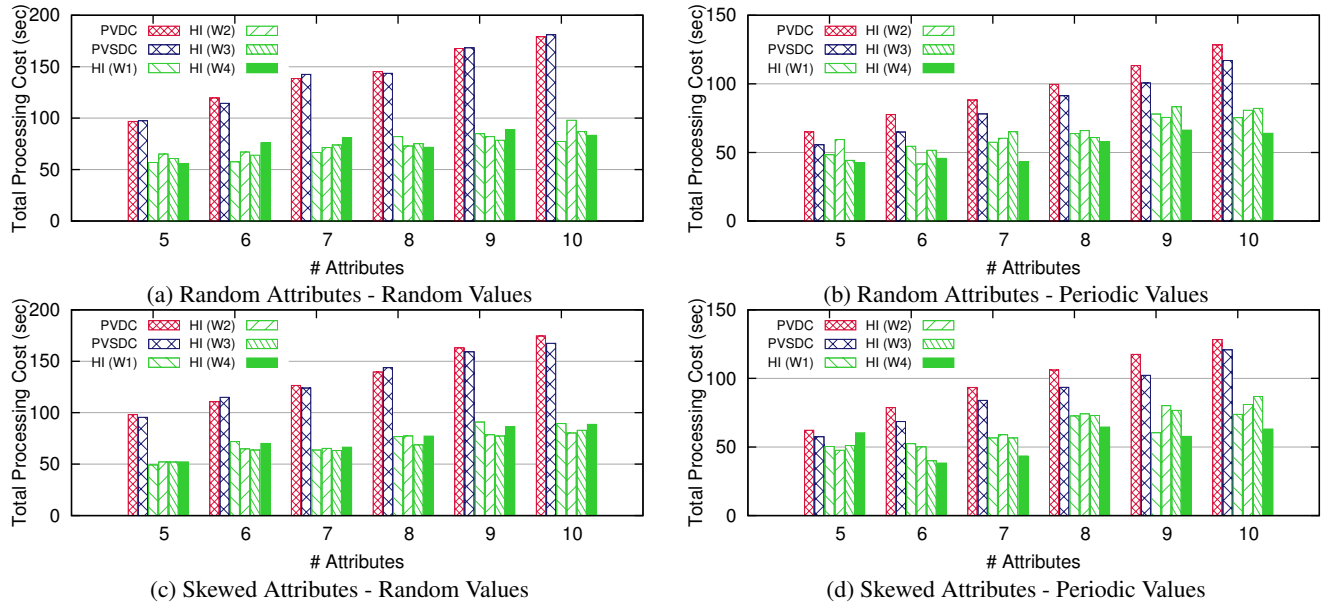


Figure 13: More performance gains for holistic indexing as more attributes exist in a schema. All strategies have similar performance.

5.3 Robustness

In our next experiment we study how holistic indexing compares to parallel database cracking and parallel stochastic cracking in terms of robustness. We show that holistic indexing maintains the good properties of adaptive indexing by utilizing the available CPU resources more effectively. Both holistic indexing and the parallel variants of adaptive indexing utilize all the available CPU cores.

We test four synthetic workloads. Each workload consists of 10^3 queries on 10 attributes (~ 100 queries/attribute). Each attribute consists of 2^{30} uniformly distributed integers. The queries follow a different pattern in each workload. The first four subfigures in Figure 10 depict those workload patterns. For each workload, the respective figure illustrates graphically how a sequence of queries touches the value domain of a single attribute.

Furthermore, we test holistic indexing in a real-life workload using data and queries from SkyServer [2]. SkyServer collects astronomical data and the database can be accessed publicly by individual users and institutions. We pose 10^4 real user queries that have been logged by the project servers on the “Photoobjall” table. The “Photoobjall” table consists of 1.2 Billion tuples. All queries access the “Ascension” attribute and are posed in exactly the same chronological order they were logged. The pattern the SkyServer queries follow is shown in Figure 10(e).

Figure 12 shows the results. For each indexing method we report the total time needed to process all queries for each workload.

Synthetic Workloads. In all synthetic workloads holistic indexing outperforms multi-core database cracking by a factor 2-10 depending on the workload. Multi-core database cracking is strictly driven by query predicates, and thus, can leave large unindexed pieces to be reorganized by future queries. For instance, in the sequential workload in Figure 10(d), each query cracks a column in a small piece and in a big piece, and then, a future query needs to crack the big piece again, resulting in a high cost. Multi-core stochastic cracking solves these robustness issues by injecting one extra random cracking action for each user query in order to distribute cracking more evenly. However, holistic indexing can materialize an even bigger advantage. This is because it is not restricted to perform auxiliary index refinement actions only during user queries but it can exploit all possible CPU cycles to refine the indices, resulting in many more actions taking place in parallel

with user queries. Moreover, holistic indexing spreads the auxiliary index refinements across the entire value domain (by choosing random pivots) without leaving big unindexed pieces. For example, in the skewed workload in Figure 10(b), both multi-core database cracking and multi-core stochastic cracking show a similar performance, because they restrict the index refinements to a small region of the domain according to user query predicates, i.e., from 800 million to 2^{30} . Future queries have to reorganize a big unindexed area, i.e., from 0 to 800 million; this area is already indexed in holistic indexing before the 800th query arrives. Thus, holistic indexing prepares the physical design better for (ad-hoc) future queries.

SkyServer. The SkyServer workload in Figure 10(e) shows the pattern logged in SkyServer for queries using the “right ascension” attribute of the “Photoobjall” table. We observe that the queries follow non-random patterns, i.e., they focus on a specific part of the sky before moving to a different part. Figure 12 shows that holistic indexing manages to significantly outperform multi-core database cracking by inducing auxiliary index refinement actions in parallel with query processing without penalizing individual user queries.

Overall, in all workloads tested, holistic indexing not only maintains the nice properties of the parallel variants of database cracking and stochastic cracking, but it also enhances the behavior further by being able to exploit all available CPU resources effectively for a better prepared physical design.

5.4 More Benefits with Complex Schemas

In this experiment, we show that as the database schema becomes more complex by containing more attributes, this brings more opportunities for holistic indexing to gain in performance; more attributes make the indexing space bigger and thus any a priori decisions are even more prone to be wrong. In addition, we test the various strategies for choosing among the candidate indices and we show that indeed making random decisions is a robust approach.

Here, we assume a gradually bigger database table which consists of 5-10 attributes. Each attribute consists of 2^{30} uniformly distributed integers. We fire select-project queries as in the previous experiments but this time we may query up to X attributes in every run. Each query touches a single attribute and we vary the frequency with which each attribute is accessed, i.e., we run both a random workload where every attribute is evenly queried as well

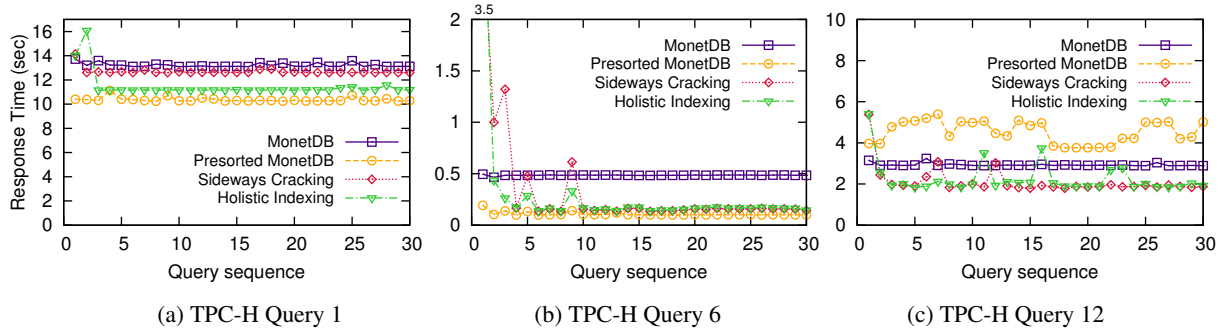


Figure 14: TPC-H results (Scale Factor 10, “pre-sorted” times exclude pre-sorting costs; Q1,6,12: 8 sec).

as a skewed workload where some attributes are queried more than others. For each workload we vary also the workload pattern followed by the queries. Here, we present the results for random and periodic workload patterns. For each case, we perform 10^3 queries.

We compare holistic indexing using one of the four strategies described in Section 4.2 against multi-core variants of database cracking and stochastic cracking. Figure 13 shows the results; for all cases, holistic indexing materializes a big benefit. As the number of attributes in the database table grows, the performance benefit for holistic indexing increases. What happens is that holistic indexing makes sure to evenly spread auxiliary index refinement actions across all attributes in parallel with user queries whenever free CPU cycles are available. Then, future queries on those attributes can exploit this refined indexing. Compared to the case where we have a few attributes, having more attributes means that more heavy indexing actions have to be performed overall in order to crack the columns into small pieces. This allows holistic indexing to materialize a bigger benefit as it performs those actions in the background as opposed to only during user queries as in multi-core database cracking and multi-core stochastic cracking.

In addition, all index choosing strategies have similar performance on workloads where attributes are queried on random values (Figures 13(a) and (c)), because indices are already fine-grained in such cases (even when some indices are refined more than others in a skewed workload). However, in case of queries on periodic values (Figures 13(b) and (d)) the random choice (W4) shows a clear performance benefit compared to the rest of the strategies, because it refines indices with big unindexed partitions, and proves to be a robust design decision.

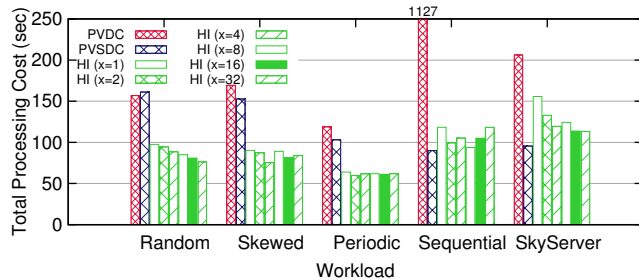


Figure 15: Performance of holistic indexing improves while the number of index refinements x increases.

5.5 Design Decisions

Holistic Worker Thread Refinements. In this experiment we demonstrate how the tuning parameter x , i.e., the number of index refinements per worker (Figure 2), affects the workload performance. We test five workloads that consist of 10^3 queries on a relation with 10 attributes as in Section 5.3 (Figure 12). We vary the

number of index refinements each holistic worker thread does from 1 to 32 and we compare holistic indexing with multi-core variants of database cracking and stochastic cracking. Figure 15 shows the results. The more index refinements each thread does, the bigger the benefit for holistic indexing because more pieces are created and thus future queries need to refine smaller pieces, touching less data. However, when we increase the number of index refinements from 16 to 32, performance does not improve significantly, because in both cases indices converge very fast to optimal ones. Thus, we use 16 as the number of index refinements that each holistic worker thread does in all our experiments.

5.6 TPC-H

In our next experiment, we evaluate holistic indexing on the standard database benchmark, TPC-H [3]. We compare against offline indexing and plain scans. We use scale factor 10 and we test with Queries 1, 6, and 12. For each query type, we created a sequence of 30 variations using the random query generator distributed with the benchmark. For offline indexing, we created the proper column-store projections by pre-sorting the data depending on each query individually, i.e., we created the perfect projection for each query. Specifically, for Queries 1 and 6 we created a copy of the Lineitem table sorted on the `l_shipdate` attribute. For Query 12 we created a copy of the Lineitem table sorted on the `l_receiptdate` attribute.

Figure 14 depicts the results. For all cases, holistic indexing brings a significant advantage, resulting in a robust and stable performance across all queries. The first query is slower as it creates the first adaptive indices which implies extra data copying but after that all queries perform significantly better than plain MonetDB. Holistic indexing matches offline indexing without having to incur the high offline indexing cost and without requiring any workload knowledge. The pre-sorting cost for all queries is 8 seconds. For Query 12, it turns out that pre-sorting does not help. This happens because even though we may improve the selection by pre-sorting the Lineitem table, it turns out we hurt the join between the Lineitem and the Orders table. This is because in the base data, the Lineitem table contains the order date ordered and this can be exploited during the join. With holistic indexing we do not face this problem, because the initial order changes only partially.

5.7 Updates

So far we tested read-only workloads. In this experiment we demonstrate that holistic indexing maintains its nice properties in workloads where read-only queries interleave with write queries. We test two scenarios. In the first scenario (High Frequency Low Volume - HFLV), 10 inserts arrive every 10 queries. In the second scenario (Low Frequency High Volume - LFHV), 100 inserts arrive every 100 queries. In both scenarios the workload consists of 500 select project range queries on a single attribute A and 500 inser-

tions in total on a single attribute A . While all queries are processed sequentially one after the other, the 11th query arrives 20 seconds after the 10th query resulting in idle time of 20 seconds in the system. Attribute A consists of 10^9 uniformly distributed integers.

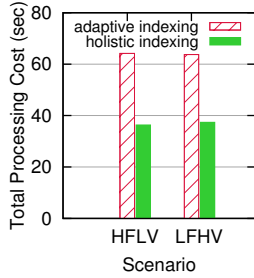


Figure 16: Updates.

adaptive indexing against holistic indexing with a single worker that refines the index only during idle time. In holistic indexing, auxiliary index refinement actions also cause insertions to be merged. In this way, holistic indexing not only refines the indices but also consumes pending insertions, which speeds up future queries even more and all that by exploiting idle CPU resources in parallel with query processing. Figure 16 shows the results. In both scenarios, holistic indexing maintains its advantage over adaptive indexing; it is not affected by updates and still provides roughly a 50% improvement.

5.8 Varying Number of Clients

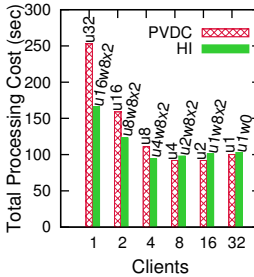


Figure 17: Varying number of clients.

In this experiment we test the performance of holistic indexing with varying number of concurrent clients in the system. Our workload consists of 1024 queries on a relation with 10 attributes. We vary the number of concurrent clients between 1 and 32, where 32 is the number of CPU cores in our machine. Figure 17 shows that holistic indexing brings a big benefit in case of a few clients. The labels on top of the bars indicate the distribution of the available threads across user queries and holistic workers (similar to Figure 7). When the number of clients increases, holistic indexing does not bring significant benefits, because it easily detects these cases as it monitors the CPU load continuously and so it is triggered only if the load is below a threshold.

6. CONCLUSIONS

Ad-hoc environments become more and more common in the big data era, where we have little workload knowledge and time to properly tune a database system. In this paper, we present holistic indexing; it continuously adapts to workload changes and makes maximum use of available CPU resources, requiring zero human administration. Holistic indexing continuously monitors the workload and refines the physical design with lightweight actions whenever spare CPU cycles are detected. We implemented holistic indexing in an open-source column-store DBMS and we demonstrate that even though holistic indexing runs in the background and on-the-fly, i.e., while processing user queries, it does not affect the running queries negatively; instead by continuously refining the physical design based on the workload and by exploiting spare CPU cycles it brings substantial improvements. A detailed experimental analysis shows that holistic indexing maintains all the good benefits and properties of adaptive indexing, while it brings an additional performance gain which is typically a factor of 2.

Holistic indexing opens a promising research path. Studying holistic indexing for traditional row stores or emerging hybrid systems is an interesting topic. While major database vendors support hybrid storage layouts, auto-tuning tools may take into account holistic indexing and depending on the workload decide which indices to build up front and which indices to assign to holistic indexing. In addition, even though we considered primarily indices for improving selections in this paper there is room to consider other indices that may be improved by holistic indexing such as join indices. Furthermore, this work and most recent work on adaptive indexing has focused on column-stores that work under the bulk processing model with the data stored in fixed width and dense arrays which allows for efficient tight for loops. However, work on partial sideways cracking [29] and more recently chunked coarse granular indexing [8] shows that even more benefits for adaptive indexing are possible if we drop the dense arrays restriction, especially since this enables effective parallelization [8]. This has side-effects of course to the whole architecture of a database system as it affects the design and access patterns of all relational operators. Thus, it is very interesting to investigate whether it is beneficial to drop certain restrictions and reconsider the whole architecture and data flow in column-stores to accommodate further adaptive indexing benefits. Finally, in this paper we focused on improving the overall response time of a workload but further interesting directions may rise once we consider energy consumption, targeting systems with the maximum performance but within specific energy bounds.

Acknowledgments We would like to thank Victor Alvarez, Felix Martin Schuhknecht, Jens Dittrich and Stefan Richter for kindly providing the source code of the chunked coarse granular adaptive indexing algorithm. We would also like to thank the anonymous reviewers for their truly insightful remarks.

7. REFERENCES

- [1] Sort-Merge Joins. <http://www.systems.ethz.ch/projects/paralleljoins>.
- [2] Sloan Digital Sky Survey (SkyServer). <http://cas.sdss.org/>.
- [3] TPC Benchmark H. <http://www.tpc.org/tpch/>.
- [4] D. Abadi, P. A. Boncz, S. Harizopoulos, S. Idreos, and S. Madden. The Design and Implementation of Modern Column-Oriented Database Systems. *Foundations and Trends in Databases*, 5(3):197–280, 2013.
- [5] S. Agrawal, S. Chaudhuri, L. Kollár, A. P. Marathe, V. R. Narasayya, and M. Syamala. Database Tuning Advisor for Microsoft SQL Server 2005. In *Proceedings of the 30th International Conference on Very Large Data Bases (VLDB)*, pages 1110–1121, 2004.
- [6] I. Alagiannis, R. Borovica, M. Branco, S. Idreos, and A. Ailamaki. NoDB: Efficient Query Execution on Raw Data Files. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 241–252, 2012.
- [7] I. Alagiannis, S. Idreos, and A. Ailamaki. H2O: A Hands-free Adaptive Store. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 1103–1114, 2014.
- [8] V. Alvarez, F. M. Schuhknecht, J. Dittrich, and S. Richter. Main Memory Adaptive Indexing for Multi-core Systems. In *Proceedings of the 10th International Workshop on Data Management on New Hardware (DaMoN)*, 2014.
- [9] C. Balkesen, G. Alonso, J. Teubner, and M. T. Özsu. Multi-Core, Main-Memory Joins: Sort vs. Hash Revisited. *Proceedings of the Very Large Data Bases Endowment (PVLDB)*, 7(1):85–96, 2013.

- [10] N. Bruno and S. Chaudhuri. An Online Approach to Physical Design Tuning. In *Proceedings of the IEEE 23rd International Conference on Data Engineering (ICDE)*, pages 826–835, 2007.
- [11] S. Chaudhuri and V. R. Narasayya. An Efficient Cost-Driven Index Selection Tool for Microsoft SQL Server. In *Proceedings of the 23rd International Conference on Very Large Data Bases (VLDB)*, pages 146–155, 1997.
- [12] S. Chaudhuri and V. R. Narasayya. AutoAdmin “What-if” Index Analysis Utility. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 367–378, 1998.
- [13] G. P. Copeland and S. N. Khoshafian. A Decomposition Storage Model. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 268–279, 1985.
- [14] B. Dageville, D. Das, K. Dias, K. Yagoub, M. Zaït, and M. Ziauddin. Automatic SQL Tuning in Oracle 10g. In *Proceedings of the 30th International Conference on Very Large Data Bases (VLDB)*, pages 1098–1109, 2004.
- [15] J. Dittrich, J.-A. Quiané-Ruiz, A. Jindal, Y. Kargin, V. Setty, and J. Schad. Hadoop++: Making a Yellow Elephant Run Like a Cheetah (Without It Even Noticing). *Proceedings of the Very Large Data Bases Endowment (PVLDB)*, 3(1):518–529, 2010.
- [16] G. Graefe, F. Halim, S. Idreos, H. A. Kuno, and S. Manegold. Concurrency Control for Adaptive Indexing. *Proceedings of the Very Large Data Bases Endowment (PVLDB)*, 5(7):656–667, 2012.
- [17] G. Graefe, F. Halim, S. Idreos, H. A. Kuno, S. Manegold, and B. Seeger. Transactional Support for Adaptive Indexing. *The Very Large Data Bases Journal (VLDB J.)*, 23(2):303–328, 2014.
- [18] G. Graefe, S. Idreos, H. A. Kuno, and S. Manegold. Benchmarking Adaptive Indexing. In *Proceedings of the 2nd Technology Conference on Performance Evaluation and Benchmarking (TPCTC)*, pages 169–184, 2010.
- [19] G. Graefe and H. Kuno. Adaptive Indexing for Relational Keys. In *Workshops Proceedings of the IEEE 26th International Conference on Data Engineering (ICDE)*, pages 69–74, 2010.
- [20] G. Graefe and H. A. Kuno. Self-Selecting, Self-Tuning, Incrementally Optimized Indexes. In *Proceedings of the 13th International Conference on Extending Database Technology (EDBT)*, pages 371–381, 2010.
- [21] F. Halim, S. Idreos, P. Karras, and R. H. C. Yap. Stochastic Database Cracking: Towards Robust Adaptive Indexing in Main-Memory Column-Stores. *Proceedings of the Very Large Data Bases Endowment (PVLDB)*, 5(6):502–513, 2012.
- [22] N. Hardavellas, M. Ferdman, B. Falsafi, and A. Ailamaki. Near-Optimal Cache Block Placement with Reactive Nonuniform Cache Architectures. In *IEEE Micro*, page 29, 2010.
- [23] S. Harizopoulos and A. Ailamaki. A Case for Staged Database Systems. In *Proceedings of the Conference on Innovative Data Systems Research (CIDR)*, 2003.
- [24] S. Harizopoulos, V. Shkapenyuk, and A. Ailamaki. QPipe: A Simultaneously Pipelined Relational Query Engine. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 383–394, 2005.
- [25] S. Idreos, I. Alagiannis, R. Johnson, and A. Ailamaki. Here are my Data Files. Here are my Queries. Where are my Results? In *Proceedings of the 5th Biennial Conference on Innovative Data Systems Research (CIDR)*, pages 57–68, 2011.
- [26] S. Idreos, F. Groffen, N. Nes, S. Manegold, K. S. Mullender, and M. L. Kersten. MonetDB: Two Decades of Research in Column-oriented Database Architectures. *IEEE Data Engineering Bulletin*, 35(1):40–45, 2012.
- [27] S. Idreos, M. L. Kersten, and S. Manegold. Database Cracking. In *Proceedings of the 3rd Biennial Conference on Innovative Data Systems Research (CIDR)*, pages 68–78, 2007.
- [28] S. Idreos, M. L. Kersten, and S. Manegold. Updating a Cracked Database. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 413–424, 2007.
- [29] S. Idreos, M. L. Kersten, and S. Manegold. Self-Organizing Tuple Reconstruction in Column-Stores. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 297–308, 2009.
- [30] S. Idreos and E. Liarou. dbTouch: Analytics at your Fingertips. In *Proceedings of the 6th Biennial Conference on Innovative Data Systems Research (CIDR)*, 2013.
- [31] S. Idreos, S. Manegold, H. A. Kuno, and G. Graefe. Merging What’s Cracked, Cracking What’s Merged: Adaptive Indexing in Main-Memory Column-Stores. *Proceedings of the Very Large Data Bases Endowment (PVLDB)*, 4(9):585–597, 2011.
- [32] H. V. Jagadish, A. Chapman, A. Elkiss, M. Jayapandian, Y. Li, A. Nandi, and C. Yu. Making Database Systems Usable. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 13–24, 2007.
- [33] R. Johnson, I. Pandis, N. Hardavellas, A. Ailamaki, and B. Falsafi. Shore-MT: A Scalable Storage Manager for the Multicore Era. In *Proceedings of the 12th International Conference on Extending Database Technology (EDBT)*, pages 24–35, 2009.
- [34] R. Kallman, H. Kimura, J. Natkins, A. Pavlo, A. Rasin, S. B. Zdonik, E. P. C. Jones, S. Madden, M. Stonebraker, Y. Zhang, J. Hugg, and D. J. Abadi. H-store: A High-Performance, Distributed Main Memory Transaction Processing System. *Proceedings of the Very Large Data Bases Endowment (PVLDB)*, 1(2):1496–1499, 2008.
- [35] A. Kemper and T. Neumann. HyPer: A Hybrid OLTP&OLAP Main Memory Database System Based on Virtual Memory Snapshots. In *Proceedings of the IEEE 27th International Conference on Data Engineering (ICDE)*, pages 195–206, 2011.
- [36] M. L. Kersten, S. Idreos, S. Manegold, and E. Liarou. The Researcher’s Guide to the Data Deluge: Querying a Scientific Database in Just a Few Seconds. *Proceedings of the Very Large Data Bases Endowment (PVLDB)*, 4(12):1474–1477, 2011.
- [37] U. Khurana and A. Deshpande. HiNGE: Enabling Temporal Network Analytics at Scale. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 1089–1092, 2013.
- [38] C. Kim, E. Sedlar, J. Chhugani, T. Kaldewey, A. D. Nguyen, A. D. Blas, V. W. Lee, N. Satish, and P. Dubey. Sort vs. Hash Revisited: Fast Join Implementation on Modern Multi-Core CPUs. *Proceedings of the Very Large Data Bases Endowment (PVLDB)*, 2(2):1378–1389, 2009.

- [39] A. Kumar, F. Niu, and C. Ré. Hazy: Making it Easier to Build and Maintain Big-Data Analytics. *Communications of the ACM*, 56(3):40–49, 2013.
- [40] N. Laptev, K. Zeng, and C. Zaniolo. Very Fast Estimation for Result and Accuracy of Big Data Analytics: The EARL System. In *Proceedings of the IEEE 29th International Conference on Data Engineering (ICDE)*, pages 1296–1299, 2013.
- [41] Y. Li and J. M. Patel. BitWeaving: Fast Scans for Main Memory Data Processing. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 289–300, 2013.
- [42] E. Liarou and S. Idreos. dbTouch in Action Database Kernels for Touch-based Data Exploration. In *Proceedings of the IEEE 30th International Conference on Data Engineering (ICDE)*, pages 1262–1265, 2014.
- [43] M. Lühring, K.-U. Sattler, K. Schmidt, and E. Schallehn. Autonomous Management of Soft Indexes. In *Proceedings of the 2nd International Workshop on Self-Managing Data Bases (SMDB)*, pages 450–458, 2007.
- [44] H. Pirk, E. Petraki, S. Idreos, S. Manegold, and M. L. Kersten. Database Cracking: Fancy Scan, not Poor Man’s Sort! In *Proceedings of the 10th International Workshop on Data Management on New Hardware (DaMoN)*, 2014.
- [45] O. Polychroniou and K. A. Ross. A Comprehensive Study of Main-memory Partitioning and Its Application to Large-scale Comparison- and Radix-sort. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 755–766, 2014.
- [46] S. Richter, J.-A. Quian-Ruiz, S. Schuh, and J. Dittrich. Towards Zero-Overhead Static and Adaptive Indexing in Hadoop. *The Very Large Data Bases Journal (VLDB J.)*, 23(3):469–494, 2013.
- [47] K. Schnaitter, S. Abiteboul, T. Milo, and N. Polyzotis. COLT: Continuous On-Line Tuning. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 793–795, 2006.
- [48] F. M. Schuhknecht, A. Jindal, and J. Dittrich. The Uncracked Pieces in Database Cracking. *Proceedings of the Very Large Data Bases Endowment (PVLDB)*, 7(2):97–108, 2013.
- [49] J. Zhou and K. A. Ross. Implementing Database Operations Using SIMD Instructions. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 145–156, 2002.
- [50] D. C. Zilio, J. Rao, S. Lightstone, G. M. Lohman, A. J. Storm, C. Garcia-Arellano, and S. Fadden. DB2 Design Advisor: Integrated Automatic Physical Database Design. In *Proceedings of the 30th International Conference on Very Large Data Bases (VLDB)*, pages 1087–1097, 2004.
- [51] K. Zoumpatianos, S. Idreos, and T. Palpanas. Indexing for Interactive Exploration of Big Data Series. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pages 1555–1566, 2014.
- [52] *MonetDB*. www.monetdb.org.