



Adaptive partitioning and indexing for in situ query processing

Matthaios Olma¹ · Manos Karpathiotakis² · Ioannis Alagiannis³ · Manos Athanassoulis⁴ · Anastasia Ailamaki¹

Received: 1 December 2018 / Revised: 1 July 2019 / Accepted: 5 October 2019
© Springer-Verlag GmbH Germany, part of Springer Nature 2019

Abstract

The constant flux of data and queries alike has been pushing the boundaries of data analysis systems. The increasing size of raw data files has made data loading an expensive operation that delays the data-to-insight time. To alleviate the loading cost, in situ query processing systems operate directly over raw data and offer instant access to data. At the same time, analytical workloads have increasing number of queries. Typically, each query focuses on a constantly shifting—yet small—range. As a result, minimizing the workload latency requires the benefits of indexing in in situ query processing. In this paper, we present an online partitioning and indexing scheme, along with a partitioning and indexing tuner tailored for in situ querying engines. The proposed system design improves query execution time by taking into account user query patterns, to (i) partition raw data files *logically* and (ii) build lightweight *partition-specific* indexes for each partition. We build an in situ query engine called Slalom to showcase the impact of our design. Slalom employs adaptive partitioning and builds non-obtrusive indexes in different partitions on-the-fly based on lightweight query access pattern monitoring. As a result of its lightweight nature, Slalom achieves efficient query processing over raw data with minimal memory consumption. Our experimentation with both microbenchmarks and real-life workloads shows that Slalom outperforms state-of-the-art in situ engines and achieves comparable query response times with fully indexed DBMS, offering lower cumulative query execution times for query workloads with increasing size and unpredictable access patterns.

Keywords Online tuning · Adaptive indexing · Logical partitioning

1 Introduction

Data-intensive applications in various domains generate and collect massive amounts of data at a rapid pace. New research fields and applications (e.g., network monitoring, sensor data management, clinical studies, etc.) emerge and require

broader data analysis functionality to rapidly gain deeper insights from the available data. In practice, analyzing such datasets become costlier as their size grows.

Big data, small queries The trend of exponential data growth due to intense data generation and data collection is expected to persist. However, recent studies of the data analysis workloads show that typically only a small subset of the data is relevant and ultimately used by analytical and/or exploratory workloads [1,17]. In addition, modern business and scientific applications require interactive data access, which is characterized by *no or little a priori workload knowledge* and constant *workload shifting* both in terms of the attributes projected and the ranges selected from the dataset.

The cost of loading, indexing, and tuning Traditional data management systems (DBMS) require the costly steps of *data loading*, *physical design decisions*, and then *index building* in order to offer interactive access over large datasets. Given the data sizes involved, any transformation, copying, and preparation steps over the data introduce substantial delays before the data can be utilized, queried, and provide useful insights [2,5,36]. The lack of a priori knowledge of

✉ Matthaios Olma
matthaios.olma@epfl.ch

Manos Karpathiotakis
manos@fb.com

Ioannis Alagiannis
ioalagia@microsoft.com

Manos Athanassoulis
mathan@bu.edu

Anastasia Ailamaki
anastasia.ailamaki@epfl.ch

¹ EPFL, Lausanne, Switzerland

² Facebook, London, UK

³ Microsoft, Redmond, WA, USA

⁴ Boston University, Boston, MA, USA

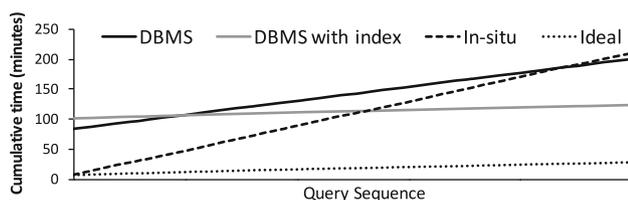


Fig. 1 Ideally, in situ data analysis should be able to retrieve only the relevant data for each query after the initial table scan (ideal—dotted line). In practice, in situ query processing avoids the costly phase of data loading (dashed line); however, as the number of the queries increases, the initial investment for full index on a DBMS pays off (the dashed line meets the gray line)

the workload makes the physical design decisions virtually impossible because cost-based advisors rely heavily on past or sample workload knowledge [3,16,22,29,68]. The workload shifts observed in the interactive setting of exploratory workloads can nullify investments toward indexing and other auxiliary data structures, since frequently, they depend on the actual data values and the knowledge generated by the ongoing analysis.

Querying raw data files is not enough Recent efforts opt to query directly raw files [2,5,12,18,31,42] to reduce the data-to-insight cost. These in situ systems avoid the costly initial data loading step and allow the execution of declarative queries over external files without duplicating or “locking” data in a proprietary database format. Further, they concentrate on reducing costs associated with raw data accesses (e.g., parsing and converting data fields) [5,18,42]. Finally, although recent scientific data management approaches index raw data files using file-embedded indexes, they do it in a workload-oblivious manner, or requiring full a priori workload knowledge [12,67]. Hence, they bring back the cost of full index building, in the raw data querying paradigm, negating part of the benefits of avoiding data loading.

Figure 1 shows what the ideal in situ query performance should be (dotted line). After the unavoidable first table scan, ideally, in situ queries need to access only data relevant to the currently executed query. The figure also visualizes the benefits of state-of-the-art in situ query processing when compared with a full DBMS. The y-axis shows the cumulative query latency, for an increasing number of queries with fixed selectivity on the x-axis. By avoiding the costly data loading phase, the in situ query execution system (dashed line) can start answering queries very quickly. On the other hand, when a DBMS makes an additional investment on full DBMS indexing (solid gray line), it initially increases significantly the data-to-query latency; however, later it pays off as the number of queries issued over the same (raw) dataset increases. Eventually, the cumulative query latency for an in situ approach becomes larger than the latency of a DBMS equipped with indexing. When operating over raw data, *ide-*

ally, we want after the initial—unavoidable—table scan to collect enough metadata to allow future queries to access only the useful part of the dataset.

Adaptive partitioning and fine-grained indexing We use the first table scan to generate partitioning and lightweight indexing hints which are further refined by the data accesses of (only a few) subconsequent queries. During this refinement process, the dataset is partially indexed in a dynamic fashion adapting to three key workload characteristics: (i) data distribution, (ii) query type (e.g., point query, range query), and (iii) projected attributes. Workload shifts lead to varying selected value ranges, selectivity, which dataset areas are relevant for a query, and projected attributes.

This paper proposes an online partitioning and indexing tuner for in situ query processing which, when plugged into a raw data query engine, offers *fast queries over raw data files*. The tuner reduces data access cost by: (i) *logically partitioning* a raw dataset to break it into smaller manageable chunks without physical restructuring and (ii) choosing *appropriate indexing strategies over each logical partition* to provide efficient data access. The tuner adapts the partitioning and indexing scheme as a side effect of executing the query workload. It continuously collects information regarding the values and access frequency of queried attributes at runtime. Based on this information, it uses a randomized online algorithm to define logical partitions. For each logical partition, the tuner estimates the cost-benefit of building partition-local index structures considering both approximate membership indexing (i.e., Bloom filters and zonemaps) and full indexing (i.e., bitmaps and B⁺ trees). By allowing fine-grained indexing decisions, our proposal defers the decision of the index shape to the level of each partition rather than the overall relation. This has two positive side effects. First, there is no costly indexing investment that might be unnecessary. Second, any indexing effort is tailored to the needs of data accesses on the corresponding range of the dataset.

Efficient in situ query processing with Slalom We integrate our online partitioning and indexing tuner to an in situ query processing prototype system, *Slalom*, which combines the tuner with a state-of-the-art raw data query executor. Slalom is further augmented with index structures and uses the tuner to decide how to partition and which index or indexes to build for each partition. In particular, Slalom logically splits raw data into partitions and selects which fine-grained index to build, per partition based on how “hot” (i.e., frequently accessed) each partition is, and what types of queries target each partition. Furthermore, Slalom populates binary caches (of data converted from raw to binary) to further boost performance. Slalom adapts to workload shifts by adjusting the current partitioning and indexing scheme using a randomized cost-based decision algorithm. Overall, the logical partitions and the indexes that Slalom builds over each partition provide

performance enhancements without requiring expensive full data indexing nor data file reorganization, all while adapting to workload changes.

Contributions This paper makes the following contributions:

- We present a logical partitioning scheme of raw data files that enables fine-grained indexing decisions at the level of each partition. As a result, lightweight per-partition indexing provides near-optimal data access.
- The lightweight partitioning allows our approach to maintain the benefits of in situ approaches. In addition, the granular way of indexing (i) brings the benefit of indexing to in situ query processing, (ii) having low index building cost, and (iii) small memory footprint. These benefits are highlighted as the partitioning and indexing decisions are refined on-the-fly using an online randomized algorithm.
- We enable both in-place and append-like updates for in situ query processing. We exploit specialized hardware (GPUs and CRC checksum units) to reduce update recognition cost and minimize changes to partitioning and indexing, overall minimizing the query execution overhead in the presence of updates.
- We integrate our partitioning and indexing tuner into our prototype state-of-the-art in situ query engine *Slalom*. We use synthetic and real-life workloads to compare the query latency of (i) *Slalom*, (ii) a traditional DBMS, (iii) a state-of-the-art in situ query processing engine, and (iv) adaptive indexing (cracking). Our experiments show that, even when excluding the data loading cost, *Slalom* offers the fastest cumulative query latency. In particular, *Slalom* outperforms (a) state-of-the-art disk-based approaches by one order of magnitude, (b) state-of-the-art in-memory approaches by $3.7\times$ (with $2.45\times$ smaller memory footprint), and (c) adaptive indexing by 19% (having $1.93\times$ smaller memory footprint). Finally, we examine the performance of *Slalom* in the presence of both in-place and append-like updates.

To our knowledge, *Slalom* is the first approach that proposes the use of a randomized online algorithm to select which workload-tailored index structures should be built per partition of the data file. This approach reduces index building time and provides minimal decision time.

Outline The remainder of this paper is organized as follows: Sect. 2 provides an overview of related work. Sect. 3 presents the architecture of *Slalom* and gives an overview of its design. Section 4 presents the online tuner and describes its partitioning and indexing cost models. Section 4.3 presents the techniques enabling efficient data updates for in situ query

processing. We experimentally demonstrate the benefits of *Slalom* in Sect. 5, and we conclude in Sect. 6.

2 Related work

In recent years, many research efforts re-design the traditional data management architecture to address the challenges and opportunities associated with dynamic workloads and interactive data access. In this section, we discuss research approaches related to *Slalom* and highlight how *Slalom* pushes the state of the art.

Queries over raw data Data loading accounts for a large fraction of overall workload execution time in both the DBMS and Hadoop ecosystems [31]. NoDB treats raw data files as native storage of the DBMS and introduces auxiliary data structures (positional maps and caches) to reduce the expensive parsing and tokenization costs of raw data access [5]. ViDa introduces code-generated access paths and data pipeline to adapt the query engine to the underlying data formats and layouts, and to the incoming queries [40–42]. Data Vaults [36,38] and SDS/Q [12] perform analysis over scientific array-based file formats. SCANRAW [18] uses parallelism to mask the increased CPU processing costs associated with raw data accesses during in situ data processing. In situ DBMS approaches either rely on accessing the data via full table scans or require a priori workload knowledge and enough idle time to create the proper indexes. The mechanisms of *Slalom* are orthogonal to these systems and can augment their design by enabling data skipping and indexed accesses while constantly adapting its indexing and partitioning schemes to queries.

Hadoop-based systems such as Hive [64] can access raw data stored in HDFS. While such frameworks internally translate queries to MapReduce jobs, other systems follow a more traditional massive parallel processing (MPP) architecture to offer SQL-on-Hadoop functionality [45,49]. Hybrid approaches such as invisible loading [2] and polybase [21] propose co-existence of a DBMS and a Hadoop cluster, transferring data between the two when needed. SQL Server PDW [24] and AsterixDB [6] propose indexes for data stored in HDFS and, in general, for external data. The key techniques of *Slalom* can also be applied in a Hadoop-based environment. SQL Server PDW and AsterixDB build secondary indexes over HDFS files. The techniques used by *Slalom*, on the other hand, improve system scalability by reducing the size of the index and building memory efficient indexes per file partition.

On the other side of raw data querying, instant Loading [51] parallelizes the loading process for main-memory DBMS, offering bulk loading at near-memory-bandwidth speed. Similarly to instant Loading, *Slalom* uses data parsing

with hardware support for efficient raw data access. Instead of loading all data, however, Slalom exploits workload locality to adaptively create a fine-grained indexing scheme over raw data and gradually reduce I/O and access costs, all while operating within a modest memory budget.

Database partitioning A table can be physically subdivided into smaller disjoint sets of tuples (partitions), allowing tables to be stored, managed, and accessed at a finer level of granularity [46].

Offline partitioning approaches [3,27,53,68] present physical design tools that automatically select the proper partition configuration for a given workload to improve performance. Online partitioning [37] monitors and periodically adapts the database partitions to fit the observed workload. Furtado et al. [23] combine physical and virtual partitioning to fragment and dynamically tune partition sizes for flexibility in intra-query parallelism. Shinobi [66] clusters hot data in horizontal partitions which it then indexes, while Sun et al. [63] use a bottom-up clustering framework to offer an approximate solution for the partition identification problem.

Physical reorganization, however, is not suitable for data file repositories due to its high cost and the immutable nature of the files. Slalom presents a non-intrusive, flexible partitioning scheme that creates logical horizontal partitions by exploiting data skew. Additionally, Slalom continuously refines its partitions during query processing without requiring a priori workload knowledge.

Database indexing There is a vast collection of index structures with different capabilities, performance, and initialization/maintenance overheads [9,10,35,44]. This paper uses representative index structures from the two categories: (i) value-position and (ii) value-existence indexes, that offer good indexing for point and range queries. Value-position indexes include the B⁺ tree and hash indexes and their variations [8]. Common value-existence indexes are Bloom filters [13], bitmap indexes [11,52,61], and zonemaps [50]. They are lightweight and can provide the information whether a value is present in a given dataset. Value-existence indexes are frequently used in scientific workloads [19,62,67]. Slalom builds main-memory auxiliary structures (i) rapidly, (ii) with small footprint, and (iii) without a priori workload knowledge. That way, it enables low data-to-insight latency without hurting the performance of long-running workloads, for which indexing is typically more useful.

Online indexing Physical design decisions made before workload execution can also be periodically re-evaluated. COLT [59] continuously monitors the workload and periodically creates new indexes and/or drops unused ones. COLT adds overhead on query execution because it obtains cost estimations from the optimizer at runtime. A “lighter” approach requiring fewer calls to the optimizer has also been proposed in recent literature [15]. Slalom also focuses on the problem

of selecting an effective set of indexes and builds indexes on partition granularity. It populates indexes during query execution in a pipelined fashion instead of triggering a stand-alone index building phase. Slalom aims to minimize the cost of index construction decisions and the complexity of the costing algorithm.

Adaptive indexing In order to avoid the full cost of indexing before workload execution, adaptive indexing incrementally refines indexes during query processing. In the context of in-memory column stores, database cracking approaches [25,32–34,56] create a duplicate of the indexed column and incrementally sort it according to the incoming workload, thus reducing memory access. HAIL proposes an adaptive indexing approach for MapReduce systems [57]. ARF is an adaptive value-existence index similar to Bloom filters, yet useful for range queries [7]. Similarly to adaptive indexing, Slalom does not index data upfront and builds indexes during query processing and continuously adapts to the workload characteristics. However, contrary to adaptive indexing that duplicates the whole indexed attribute upfront, Slalom’s gradual index building allows its indexes to have small memory footprint by indexing both the targeted value ranges and the targeted attributes.

3 The SLALOM system

Slalom uses adaptive partitioning and indexing to provide inexpensive index support for in situ query processing while adapting to workload changes. Slalom accelerates query processing by skipping data and minimizes data access cost when this access is unavoidable. At the same time, it operates directly on the original data files without need for physical restructuring (i.e., copying, sorting).

Slalom incorporates state-of-the-art in situ querying techniques and enhances them with logical partitioning and fine-grained indexing, thereby reducing the amounts of accessed data. To remain effective despite workload shifts, Slalom introduces an online partitioning and indexing tuner, which calibrates and refines logical partitions and secondary indexes based on data and query statistics. Slalom treats data files as relational tables to facilitate the processing of read-only and append-like workloads. The rest of this section focuses on the architecture and implementation of Slalom.

3.1 Architecture

Figure 2 presents the architecture of Slalom. Slalom combines an online partitioning and indexing tuner with a query executor featuring in situ querying techniques. The core components of the tuner are the *Partition Manager*, which is responsible for creating logical partitions over the data files,

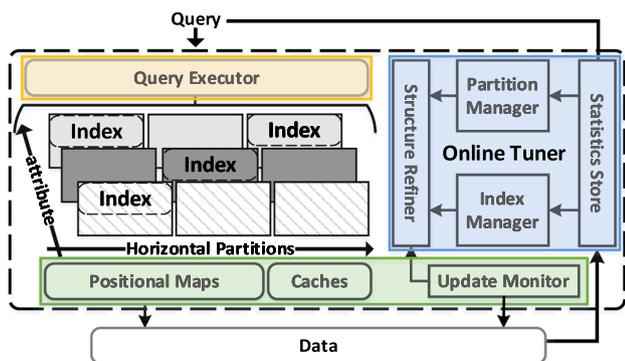


Fig. 2 The architecture of Slalom

and the *Index Manager*, which is responsible for creating and maintaining indexes over partitions. The tuner collects statistics regarding the data and query access patterns and stores them in the *Statistics Store*. Based on those statistics, the *Structure Refiner* evaluates the potential benefits of alternative configurations of partitions and indexes. Furthermore, Slalom uses in situ querying techniques to access data. Specifically, Slalom uses auxiliary structures (i.e., positional maps and caches) which minimize raw data access cost. During query processing, the *Query Executor* utilizes the available data access paths and orchestrates the execution of the other components. Finally, the *Update Monitor* examines whether a file has been modified and adjusts the data structures of Slalom accordingly.

Slalom scope The techniques of Slalom are applicable to any tabular dataset. Specifically, the scan operator of Slalom uses a different specialized parser for each underlying data format. This work concentrates on queries over delimiter-separated textual CSV files, because CSV is the most popular structured textual file format. Still, the yellow- and blue-coded components of Fig. 2 are applicable over binary files, which are the typical backend of databases and are also frequently used in scientific applications (e.g., high-energy physics, DNA

sequencing, GIS). Regarding query types, Slalom concentrates on efficient access of different raw data files and enables queries containing filters on different attributes. Slalom, in its current format, does not support arbitrary joins and nested SQL queries. However, we assume that any query involving nested queries or joins can be flattened and Slalom can perform filtering over the resulting underlying data. We discuss further Slalom’s extensibility in Sect. 3.4.

Reducing data access cost Slalom launches queries directly over the original raw data files, without altering or duplicating the files by ingesting them in a DBMS. That way, Slalom avoids the initialization cost induced by loading and offers instant data access. Similarly to state-of-the-art in situ query processing approaches [5,18], Slalom mitigates the overheads of parsing and tokenizing textual data with positional maps [5] and partial data caching.

PMs are populated on-the-fly and maintain *structural* information about an underlying textual file; they keep the positions of various file attributes. This information is used during query processing to “jump” to the exact position of an attribute or as close as possible to an attribute, significantly reducing the cost of tokenizing and parsing when a tuple is accessed. Furthermore, Slalom builds binary caches of fields that are already converted to binary to reduce parsing and data type conversion costs of future accesses.

Statistics Store Slalom collects statistics during query execution and utilizes them to (i) detect workload shifts and (ii) enable the tuner to evaluate partitioning and index configurations. Table 1 summarizes the statistics about *Data* and *Queries* that Slalom gathers per data file. *Data statistics* are updated after every partitioning action and include the per-partition standard deviation (dev_i) of values, mean (m_i), max (max_i), and min (min_i) values. Additionally, Slalom keeps as global statistics the physical page size ($Size_{page}$) and file size ($Size_{file}$). Regarding *Query statistics*, Slalom maintains the number of queries since the last access (LA_i), the percentage of queries accessing each partition (access frequency AF_i),

Table 1 Statistics collected by Slalom per data file during query processing and used to decide (i) which logical partitions to create and (ii) select the appropriate matching indexes

Data (partition i)	m_i	Mean value
	min_i	Min value
	max_i	Max value
	dev_i	Standard deviation
	DV_i	#distinct values
Data (global)	$Size_{page}$	Physical page size
	$Size_{file}$	File size
Queries (partition i)	$C_{i_{build}}$	Index building cost
	$C_{i_{fullscan}}$	Full scan cost
	LA_i	#queries since last access
	AF_i	Partition access frequency
	sel_i	Average selectivity (0.0–1.0)

and the average query selectivity (sel_i). Finally, the full scan cost over a partition ($C_{i_{fullscan}}$) and the indexing cost for a partition ($C_{i_{build}}$) are calculated by considering the operator's data accesses.

Partition Manager The Partition Manager recognizes patterns in the dataset and logically divides the file into contiguous non-overlapping chunks to enable fine-grained access and indexing. The Partition Manager specifies a logical partitioning scheme for each attribute in a relation. Each partition is internally represented by its starting and ending byte within the original file. The logical partitioning process starts the first time a query accesses an attribute. The Partition Manager triggers the Structure Refiner to iteratively fine-tune the partitioning scheme with every subsequent query. All partitions progressively reach a state in which there is no benefit from further partitioning. The efficiency of a partitioning scheme depends highly on the data distribution and the query workload. Therefore, the Partition Manager adjusts the partitioning scheme based on value cardinality (details in Sect. 4.1).

Index Manager The Index Manager estimates the benefit of an index over a partition and suggests the most promising combination of indexes for a given attribute/partition. For every new index configuration, the Index Manager invokes the Structure Refiner to build the selected indexes during the execution of the next query. Every index corresponds to a specific data partition. Depending on the access pattern of an attribute and the query selectivity, a single partition may have multiple indexes. Slalom chooses indexes from two categories based on their capabilities: (i) *value-existence* indexes, which respond whether a value exists in a dataset, and (ii) *value-position* indexes, which return the positions of a value within the file. The online nature of Slalom imposes a significant challenge not only on which indexes to choose but also on when and how to build them with low cost. The Index Manager monitors previous queries to decide which indexes to build and when to build them; timing is based on an online randomized algorithm which considers (i) statistics on the cost of full scan ($C_{i_{fullscan}}$), (ii) statistics on the cost of building an index ($C_{i_{build}}$), and (iii) partition access frequency (AF_i), further explained in Sect. 4.2.

Update Monitor The main focus of Slalom is read-only and append workloads. Still, to provide query result consistency, the Update Monitor checks the input files for both appends and in-place updates at real time. Slalom enables append-like updates without disturbing query execution by dynamically adapting its auxiliary data structures. Specifically, Slalom creates a partition at the end of the file to accommodate the new data and builds binary caches, PMs, and indexes over them during the first post-update query. In-place updates require special care in terms of positional map and index maintenance because they can change the internal

file structure. Slalom reacts to in-place updates during the first post-update query by identifying the updated partitions, updating the positional map, and recreating the other corresponding structures. We discuss in detail how Slalom deals with updates in Sect. 4.3.

3.2 Implementation

We implement Slalom from scratch in C++. Slalom's query engine uses tuple-at-a-time execution based on the Volcano iterator model [26]. The rest of the components are implemented as modules of the query engine. Specifically, the *Partitioning* and *Indexing* Managers as well as the *Structure Refiner* attach to the Query Executor. Furthermore, the *Statistics Store* runs as a daemon, gathering the data and query statistics and persisting them in a catalog.

Slalom reduces raw data access cost by using vectorized parsers, binary caches, and positional maps (PMs). The CSV parser uses SIMD instructions; it consecutively scans a vector of 256 bytes from the input file and applies a mask over it using SIMD execution to identify delimiters. Slalom populates a PM for each CSV file accessed. To reduce memory footprint, the PM stores only delta distances for each tuple and field. Specifically, to denote the beginning of a tuple, the PM stores the offset from the preceding tuple. Furthermore, for each field within a tuple, the PM stores only the offset from the beginning of the tuple. The Partition Manager maintains a mapping between partitions and their corresponding PM portions.

Slalom populates binary caches at a partition granularity. When a query accesses an attribute for the first time, Slalom consults the positional map to identify the attribute's position and then caches the newly converted values. To improve insertion efficiency, Slalom stores the converted fields of each tuple as a group of columns. If Slalom opts to convert an additional field during a subsequent query, it appends the converted value to the current column group.

Slalom also populates secondary indexes at a partition granularity; for each attribute, the indexes store its position in the file and its position in the binary cache (when applicable). Slalom uses a cache friendly in-memory B⁺ tree implementation. It uses nodes of 256 bytes that are kept 60% full. To minimize the size of inner nodes and make them fit in a processor cache line, the keys in the nodes are stored as deltas. Furthermore, to minimize tree depth, the B⁺ tree stores all appearances of a single value in one record.

The Structure Refiner monitors the construction of all auxiliary structures and is responsible for memory management. Slalom works within a memory area of predefined size. The indexes, PMs, and caches are placed in the memory area. However, maintaining caches of the entire file and all possible indexes is infeasible. Thus, the Structure Refiner dynamically decides, on a partition basis, which structure to drop

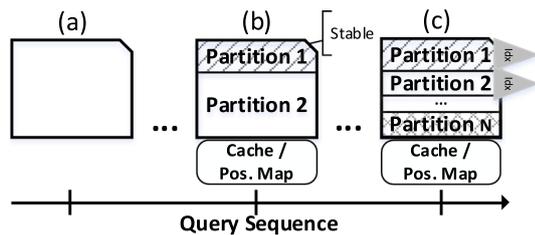


Fig. 3 Slalom execution

so Slalom can operate under limited resources (details in Sect. 4.2).

3.3 Query execution

Figure 3 presents an overview of a query sequence execution over a CSV file. During each query, Slalom analyzes its current state in combination with the workload statistics and updates its auxiliary structures. In the initial state (a), Slalom has no data or query workload information. The first query accesses the data file without any support from auxiliary structures; Slalom thus builds a PM, accesses the data requested, and places them in a cache. During each subsequent query, Slalom collects statistics regarding the data distribution of the accessed attributes and the average query selectivity to decide whether logical partitioning would benefit performance. If a partition has not reached its *stable* state (i.e., further splitting will not provide benefit), Slalom splits the partition into subsets as described in Sect. 4.1. In state (b), Slalom has already executed some queries and has built a binary cache and a PM on the accessed attributes. Slalom has decided to logically partition the file into two chunks, of which the first (partition 1) is declared to be in a *stable* state. Slalom checks stable partitions for the existence of indexes; if no index exists, Slalom uses the randomized algorithm described in Sect. 4.2 to decide whether to build one. In state (c), Slalom has executed more queries, and based on the query access pattern, it decided index partition 1. In this state, partition 2 of state (b) has been further split into multiple partitions of which partition 2 was declared *stable* and an index was built on it.

3.4 Extensibility of Slalom

To address the increasing data format heterogeneity, Slalom queries over a variety of data formats by adding the corresponding parsers and adjusting the online tuner partitioning algorithm.

The parser transforms all underlying data to a common representation, which is then passed to the query engine. In that way, Slalom supports multiple data formats by requiring a parser for each input data format (e.g., CSV, JSON, binary). Slalom uses as common representation binary tuples stored

in fixed length slots. Similarly, irrespective to data format, Slalom's binary cache has the same format.

For each new data format, the online tuner applies the same principled techniques of logical horizontal partitioning and indexing, however, must be adjusted slightly depending on the format. Specifically, for data formats that store records sequentially (e.g., CSV, binary, XML, and JSON) Slalom follows the same technique of partitioning and indexing by creating sequential logical partitions by keeping the first and last byte of each partition within the file. For data formats that store records in a PAX-like format [4] (e.g., parquet, SAM-BAM), the partitioning approach must make sure that the partitions complete full mini-pages. Slalom supports executing queries over CSV, binary, and XML files.

4 Continuous partition and index tuning

Slalom provides performance enhancements without requiring expensive full data indexing nor data file reorganization, all while adapting to workload changes. Slalom uses an online partitioning and indexing tuner to minimize the accessed data by (i) logically partitioning the raw dataset, and (ii) choosing appropriate indexing strategies over each partition. To enable online adaptivity, all decisions that the tuner makes must have minimal computational overhead. The tuner employs a Partition Manager which makes all decision considering the partitioning strategy, and an Index Manager which makes all decisions considering indexing. This section presents the design of the Partition and Index Managers as well as the mathematical models they are based on.

4.1 Raw data partitioning

The optimal access path may vary across different parts of a dataset. For example, a filtering predicate may be highly selective in one part of a file and thus benefit from index-based query evaluation, whereas another file part may be better accessed via a sequential scan. As such, any optimization applied on the entire file may be suboptimal for parts of the file. To this end, the Partition Manager of Slalom splits the original data into more manageable subsets; the minimum partition size is a physical disk page. The Partition Manager opts for horizontal logical partitioning as physical partitioning would require manipulating physical storage—a breaking point for many of the use cases that Slalom targets.

Why logical partitions Slalom uses logical partitioning to virtually break a file into more manageable chunks without physical restructuring. The goal of logical partitioning is twofold: (i) enable partition filtering, i.e., try to group relevant data values together so that they can be skipped for some queries, and (ii) allow for more fine-grained index tuning. The

efficiency of logical partitioning in terms of partition filtering depends mainly on data distribution and performs best with clustered or sorted data. Still, even in the worst case of uniformly distributed data, although few partitions will be skippable, the partitioning scheme facilitates fine-grained indexing. Instead of populating deep B⁺ tree that cover the entire dataset, the B⁺ tree of Slalom are smaller and target only “hot” subsets of the dataset. Thus, Slalom can operate under limited memory budget, has a minimal memory footprint, and provides rapid responses.

The Partition Manager performs partitioning as a by-product of query execution and chooses between two partitioning strategies depending on the cardinality of an attribute. For candidate key attributes, where all tuples have distinct values, the Partition Manager uses *query-based partitioning*, whereas for other value distributions, it uses *homogeneous partitioning*. Ideally, the Partition Manager aims to create partitions such that: (i) each partition contains uniformly distributed values, and (ii) partitions are pairwise disjoint (e.g., partition 1 has values 12, 1, 8 and partition 2 has values 19, 13, 30). Uniformly distributed values in a partition enable efficient index access for all values in a partition, and creating disjoint partitions improves partition skipping.

4.1.1 Homogenous partitioning

Homogeneous partitioning aims to create partitions with uniformly distributed values and maximize average selectivity within each partition. Increasing query selectivity over partitions implies that for some queries, some of the newly created partitions will contain a high percentage of the final results, whereas other partitions will contain fewer or zero results and will be skippable. Computing the optimal set of contiguous uniformly distributed partitions has exponential complexity and thus is prohibitive for online execution. Instead, to minimize the overhead of partitioning, the Partition Manager iteratively splits a partition into multiple equi-size partitions. In every iteration, the tuner decides on (i) when to stop splitting and (ii) into how many subsets to split a given partition.

The Partition Manager splits incrementally a partition until it reaches a *stable* state (i.e., a state where the tuner estimates no more gains can be achieved from further splitting). After each partition split, the tuner relies on two conditions to decide whether a partition has reached a stable state. The tuner considers whether (i) the variance of values in the new partition and the excess kurtosis [54] of the value distribution have become smaller than the variance and kurtosis in the parent partition, and (ii) the number of distinct values has decreased. Specifically, as variance and excess kurtosis decrease, outliers are removed from the partition and the data distribution of the partition in question becomes more uniform. As the number of distinct values per partition iteratively decreases, the probability of partition disjointness increases.

If any of these metrics increases or remains stable by partitioning, then the partition is declared stable. We use the combination of variance and excess kurtosis as a metric for uniformity, because their calculation has a constant complexity and can be performed in an incremental fashion during query execution. An alternative would be using a histogram or Chi-square estimators [54], but that would require building a histogram as well as an additional pass over the data.

Making partitioning decisions The number of subpartitions to which an existing partition is divided depends on the average selectivity of the past queries accessing the partition and the size of the partition in number of tuples. The goal of the tuner is to maximize selectivity in new partitions, thereby increasing the number of prospective skipped partitions. We assume that the rows of the partition that have been part of query results within the partition are randomly distributed. We model the partitioning problem as randomly choosing tuples from the partition with the goal to have at least 50% of the new partitions exhibit higher selectivity than the original partition. The intuition is that by decreasing selectivity in a subset of partitions will enhance partition skipping in the rest. The more results tuples in some partitions, the better candidates for skipping are the rest.

We model this problem with the hypergeometric distribution. Our goal is to choose m partitions by picking randomly n tuples, and we want each partition to contain at least k result tuples. The hypergeometric distribution is a discrete probability distribution that describes the probability of k random draws in n draws, without replacement. Thus, assuming that N represents all the tuples in the file, K represents the tuples appearing in the result, and $N - K$ all other tuples. The equation describing the CDF of hypergeometric distribution is the following:

$$P(X \geq k) \approx \sum_{i=k}^n \frac{\binom{K}{i} \binom{N-K}{n-i}}{\binom{N}{n}} \quad (1)$$

The calculation of the hypergeometric distribution requires the calculation of a factorial and has computational complexity $O(\log(\log(nM(n \log n))))$, where $M(n)$ is the complexity of multiplying two n -digit numbers [14]. Such a computational complexity is unacceptable for Slalom as this operation is executed for each query for the majority of partition numerical times and for large partition sizes.

Slalom approximates the hypergeometric distribution using the binomial distribution. Prior work shows that when $p \leq 0.1$ and $N \geq 60$ binomial is a good approximation of hypergeometric [47], and since the sizes of partitions are large in comparison with selectivity (small selectivity ≤ 0.1 and $N \geq 1000$), Slalom can exploit this observation.

$$P(X \geq k) = \sum_{i=k}^n \binom{n}{i} p^i (1-p)^{n-i} \quad (2)$$

The binomial distribution requires the calculation of the binomial coefficient $\binom{n}{i}$ which similarly to the hypergeometric distribution requires the calculation of factorial. To overcome this problem, we further approximate the binomial coefficient calculation by using the following equation [20]:

$$\binom{n}{k} = \frac{(n/k - 0.5)^k \cdot e^k}{\sqrt{2 \cdot \pi \cdot k}} \quad (3)$$

We combine Eqs. 2 and 3, we use $p = K/N$ and $n = N/m$, and we solve for m to get the equation that the Partition Manager uses to calculate the number of subpartitions created for every split:

$$m = \frac{N \cdot (sel + \log_b(1 - sel))}{\log_b \frac{\sqrt{2 \cdot \pi \cdot sel \cdot N}}{e}} \quad (4)$$

where $b = \frac{e}{sel \cdot (1 - sel)}$

The tuner chooses this set of partitions with the minimal overhead and number of iterations. The number of distinct values is calculated during the next query after each partition split, whereas the variance and the kurtosis are calculated incrementally; thus, the partitioning algorithm creates negligible overheads. To achieve that, Slalom uses a set of one-pass algorithms for calculating common statistics [48].

4.1.2 Query-based partitioning

Query-based partitioning targets candidate keys or attributes that are *implicitly clustered* (e.g., increasing time stamps). For such attributes, homogeneous partitioning will lead to increasingly small partitions as the number of distinct values and variance will be constantly decreasing with smaller partitions. Thus, the tuner decides upon a static number of partitions to split the file. Specifically, the number of partitions is decided based on the selectivity of the first range query using the same mechanism as in homogeneous partitioning. If the partition size is smaller than the physical disk page size, the tuner creates a partition per disk page. By choosing its partitioning approach based on the data distribution, Slalom improves the probability of data skipping and enables fine-grained indexing.

4.2 Adaptive indexing in Slalom

The tuner of Slalom employs the Index Manager to couple logical partitions with appropriate indexes and thus decrease

the amount of accessed data. The Index Manager uses *value-existence* and *value-position* indexes; it takes advantage of the capabilities of each category in order to reduce execution overhead and memory footprint. To achieve these goals, the Index Manager enables each partition to have multiple value-existence and value-position indexes.

Value-existence indexes Value-existence indexes are the basis of partition skipping for Slalom; once a partition has been set as stable, the Index Manager builds a value-existence index over it. Value-existence indexes allow Slalom to avoid accessing some partitions. The Index Manager uses Bloom filters, bitmaps, and zonemaps (min–max values) as value-existence indexes. Specifically, the Index Manager uses bitmaps only when indexing Boolean attributes, because they require a larger memory budget than Bloom Filters for other data types. The Index Manager also uses zonemaps on all partitions because they have small memory overhead and provide sufficient information for value existence on partitions with small value variation. For all other data types, the Index Manager favors Bloom filters because of their high performance and small memory footprint. Specifically, the memory footprint of a Bloom filter has a constant factor, yet it also depends on the number of distinct values it will store and the required false positive probability. To overcome the inherent false positives that characterize Bloom filters, the Index Manager adjusts the Bloom filter’s precision by calculating the number of distinct values to be indexed and the optimal number of bytes required to model them [13].

Value-position indexes The Index Manager builds a value-position index (B⁺ tree) over a partition to offer fine-grained access to tuples. As value-position indexes are more expensive to construct compared to value-existence indexes, both in terms of memory and time, it is crucial for the index to pay off the building costs in future query performance. The usefulness and performance of an index depend highly on the type and selectivity of queries and the distribution of values in the dataset. Thus, for workloads of shifting locality, the core challenge is deciding *when* to build an index.

When to build a value-position index The Index Manager builds a value-position index over a partition if it estimates that there will be enough subsequent queries accessing that partition to pay off the investment (in execution time). As the tuner is unaware of the future workload trends, decisions for building indexes are based on the past query access patterns. To make these decisions, the Index Manager uses an online randomized algorithm which considers the cost of indexing the partition ($C_{i_{build}}$), the cost of full partition scan ($C_{i_{fullscan}}$), and the access frequency on the partition (AF_i). These values depend on the data type and the size of the partition, so they are updated accordingly in case of a partition split or an append to the file. The tuner stores the average cost of an access to a file tuple as well as the average cost of an inser-

tion to every index for all data types and uses these metrics to calculate the cost of accessing and building an index over a partition. In addition, the tuner calculates the cost of an index scan ($C_{i_{indexscan}}$) based on the cost of a full partition scan and the average selectivity. For each future access to the partition, the Index Manager uses these statistics to generate online a probability estimate calculating whether the index will reduce execution time for the rest of the workload. Given this probability, the Index Manager decides whether to build the index.

The Index Manager calculates the index building probability using a randomized algorithm based on the randomized solution of the snoopy caching problem [39]. In the snoopy caching problem, two or more caches share the same memory space which is partitioned into blocks. Each cache writes and reads from the same memory space. When a cache writes to a block, caches that share the block spend 1 bus cycle to get updated. These caches can invalidate the block to avoid the cost of updating. When a cache decides to invalidate a block which ends up required shortly after, there is a penalty of p cycles. The optimization problem lies in finding when a cache should invalidate and when to update the block. The solution to the index building problem in this work involves a similar decision. The indexing mechanism of the tuner of Slalom decides whether to pay an additional cost per query (“updating a block”) or invest in building an index, hoping that the investment will be covered by future requests (“invalidating a block”). Specifically, in cases where the cost of using an index is negligible compared to the cost of full data scan, deciding on index construction can be directly mapped to the snoopy caching problem.

The performance measure of randomized algorithms is the *competitive ratio (CR)*: the ratio between the expected cost incurred when the online algorithm is used and that of an optimal offline algorithm that we assume has full knowledge of the future. When index access cost is negligible, the randomized algorithm of the tuner guarantees optimal CR ($\frac{e}{e-1}$). The tuner uses a randomized algorithm in order to avoid the high complexity of what-if analysis [59] and to improve the competitive ratio offered by the deterministic solutions [15].

Cost model Assume query workload W . At a given query q of the workload, a partition is in one of the two states: it either has an index or it does not. A state is characterized by the pair (C_{build}, C_{use}) where C_{build} is the cost to enter the state (e.g., build the index) and C_{use} the cost to use the state (e.g., use the index). The initial state is the state with no index (i.e., full scan) $(C_{build,fs}, C_{use,fs})$ where $C_{build,fs} = 0$. In the second state $(C_{build,idx}, C_{use,idx})$, the system has an index. We assume that the relation between the costs for the two states is $C_{build,idx} > C_{build,fs}$ and $C_{use,idx} < C_{use,fs}$ and $C_{build,idx} > C_{use,fs}$.

Given a partition i , the index building cost over that partition ($C_{i_{build}}$), the full partition scan cost ($C_{i_{fullscan}}$), the index partition scan cost ($C_{i_{indexscan}}$) and a sequence of queries $Q : [q_1, \dots, q_T]$ access the partition. Assume that q_T is the last query that accesses the partition (*and is not known*). At the arrival time of $q_k, k < T$, we want to decide whether the Index Manager should build the index or perform full scan over the partition to answer the query.

To make the decision, we need a probability estimate p_i for building the index at moment i based on the costs of building the index or not. In order to calculate p_i , we initially define the overall expected execution cost of the randomized algorithm that depends on the probability p_i . The expected cost E comprises three parts:

- (i) the cost of using the index, which corresponds to the case where the index has already been built.
- (ii) the cost of queries doing full partition scan, which corresponds to the case for which the index has not been built.
- (iii) the cost of building the index, which corresponds to the case where the building of the index will take place at time i . Index construction takes place as a by-product of query execution and includes the cost of the current query.

$$E = \sum_{i=1}^T \left(\sum_{j=1}^{i-1} p_j \cdot C_{use,idx} + \left(1 - \sum_{j=1}^{i-1} p_j \right) \cdot \left(p_i \cdot C_{build,idx} + (1 - p_i) \cdot C_{use,fs} \right) \right)$$

Knowing the expected cost, we minimize and we solve for p_i :¹

$$p_i = \frac{C_{use,fs} - C_{use,idx}}{C_{build,idx} - C_{use,fs}} \cdot (T - i) - \left(1 - \sum_{j=1}^{i-1} p_j \right) \quad (5)$$

Based on our model, performing a full scan over the complete data file should be always cheaper than an index access and the amortized extra cost of building the index (over T queries).

Eviction policy The tuner works within a predefined memory budget to minimize memory overhead. If the memory budget is fully consumed and the Index Manager attempts

¹ Details on how this formula is derived are found in “Appendix.”

to build a new index, then it defers index construction for the next query and searches indexes to drop to make the necessary space available. The Index Manager keeps all value-existence indexes once built, because their size is minimal and they are the basis of partition skipping. Furthermore, the Index Manager prioritizes binary caches over indexes, because (i) using a cache improves the performance of all queries accessing a partition, and (ii) accessing the raw data file is typically more expensive than rebuilding an index for large partitions. Deciding which indexes from which partitions to drop is based on index size ($Size_{index_i}$), number of queries since last access (LA_i), and average selectivity (sel_i) in a partition. To compute the set of indexes to drop, the Index Manager uses a greedy algorithm which gathers the least accessed indexes with cumulative size ($\sum_i Size_{index_i}$) equal to the size of the new index. Specifically, to discover the least accessed indexes, the Index Manager keeps a bitmap of accesses for each partition. During a query predicate evaluation on a partition and depending on whether the current query touches the partition, the Index Manager shifts the partition's bitmap to the left and appends a bit to it: 1 (yes) or 0 (no). When calculating the candidate indexes to drop, the Index Manager uses SIMD instructions to evaluate the set of least accessed partitions. Specifically, each bitmap is an 8-byte unsigned integer which stores the past 64 queries. In a 256-byte wide CPU register, the Index Manager uses a bitmask operation to check the occupancy of 32 partitions simultaneously. When all indexes are used with the same frequency, the tuner uses the average selectivity of queries on each partition as a tie-breaker condition. The less selective the queries are, the smaller the gap between index and full scan performance; therefore, the Index Manager victimizes partitions touched by non-selective queries.

4.3 Handling file updates

Slalom supports both append-like and in-place updates directly over the raw data file and ensures consistent results. In order to achieve efficient data access and correct results despite updates, Slalom continuously monitors the queried files for any write operation and stores summaries of the queried files representing their current state. If a file is updated, Slalom compares its existing summary, with the stored state, identifies the changes, and updates any dependent data structures.

In this section, we describe in detail how Slalom: (i) monitors its input files for updates at real time, (ii) calculates and stores a summary of the most recent consistent state for reference, (iii) identifies the updated file subsets, and (iv) updates its internal data structures.

4.3.1 Monitoring Files

In order to recognize whether an input file has been updated by another application (e.g., vim), Slalom uses OS support (i.e., inotify [43]). Specifically, Slalom initializes a watchdog, over the queried file, which is triggered when the file is written upon and adds a log entry into a queue. This queue contains all updates that have not been addressed by Slalom yet. Slalom checks the queue for new updates both at the beginning of every query as well as during execution. During a running query, Slalom checks for any updates that happened in data that has been already scanned. If such an update has taken place, Slalom re-executes the query as results might be invalid if the records processed come from different file versions.

4.3.2 Calculating and Storing State

In order to be able to discover the updated rows in the file and the type of update (append or in-place), Slalom exploits its logical partitioning scheme. For each partition, Slalom stores a checksum encoding the contents within that partition and the starting and ending positions of the partition in the file. This information is collected during the first query accessing a partition. The collected information summarizes the size as well as the content of the partition and thereby is sufficient to identify the existence of an update. As the checksum calculation is part of the critical path of query execution, it increases the query runtime. To alleviate this cost, Slalom exploits specialized hardware that offers high throughput in checksum calculation. Furthermore, the performance and accuracy of checksum algorithms depend highly on the size of data they summarize; thus, Slalom varies the checksum algorithm depending on partition size. Currently, Slalom supports two checksum algorithms: (i) MD5 and (ii) CRC; these algorithms are widely used in a variety of applications based on their reliability and performance.

MD5 algorithm MD5 [58] is a cryptographic hash function and widely used data integrity verification checksum [65]. Given input of arbitrary size, MD5 algorithm produces a 128-bit output, which is usually represented in 32 hexadecimal digits. MD5 uses four nonlinear functions, and it can deal with data of arbitrary length. MD5 serves as a good candidate for detecting file updates; however, its calculations on a CPU are expensive. Thus, we design a parallelization scheme for MD5. MD5 is an irreversible transformation transforming a set of data of any length into a hash value of 128-bit length. MD5 is a consecutive processing method as the original algorithm processes the input data incrementally in 512-bit groups and combines them with the result coming from the processing of prior groups. In order to parallelize the computation of MD5, we compute in parallel different portions

of the checksum. We divide the input data into small blocks of equal size. Subsequently, we perform the standard MD5 algorithm on each data block, in parallel, and we store the calculated checksums. Finally, the resulting checksums are combined until the result is 128-bit long. The checksum computed by this approach is not identical to the standard MD5 checksum and however has equal encryption strength [30]. As the algorithm is inherently suitable for multi-threading, and to further improve the performance of MD5 checksum calculation, we implemented the parallel MD5 over NVIDIA CUDA and calculate checksums over NVIDIA GPUs.

CRC Cyclic redundancy codes are used to mostly detect errors in network packets [55]. As this operation is latency sensitive, modern processors have added CPU instructions, `_mm_crc32_u64`, for calculating 32-bit CRC code to its SSE4.2 instruction set. To obtain m -bit CRC code, the n -bit input data are first appended with m zeros. Then, it is XORed with a polynomial divisor of the size of $(n + 1)$ bit from left to right. The last m bits are the final resulting code.

Typically a n -bit CRC applied to a data block of arbitrary length will detect any single error burst not longer than n bits and will detect a fraction $\frac{1}{(1-2^{-n})}$ of all longer error bursts. As partitions used by Slalom can be of arbitrary size, Slalom calculates the 32-bit CRC value for each 1024-byte block in the partition and then adds up all computed values to give the final verification code. This code has the same detection ability, namely detecting changes no longer than 4 bytes and almost all longer changes.

4.3.3 Recognizing Update Type and Updating Data Structures

In order to provide efficient data access, Slalom builds a set of data structures which are built based on the existing state of the queried file. Updates may change that state thus making the prior investments obsolete. Specifically, indexes and PMs are sensitive to the specific location of attributes and number of tuples within the file. Similarly, caches and Bloom filters become obsolete with any change in a partition. To overcome this issue, Slalom updates its data structures accordingly depending on the update type.

To identify the type of update, Slalom compares the current state of each partition with the stored one. Thus, Slalom checks whether the partition beginning and ending character has changed or if the checksum has changed. If the state of each partition matches with the existing one, then the update type is an append. Otherwise, it is an in-place update.

Append-like updates Slalom supports updates in an append-like scenario without disturbing query execution and by dynamically extending auxiliary data structures. In append-like scenarios, Slalom creates a new partition at the end of the file to accommodate the new data. Depending on the

partitioning approach, Slalom either accumulates updates to create partitions of equal size (i.e., query-based partitioning) or dynamically repartitions the fresh data. Once Slalom has organized the new data in partitions, it treats them similarly to a first time input. Thus, during the first query after an update, Slalom builds binary caches and positional maps over the new data. When the new partitions are declared *stable*, Slalom builds indexes on top of them.

In-place updates In-place updates correspond to random changes in the file by another application, such as updating values of specific fields or adding additional rows in the middle of the file. In-place updates are more challenging, especially when considering the case of the positional map and indexes. A change in a position of an attribute in the data file might require significant reorganization in all generated data structures.

Updating positional maps To update the positional map for a modified partition, Slalom scans character by character each field to narrow down the updated parts. Once the updated section has been identified, Slalom stores the difference in byte offsets between the old and new fields into a *delta list*. All new changes are appended to the list, and any possible changes in previous offset differences are being integrated as well. The delta list adds additional computational overhead when using the positional map as for every access Slalom must access the delta list to check whether the position has been altered by an update. As the delta list is growing, the complexity of position computation is growing as well. Thus to reduce the query cost, the delta list is incorporated into the original positional map every ten updates. Specifically, to incorporate the delta list into the positional map, Slalom scans over the delta list and adds the offsets to the existing indexes in the positional map. This way, it does not have to completely reconstruct the positional map while reducing the delta list.

Updating caches and indexes In order to keep minimal memory footprint, Slalom does not store a replica of the original file to be able retrieve old values for each updated field. Hence, Slalom is unable to update indexes and caches. Rather, it invalidates and re-builds them.

5 Experimental Evaluation

In this section, we present an analysis of Slalom. We analyze its partitioning and indexing algorithm and compare it against state-of-the-art systems over both synthetic and real-life workloads.

Methodology We compare Slalom against DBMS-X, a commercial state-of-the-art in-memory DBMS that stores records in a row-oriented manner and the open-source DBMS PostgreSQL (version 9.3). We use DBMS-X and PostgreSQL

with two different configurations: (i) fully loaded tables and (ii) fully loaded, indexed tables. We also compare Slalom with the in situ DBMS PostgresRaw [5]. PostgresRaw is an implementation of NoDB [5] over PostgreSQL; PostgresRaw avoids data loading and executes queries by performing full scans over CSV files. In addition, PostgresRaw builds positional maps on-the-fly to reduce parsing and tokenization costs. Besides positional maps, PostgresRaw uses caching structures to hold previously accessed data in a binary format. Furthermore, to compare Slalom with other adaptive indexing techniques we integrate into Slalom two variations of database cracking: (i) standard cracking [32] and (ii) the MDDIR variant of stochastic cracking [28]. We chose MDDIR as it showed the best overall performance in [60]. We integrated the cracking techniques by disabling the Slalom tuner and setting cracking as the sole access path. Thus, Slalom and cracking use the same execution engine and have the same data access overheads.

Slalom's query executor pushes predicate evaluation down to the access path operators for early tuple filtering, and results are pipelined to the other operators of a query (e.g., joins). Thus, in our analysis, we focus on scan intensive queries. We use select–project–aggregate queries to minimize the number of tuples returned and avoid any overhead from the result tuple output that might affect the measured times. Unless otherwise stated, the queries are of the following template ($OP : \{<, >, =\}$):

```
SELECT agg(A), agg(B), . . . , agg(N) FROM R
WHERE A OP X (AND A OP Y)
```

Experimental Setup The experiments are conducted in a Sandy Bridge server with a dual socket Intel(R) Xeon(R) CPU E5-2660 (8 cores per socket @ 2.20 Ghz), equipped with 64 KB L1 cache and 256 KB L2 cache per core, 20 MB L3 cache shared, and 128GB RAM running Red Hat Enterprise Linux 6.5 (Santiago—64 bit) with kernel version 2.6.32. The server is equipped with a RAID-0 of 7 250 GB 7500 RPM SATA disks.

5.1 Adapting to workload shifts

Slalom adapts efficiently to workload shifts despite (i) changes in data distribution, (ii) changes in query selectivity, and (iii) changes in query locality—both vertical (i.e., different attributes) and horizontal (i.e., different records). We demonstrate the adaptivity experimentally by executing a dynamic workload with varying selectivity and access patterns over a synthetic dataset.

Methodology To emulate the worst possible scenario for Slalom, we use a relation of 640 million tuples (59GB), where each tuple comprises 25 unsigned integer attributes with uniformly distributed values ranging from 0 to 1000.

Slalom is unable to find a value clustering in the file because all values are uniformly distributed; thus, Slalom applies homogeneous partitioning. Slalom, cracking, and PostgresRaw operate over the CSV data representation, whereas PostgreSQL and DBMS-X load the raw data prior to querying. In this experiment, we limit the index memory budget for Slalom to 5 GB and the cache budget to 10 GB. All other systems are free to use all available memory. Specifically, for this experiment DBMS-X required 98 GB of RAM to load and fully build the index.

We execute a sequence of 1000 point and range select–project–aggregation queries following the template from Sect. 5. The selection value is randomly selected from the domain of the predicate attribute. Point query selectivity is 0.1%, and range query selectivity varies from 0.5 to 5%. To emulate workload shifts and examine system adaptivity, in every 100 queries, queries 1–30 and 61–100 use a predicate on the first attribute of the relation and queries 31–60 use a predicate on the second attribute.

The indexed variations of PostgreSQL and DBMS-X build a clustered index only on the first attribute. It is possible to build indexes on more columns for PostgreSQL and DBMS-X; however, it requires additional resources and would increase data-to-query time. In addition, choosing which attributes to index requires a priori knowledge of the query workload, which is unavailable in the dynamic scenarios that Slalom considers. Indicatively, building an secondary index on a column for PostgreSQL for our experiment takes ~25 minutes. Thus, by the time PostgreSQL finishes indexing, Slalom will have finished executing the workload (Fig. 6).

Slalom Convergence Figure 4 shows the response time of each query of the workload for the different system configurations. For clarity, we present the results for the first 100 queries. To emulate the state of DBMS systems immediately after loading, all systems run from a hot state where data are resting in the OS caches. Figure 4 plots only query execution time and does not show data loading or index building for PostgreSQL and DBMS-X.

The runtime for the first query of Slalom is 20× slower than its average query time, because during that query it builds a positional map and a binary cache. In subsequent queries (queries 2–7), Slalom iteratively partitions the dataset and builds B⁺ tree. After the initial set of queries (queries 1–6), Slalom has comparable performance to the one of PostgresRaw over fully indexed data. During the third query, multiple partitions stabilize simultaneously, and thus, Slalom builds many B⁺ tree and Bloom filter indexes, adding considerable overhead. When Slalom converges to its final state, its performance is comparable to indexed DBMS-X. When the queried attribute changes (query 31), Slalom starts partitioning and building indexes on the new attribute. After query

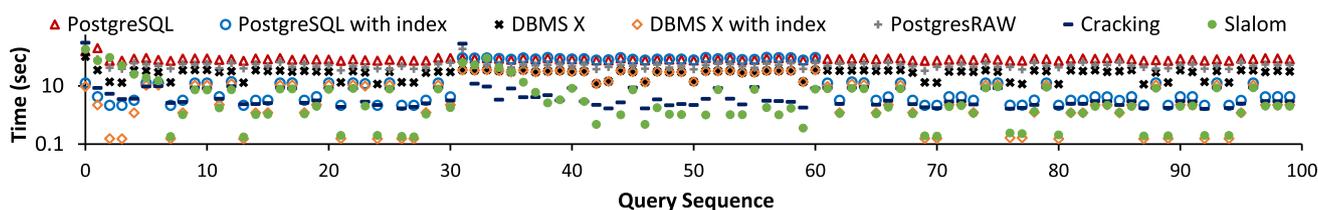


Fig. 4 Sequence of 100 queries. Slalom dynamically refines its indexes to reach the performance of an index over loaded data

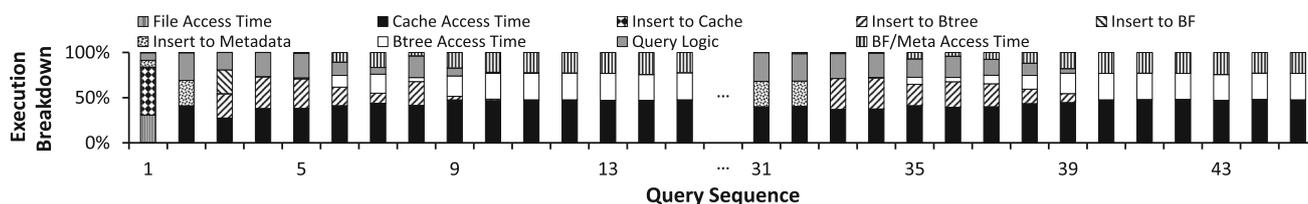


Fig. 5 A breakdown of the operations taking place for Slalom during the execution of a subset of the 100 point query sequence

60, when the workload filters data based on the first attribute again (for which the partitioning is already stable), Slalom re-uses the pre-existing indexes.

PostgreSQL with no indexes demonstrates a stable execution time as it has to scan all data pages of the loaded database regardless of the result size. Due to the queries being very selective, when an index is available for PostgreSQL, the response times are $\sim 9\times$ lower when queries touch the indexed attribute. DBMS-X keeps all data in memory and uses memory-friendly data structures, so it performs on average $3\times$ better than PostgreSQL. The difference in performance varies with query selectivity. In highly selective queries, DBMS-X is more efficient in data access, whereas for less selective queries the performance gap is smaller. Furthermore, for very selective queries, indexed DBMS-X is more efficient than Slalom as its single B⁺ tree traverses very few results nodes.

During query 1, PostgresRaw builds auxiliary structures (cache, positional map) and takes $3\times$ more time (180s) than its average query runtime. PostgresRaw becomes faster than the unindexed PostgreSQL variation as its scan operators use vector-based (SIMD) instructions and exploit compact caching structures.

Similarly, during query 1, cracking builds a binary cache and populates the cracker column it uses for incremental indexing. The runtime of its first query is $4\times$ slower than the average query time for PostgreSQL without indexes. When it touches a different attribute (query 31), it also populates a cracker column for the second attribute. Despite the high initialization cost, cracking converges efficiently and reaches its final response time after the fourth query. The randomness in the workload benefits cracking as it splits the domain into increasingly smaller pieces. After converging, cracking performance is comparable to the PostgreSQL with index. Slalom requires more queries to converge than

cracking. However, after it converges, Slalom is $\sim 2\times$ faster than cracking. This difference stems from cracking execution overheads. Cracking sorts the resulting tuples based on their memory location and enforces sequential memory access. This sorting operation adds an overhead, especially for less selective queries.

Execution breakdown Slalom aims to build efficient access paths with minimal overhead. Figure 5 presents the breakdown of query execution for the same experiment as before. For clarity, we present only queries Q1–15 and Q31–45 as Q16–30 show the same pattern as Q11–15. Queries Q1–15 have a predicate on the first attribute, and queries Q31–45 have a predicate on the second attribute.

During the first query, Slalom scans through the original file and creates the cache. During Q2 and Q3, Slalom is actively partitioning the file and collects data statistics (i.e., distinct value counts) per partition; Slalom bases the further partitioning and indexing decisions on these statistics. Statistics gathering cost is represented in Fig. 5 as “Insert to Metadata.” During queries Q2 and Q3, as the partitioning scheme stabilizes, Slalom builds Bloom filters and B⁺ trees. Q3 is the last query executed using a full partition scan, and since it also incurs the cost of index construction, there is a local peak in execution time. During Q4 through Q8, Slalom increasingly improves performance by building new indexes. After Q31, the queries use the second attribute of the relation in the predicate, and thus, Slalom repeats the process of partitioning and index construction. In total, even after workload shifts, Slalom converges into using index-based access paths over converted binary data.

Full workload: from raw data to results Figure 6 presents the full workload of 1000 queries, this time starting with cold OS caches and no loaded data to include the cost of the first access to raw data files for all systems. We plot the aggregate execution time for all approaches described earlier, including

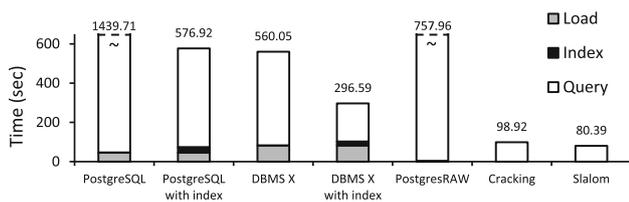


Fig. 6 Sequence of 1000 queries. Slalom does not incur loading cost and dynamically builds indexes

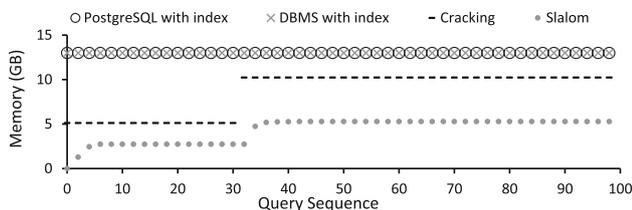


Fig. 7 Memory consumption of Slalom vs. a single fully built B⁺ tree for PostgreSQL and DBMS-X. Slalom uses less memory because its indexes only target specific areas of a raw file

the loading and indexing costs for PostgreSQL and DBMS-X.

PostgresRaw, Slalom, and cracking incur no loading and indexing cost and start answering queries before the other DBMS load data and before the indexed approaches finish index building. Unindexed PostgreSQL incurs data loading cost as well as a total query aggregate greater than PostgresRaw. Indexed PostgreSQL incurs both indexing and data loading cost, and due to some queries touching a non-indexed attribute, its aggregate query time is greater than the one of Slalom. Unindexed DBMS-X incurs loading cost; however, thanks to its main-memory friendly data structures and execution engine, it is faster than the disk-based engine of PostgreSQL.

After adaptively building the necessary indexes, Slalom has comparable performance with a conventional DBMS which uses indexes. cracking converges quickly and adapts to the workload efficiently. However, creating the cracker columns incurs a significant cost. Overall, cracking and Slalom offer comparable raw-data-to-results response time for this workload, while Slalom requires 0.5× memory. We compare in detail cracking and Slalom in Sect. 5.3.

Memory consumption Figure 7 plots the memory consumption of (i) the fully built indexes used for DBMS-X and PostgreSQL, (ii) the cracker columns for cracking, and (iii) the indexes of Slalom. Figure 7 excludes the size of the caches used by Slalom and cracking or the space required by DBMS-X after loading. The traditional DBMS require significantly more space for their indexes. Orthogonally to the index memory budget, DBMS-X required 98 GB of memory in total, whereas the cache of Slalom required 9.7GB. Cracking builds its cracker columns immediately

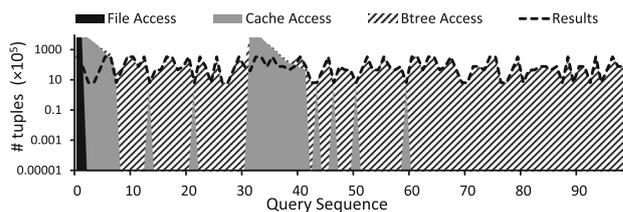


Fig. 8 Number of accessed tuples using file, cache, or B⁺ tree corresponding to the 100 queries of synthetic workload

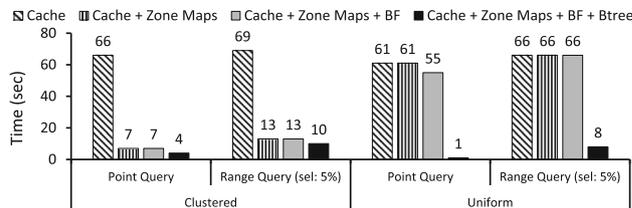


Fig. 9 Effect of different indexes on point and range queries over uniform and clustered datasets

when accessing a new attribute. The cracker column requires storing the original column values as well as pointers to the data; thus, it has a large memory footprint even for low value cardinality. Regarding the indexes of Slalom, when the focus shifts to another filtering attribute (Q31), Slalom increases its memory consumption, as during Q31–34 it creates logical partitions and builds Bloom filters and B⁺ tree indexes on the newly accessed attribute. By building and keeping only the necessary indexes for a query sequence, Slalom strikes a balance between query performance and memory utilization.

Minimizing data access The performance gains of Slalom are a combination of data skipping based on partitioning, value-existence indexes, and value-position indexes, all of which minimize the number of tuples Slalom has to access. Figure 8 presents the number of tuples that Slalom accesses for each query in this experiment. We observe that as the partitioning and indexing schemes of Slalom converge, the number of excess tuples accessed is reduced. Since the attribute participating in the filtering predicate of queries Q31–60 has been cached, Slalom accesses the raw data file only during the first query. Slalom serves the rest of the queries utilizing only the binary cache and indexes. For the majority of queries, Slalom responds using an index scan. However, there are queries where it responds using a combination of partition scan and index scan.

Figure 9 presents how the minimized data access translates to reduced response time and the efficiency of data skipping and indexing for different data distribution and different query types. Specifically, it presents the effect of zonemaps, Bloom filters, and B⁺ trees on query performance for point queries and range queries with 5% selectivity over uniform and clustered datasets. The clustered dataset contains mutually disjointed partitions (i.e., subsets of the file



Fig. 10 Slalom performance using different memory budgets

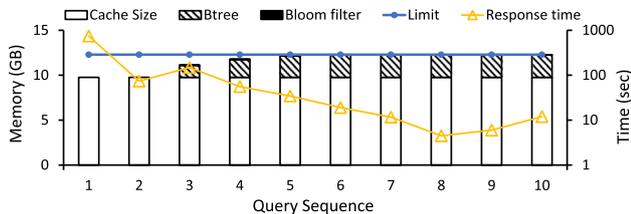


Fig. 11 Slalom memory allocation (12 GB memory budget)

contain values which do not appear in the rest of the file). The workload used is the same used for Fig. 4. Zonemaps are used for both range and point queries and are most effective when used over clustered data. Specifically, they offer a $\sim 9\times$ better performance than full cache scan. Bloom filters are useful only for point queries. As the datasets have values in the domain $[1, 1000]$, point queries have low selectivity making Bloom filters ineffective. Finally, B^+ trees improve performance for both range and point queries. The effect of B^+ tree is seen mostly for uniform data where partition skipping is less effective. Slalom stores all indexes in-memory; thus, by skipping a partition Slalom avoids full access of the partition and reduces memory access or disk I/O if the partition is cached or not, respectively.

Summary We compare Slalom against (i) a state-of-the-art in situ querying approach, (ii) a state-of-the-art adaptive indexing technique, (iii) a traditional DBMS, and (iv) a state-of-the-art in-memory DBMS. Slalom gracefully adapts to workload shifts using an adaptive algorithm with negligible execution overhead. Slalom offers performance comparable with a DBMS which uses indexes, while also being more conservative in memory space utilization.

5.2 Working under memory constraints

As described in Sect. 4.2, Slalom efficiently uses the available memory budget to keep the most beneficial auxiliary structures. We show this experimentally by executing the same workload under various memory utilization constraints. We run the 20 first queries—a mix of point and range queries. We consider three memory budget configurations with 10 GB, 12 GB, and 14 GB of available memory, respectively. The budget includes both indexes and caches.

Figure 10 presents the query execution times for the workload given the three different memory budgets. The three memory configurations build a binary cache and create the same logical partitioning. Slalom requires 13.5 GB in total for this experiment; given an 14 GB memory budget, it can build all necessary indexes, leading to the best performance for the workload. For the 10 GB and 12 GB memory budgets, there is insufficient space to build all necessary indexes; thus, these configurations experience a performance drop. We observe that configurations with 10 GB and 12 GB memory budgets outperform the configuration with 14 GB of memory budget for individual queries (i.e., Q3 and Q5). The reason is that the memory-limited configurations build fewer B^+ trees during these queries than the configuration with 14 GB of available memory. However, future queries benefit from additional B^+ trees, amortizing the extra overhead over a sequence of queries.

Figure 11 presents the breakdown of memory allocation for the same query sequence when Slalom is given a 12 GB memory budget. We consider the space required for storing caches, B^+ trees, and Bloom filters. The footprint of the statistics and metadata Slalom collects for the cost model and zonemaps is negligible; thus, we exclude them from the breakdown. Slalom initially builds the binary cache and logically partitions the data until some partitions become stable (Q1, Q2). At queries Q3, Q4, and Q5, Slalom starts building B^+ trees, and it converges to a stable state at query Q7 where all required indexes are built. Thus, from Q7–Q10 Slalom stabilizes performance. Overall, this experiment shows that Slalom can operate under limited memory budget gracefully managing the available resources to improve query execution performance.

5.3 Adaptivity efficiency

Slalom adapts to query workloads as efficiently as state-of-the-art adaptive indexing techniques while working with less memory. Furthermore, it exploits any potential data clustering to further improve its performance. We demonstrate this by executing a variety of workloads. We use datasets of 480M tuples (55 GB on disk); each tuple comprises 25 unsigned integer attributes whose values belong to the domain $[1, 10,000]$. Queries in all workloads have equal selectivity to alleviate the noise from data access; all queries have 0.1% selectivity, i.e., select ten consecutive values.

Methodology Motivated by related work [60], we compare Slalom against cracking and stochastic cracking in three cases.

Random workload over Uniform dataset We execute a sequence of range queries which access random ranges throughout the domain to emulate the best-case scenario for cracking. As subsequent queries filter on random values and

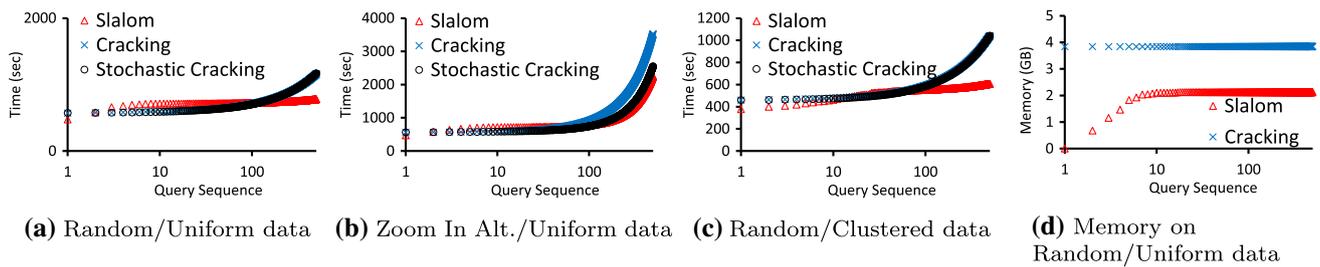


Fig. 12 Comparing cracking techniques with Slalom

the data are uniformly distributed in the file, cracking converges and minimizes data access.

“Zoom In Alternate” over Uniform dataset To emulate the effect of patterned accesses, we execute a sequence of queries that access either part of the domain in alternate, i.e., first query: [1, 10], second query: [9991, 10,000], third query: [11, 20], etc. This access pattern is one of the scenarios where the original cracking algorithm underperforms [28]. Splits are only query driven, and every query splits data into a small piece and the rest of the file. Thus, the improvements in performance with subsequent queries are minimal. Stochastic cracking alleviates the effect of patterned accesses by splitting in more pieces apart from the ones based on queries.

Random workload over Clustered dataset This setup examines how adaptive indexing techniques perform on datasets where certain data values are clustered together, for example data clustered on time stamp or sorted data. The clustered dataset we use in the experiment contains mutually disjoint partitions, i.e., subsets of the file contain specific values which appear solely in those locations and do not appear in the rest of the file.

Figure 12a demonstrates the cumulative execution time for cracking, stochastic cracking, and Slalom for the random workload over uniform data. All approaches start from a cold state, thus during the first query they parse the raw data file and build a binary cache. Stochastic cracking and cracking incur an additional cost of cracker column initialization during the first query, but reduce execution time with every subsequent query. During the first three queries, Slalom creates its partitions; during the following six queries, Slalom builds the required indexes and finally converges to a stable state at query 10. Due to its fine-grained indexing and local memory accesses, Slalom provides $\sim 8\times$ lower response time than cracking and their cumulative execution time is equalized during query 113. Furthermore, Fig. 12d demonstrates the memory consumption of the cracking approaches and Slalom for the same experiment. The cracking approaches have the same memory footprint; they both duplicate the full indexed column along with pointers to the original data. On the other hand, the cache-conscious

B^+ trees of Slalom stores only the distinct values along with the positions of each value, thus reducing the memory footprint. In addition, Slalom allocates space for its indexes gradually, offering efficient query execution even with limited resources.

Figure 12b shows the cumulative execution time for cracking, stochastic cracking, and Slalom for the “Zoom In Alternate” workload over uniform data. Cracking needs more queries to converge to its final state as it is cracking only based on query-driven values. Stochastic cracking converges faster because it cracks based on more values except the ones found in queries. Slalom uses a combination of data and query-driven optimizations. Slalom requires an increased investment during the initial queries to create its partitioning scheme and index the partitions, but ends up providing $7\times$ lower response time, and equalizes cumulative execution time with cracking at query 53 and stochastic cracking at query 128.

Figure 12c presents the cumulative execution time of cracking, stochastic cracking, and Slalom for the random workload over implicitly clustered data. In this situation, Slalom exploits the clustering of the underlying data early on (from the second query) and skips the majority of data. For the accessed partitions, Slalom builds indexes to further reduce access time. Similarly to Fig. 12a, the cracking approaches crack only based on the queries and are agnostic to the physical organization of the dataset.

Summary Slalom converges comparably to the best cracking variation when querying uniform data over both random and “Zoom In Alternate” workloads. Furthermore, when Slalom operates over clustered data, it exploits the physical data organization and provides minimal data-to-query time. Finally, as Slalom builds indexes gradually and judiciously, it requires less memory than the cracking approaches, and it can operate under a strict memory budget.

5.4 Slalom over real data

In this experiment, we demonstrate how Slalom serves a real-life workload. We use a smart home dataset (SHD) taken from an electricity monitoring company. The dataset con-

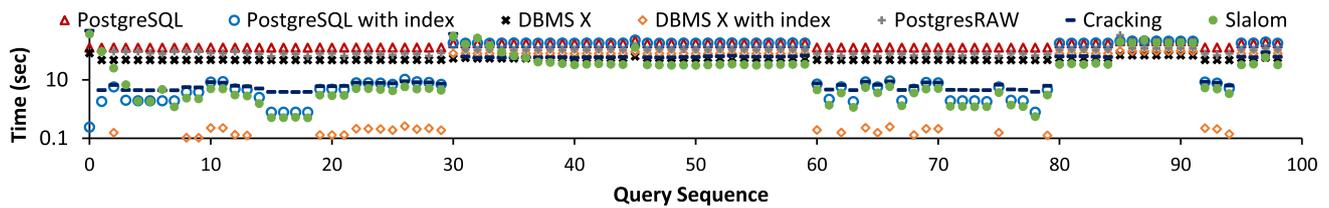


Fig. 13 Sequence of SHD analytics workload. Slalom offers consistently comparable performance to in-memory DBMS

Table 2 Cost of each phase of a smart-meter workload

System	Loading (s)	Index build (s)	Queries (s)	Total (s)
Slalom	0	0	4301	4301
Cracking	0	0	6370	6370
PostgresRaw	0	0	10077	10,077
PostgreSQL (with index)	2559	1449	9058	13,066
PostgreSQL (no index)	2559	0	15,379	17,938
DBMS-X (with index)	6540	1207	3881	11,628
DBMS-X (no index)	6540	0	5243	11783

tains time-stamped information about sensor measurements such as energy consumption and temperature, as well as a sensor ID for geographical tracking. The time stamps are in increasing order. The total size of the dataset is 55 GB in CSV format. We run a typical workload of an SHD analytics application. Initially, we ask a sequence of range queries with variable selectivity, filtering data based on the time-stamp attribute (Q1–29). Subsequently, we ask a sequence of range queries which filter data based on energy consumption measurements to identify a possible failure in the system (Q30–59). We then ask iterations of queries that filter results based on the time-stamp attribute (Q60–79, Q92–94), the energy consumption (Q80–84, Q95–100), and the sensor ID (Q85–91), respectively. Selectivity varies from 0.1 to 30%. Queries focusing on energy consumption are the least selective.

Figure 13 shows the response time of the different approaches for the SHD workload. All systems run from a hot state, with data resting in the OS caches. The indexed versions of PostgreSQL and DBMS-X build a B⁺ tree on the time-stamp attribute. The figure plots only query execution time and does not show the time for loading or indexing for PostgreSQL and DBMS-X. For other systems, where building auxiliary structures takes place during query execution, execution time contains the total cost.

PostgreSQL and DBMS-X without indexes perform full table scans for each query. Q30–60 are more expensive because they are not selective. For queries filtering on the time stamp, indexed PostgreSQL exhibits 10× better performance than PostgreSQL full table scan. Similarly, indexed DBMS-X exhibits 17× better performance compared to DBMS-X full table scan. As the queries using the index become more selective, response time is reduced. For the

queries that do not filter data based on the indexed field, the optimizer of DBMS-X chooses to use the index despite the predicate involving a different attribute. This choice leads to response time slower than the DBMS-X full scan.

PostgresRaw is slightly faster than PostgreSQL without indexes. The runtime of the first query that builds the auxiliary structures (cache, positional map) is 8× slower (374 s) than the average query runtime. For the rest of the queries, PostgresRaw behaves similar to PostgreSQL and performs a full table scan for each query.

After the first query, Slalom identifies that the values of the time-stamp attribute are unique. Thus, it chooses to statically partition the data following the cost model for query-based partitioning (Sect. 4.1) and creates 1080 partitions. Slalom creates the logical partitions during the second query and calculates statistics for each partition. Thus, the performance of Slalom is similar to that of PostgresRaw for the first two queries. During the third query, Slalom takes advantage of the implicit clustering of the file to skip the majority of the partitions and decides whether to build an index for each of the partitions. After Q5, when Slalom has stabilized partitions and already built a number of indexes over them, the performance is better than that of the indexed PostgreSQL variation.

Queries Q2–Q30 represent a best-case scenario for DBMS-X: data reside in memory, and its single index can be used; therefore, DBMS-X is faster than Slalom. After Q29, when queries filter on a different attribute, the performance of Slalom becomes equal to that of PostgresRaw until Slalom builds indexes. Because the energy consumption attribute has multiple appearances of the same value, Slalom decided to use homogeneous partitioning. Q30 to Q59 are not selective, and thus, execution times increase for all systems.

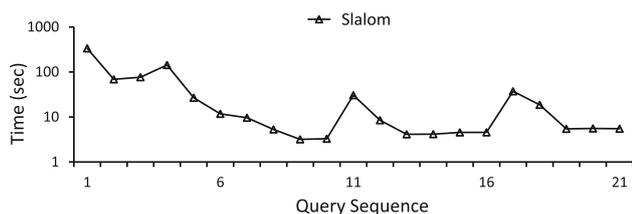


Fig. 14 Slalom executing workload with append-like updates

Table 2 shows the costs for loading and indexing as well as the aggregate query costs for the same query workload of 100 queries, for all the systems. Due to the queries being non-selective, the indexed and non-indexed approaches of DBMS-X have similar performance; thus, in total Slalom exploits its adaptive approach to offer competitive performance to the fully indexed competitors.

Summary Slalom serves a real-world workload which involves fluctuations in the areas of interest and queries of great variety in selectivity. Slalom serves the workload efficiently due to its low memory consumption and its adaptivity mechanisms which gradually lower query response times despite workload shifts.

5.5 Slalom handling file updates

In this section, we demonstrate Slalom’s update efficiency for append-like and in-place updates.

5.5.1 Append-like Updates

Slalom monitors changes in the queried files and dynamically adapts its data structures. In this experiment, we execute a sequence of 20 point queries following the template from Sect. 5 with selectivity 0.1%. Q1 to Q10 run on the original relation of 18 million tuples (22 GB). Between queries Q10 and Q11, we append to the CSV dataset 6 GB of additional uniformly distributed data. Slalom detects the change in the structure of the file and iteratively creates new logical partitions for the new tuples and creates Bloom filters and B⁺ trees during Q11, Q12, and Q13. Between Q16 and Q17, we append again 6 GB of data to the end of the CSV dataset. Slalom again dynamically partitions and builds indexes. Figure 14 shows the execution time for each of the queries in the sequence. Q11 and Q17 execute immediately after the appends; thus, we see higher execution time because Slalom (i) accesses raw data and (ii) builds auxiliary structures—positional maps and binary caches—over them. After this update-triggered spike in execution time, Slalom’s partitioning and indexing schemes converge and the execution time becomes lower and stabilizes.

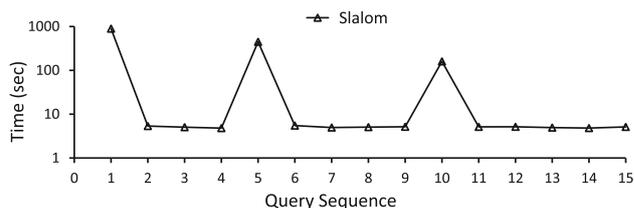


Fig. 15 Slalom executing workload with in-place updates

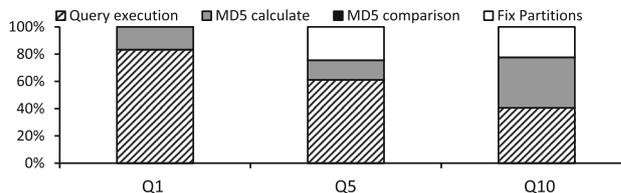


Fig. 16 Time breakdown of query executing with in-place updates

5.5.2 In-place Updates

We now show that Slalom handles in-place updates. We execute a sequence of 15 point queries following the template from Sect. 5 with selectivity 0.01%, run on a 25 million tuple relation (27 GB). We query on a candidate key field to make Slalom use the query-based partitioning strategy and observe solely the effect of updates on a partition. To evaluate update efficiency, we develop a random update generator which updates fields and rows within a file in random places. Before Q5, the update generator updates eight random rows, and before Q10, it updates three random rows. Figure 15 shows the execution time for each of the queries in the sequence. During Q1, Slalom creates 345 partitions and builds the positional map and indexes. During Q5 and Q10, the Update Monitor detects that the file has been updated. Slalom compares the state of all partitions to identify the updated partitions, performs the required corrections to the positional map, and re-builds the indexes. Figure 16 shows this process and presents the breakdown of query execution for Q1, Q5, and Q10. During Q1, along with query execution, Slalom calculates the MD5 codes for all partitions. The update before Q5 touched more partitions than the second update at Q10. Thus, Q5 has more partition data structures to fix. As the query execution progresses, the increasing number of partitions increases the number of checksum calculations.

5.5.3 Speed-up Checksum Calculation

This experiment examines the effect of using GPU and CRC accelerators for the calculation of the partition checksums. We execute 3 point queries following the template from Sect. 5 with selectivity 0.01%, over a 25 million tuple relation (27 GB). To examine the efficiency of GPU and CRC checksum calculation, we vary the number of partitions cre-

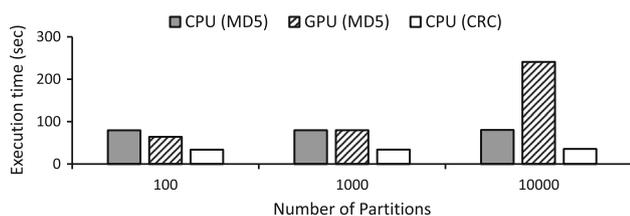


Fig. 17 Checksum calculation using different accelerators with different partition sizes

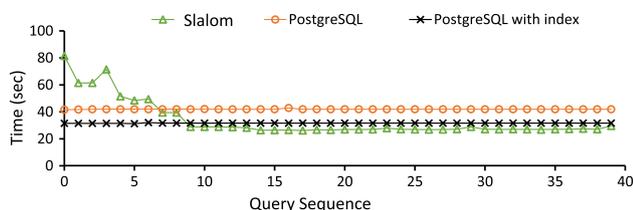


Fig. 18 Sequence of 40 queries over a binary file

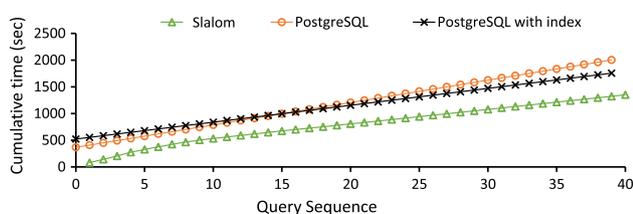


Fig. 19 Cumulative execution time of 40 queries over a binary file

ated by Slalom. The first query breaks the file into 100 equally sized partitions, the second query into 1000 partitions, and the third into 10,000 partitions. Before each query, we make a random update in the file to activate the re-calculation of checksums. Figure 17 shows the checksum calculation cost for the three queries using the three different approaches. When using the CPU (either the dedicated CRC instructions or MD5 calculation), the cost of calculation remains constant. On the other hand, when using the GPU, the checksum calculation is slower when the number of partitions is increasing. The best approach for calculating checksums is using the CRC. However, as CRC is able to compute checksums over input of 1024-byte blocks, it generates a large number of checksums for each partition. Thus, making checksum comparison is more time-consuming.

5.6 Additional data formats: binary data

This section shows that, besides CSV data, Slalom can also operate efficiently over binary datasets. To accommodate binary data, Slalom employs the same techniques as when running over CSV files, with two exceptions. It tunes the cost model to reduce the access cost equations previously associated with text-based data accesses and does not have to build a positional map. Figure 18 compares the execution

time of Slalom and PostgreSQL with and without indexes. We use a binary flat file with 100 million uniformly distributed tuples, each having 30 columns (12 GB), and we run range queries with selectivity 1%. For Slalom, the initial data access is faster than that in the case of CSV data because (i) no parsing is involved and (ii) the binary representation is more compact than the CSV one. During the first nine queries, Slalom fine-tunes its partitioning. During Q3, multiple partitions happened to stabilize, thus triggering the construction of multiple indexes and leading to increased execution overhead. Both PostgreSQL configurations have stable execution times as the selectivity remains stable. Eventually, Slalom and indexed PostgreSQL converge and have similar performance. Figure 19 presents the cumulative execution time for loading, index building, and query execution for the three systems over binary files. PostgreSQL using indexes requires more pre-processing time due to index building, and it takes 13 queries to pay off the cost of building the index. Slalom requires seven queries to start outperforming PostgreSQL, and after ten queries, it offers comparable performance to PostgreSQL with indexes. Table 3 presents separately the time required for loading, index building, and query execution for the three systems. The additional file adapters enable Slalom to efficiently and transparently operate on top of additional data formats.

6 Conclusion

In situ data analysis over large and, crucially, growing datasets faces performance challenges as more queries are issued. State-of-the-art in situ query execution reduces the data-to-insight time. However, as the number of issued queries is increasing and, more frequently, queries are changing access patterns (having variable selectivity, projectivity and are of interest in the dataset), in situ query execution cumulative latency increases.

To address this, we bring the benefits of indexing to in situ query processing. We present *Slalom*, a system that combines an in situ query executor with an online partitioning and indexing tuner. Slalom takes into account user query patterns to reduce query time over raw data by partitioning raw data files *logically* and building for each partition lightweight *partition-specific* indexes when needed. The tuner further adapts its decisions on-the-fly to follow any workload changes and maintains a balance between the potential performance gains, the effort needed to construct an index, and the overall memory consumption of the indexes built.

Acknowledgements We would like to thank the reviewers for their valuable comments. This work is partially funded by the EU FP7 programme (ERC-2013-CoG), Grant No. 617508 (ViDa), the EU FP7

Table 3 Cost of each phase of the 40 query sequence on binary file

System	Loading (s)	Index build (s)	Queries (s)	Total (s)
Slalom	0	0	1352	1352
PostgreSQL (with index)	325	165	1264	1754
PostgreSQL (no index)	325	0	1677	2002

Collaborative project Grant No. 317858 (BigFoot), NSF under Grant No. IIS-1850202, and EU Horizon 2020 research and innovation programme Grant No. 650003 (Human Brain project).

Appendix: Derivation of index construction probability formula

This section provides detailed description of how we derive the probability function for deciding to build an index over a logical partition. We expect this section to be useful for achieving a deeper understanding of the tuning decisions of Slalom. The derivation begins with the expected cost formula.

$$E = \sum_{i=1}^T \left(\sum_{j=1}^{i-1} p_j \cdot C_{use,idx} + \left(1 - \sum_{j=1}^{i-1} p_j\right) \cdot \left(p_i \cdot C_{build,idx} + (1 - p_i) \cdot C_{use,fs} \right) \right)$$

We exchange $C_{build,idx}$ with $C_{use,fs} + \delta$ as building the index will cost at least as much as a full scan.

$$E = T \cdot C_{use,fs} - \left(C_{use,fs} - C_{use,idx} \right) \cdot \left(\sum_{i=1}^T \sum_{j=1}^{i-1} p_j \right) + \delta \cdot \left(\sum_{i=1}^T p_i - \sum_{i=1}^T p_i \cdot \sum_{j=1}^{i-1} p_j \right) \quad (6)$$

We take the first partial derivative of this formula for p_i .

$$\frac{\partial E}{\partial p_i} = - \left(C_{use,fs} - C_{use,idx} \right) \cdot \frac{\partial \left(\sum_{i=1}^T \sum_{j=1}^{i-1} p_j \right)}{\partial p_i} + \delta \cdot \left(\frac{\partial \sum_{i=1}^T p_i}{\partial p_i} - \frac{\partial \left(\sum_{i=1}^T p_i \cdot \sum_{j=1}^{i-1} p_j \right)}{\partial p_i} \right) \quad (7)$$

We calculate that:

$$\frac{\partial \left(\sum_{i=1}^T \sum_{j=1}^{i-1} p_j \right)}{\partial p_i} = (T - i) \quad (8)$$

and

$$\frac{\partial \left(\sum_{i=1}^T p_i \cdot \sum_{j=1}^{i-1} p_j \right)}{\partial p_i} = \sum_{j=1}^{T-1} p_j - p_i \quad (9)$$

Thus, the final derivative becomes:

$$\frac{\partial E}{\partial p_i} = - \left(C_{use,fs} - C_{use,idx} \right) \cdot (T - i) + \delta \cdot \left(1 - \sum_{j=1}^{T-1} p_j - p_i \right) \quad (10)$$

To minimize the expected cost, we solve the equation and we solve for p_i .

$$\frac{\partial E}{\partial p_i} = 0 \Rightarrow p_i = \frac{C_{use,fs} - C_{use,idx}}{\delta} \cdot (T - i) - \left(1 - \sum_{j=1}^{T-1} p_j \right) \quad (11)$$

References

1. Abad, C.L., Roberts, N., Lu, Y., Campbell, R.H.: A storage-centric analysis of MapReduce workloads: file popularity, temporal locality and arrival patterns. In: Proceedings of the IEEE International Symposium on Workload Characterization (IISWC), pp. 100–109 (2012)
2. Abouzied, A., Abadi, D.J., Silberschatz, A.: Invisible loading: access-driven data transfer from raw files into database systems. In: Proceedings of the International Conference on Extending Database Technology (EDBT), pp. 1–10 (2013)
3. Agrawal, S., Narasayya, V.R., Yang, B.: Integrating vertical and horizontal partitioning into automated physical database design. In: Proceedings of the ACM SIGMOD International Conference on Management of Data, pp. 359–370 (2004)
4. Ailamaki, A., DeWitt, D.J., Hill, M.D., Skounakis, M.: Weaving relations for cache performance. In: Proceedings of the International Conference on Very Large Data Bases (VLDB), pp. 169–180 (2001)
5. Alagiannis, I., Borovica, R., Branco, M., Idreos, S., Ailamaki, A.: NoDB: efficient query execution on raw data files. In: Proceedings of the ACM SIGMOD International Conference on Management of Data, pp. 241–252 (2012)
6. Alamoudi, A.A., Grover, R., Carey, M.J., Borkar, V.R.: External data access and indexing in AsterixDB. In: Proceedings of the ACM International Conference on Information and Knowledge Management (CIKM), pp. 3–12 (2015)

7. Alexiou, K., Kossmann, D., Larson, P.-Å.: Adaptive range filters for cold data: avoiding trips to siberia. *Proc. VLDB Endow.* **6**(14), 1714–1725 (2013)
8. Athanassoulis, M., Ailamaki, A.: BF-Tree: approximate tree indexing. *Proc. VLDB Endow.* **7**(14), 1881–1892 (2014)
9. Athanassoulis, M., Idreos, S.: Design tradeoffs of data access methods. In: *Proceedings of the ACM SIGMOD International Conference on Management of Data, Tutorial* (2016)
10. Athanassoulis, M., Kester, M.S., Maas, L.M., Stoica, R., Idreos, S., Ailamaki, A., Callaghan, M.: Designing access methods: the RUM conjecture. In: *Proceedings of the International Conference on Extending Database Technology (EDBT)*, pp. 461–466 (2016)
11. Athanassoulis, M., Yan, Z., Idreos, S.: UpBit: scalable in-memory Updatable Bitmap indexing. In: *Proceedings of the ACM SIGMOD International Conference on Management of Data* (2016)
12. Blanas, S., Wu, K., Byna, S., Dong, B., Shoshani, A.: Parallel data analysis directly on scientific file formats. In: *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pp. 385–396 (2014)
13. Bloom, B.H.: Space/time trade-offs in hash coding with allowable errors. *Commun. ACM* **13**(7), 422–426 (1970)
14. Borwein, P.B.: On the complexity of calculating factorials. *J. Algorithms* **6**(3), 376–380 (1985)
15. Bruno, N., Chaudhuri, S.: An online approach to physical design tuning. In: *Proceedings of the IEEE International Conference on Data Engineering (ICDE)*, pp. 826–835 (2007)
16. Chaudhuri, S., Narasayya, V.R.: An efficient cost-driven index selection tool for microsoft SQL server. In: *Proceedings of the International Conference on Very Large Data Bases (VLDB)*, pp. 146–155 (1997)
17. Chen, Y., Alspaugh, S., Katz, R.H.: Interactive analytical processing in big data systems: a cross-industry study of MapReduce workloads. *Proc. VLDB Endow.* **5**(12), 1802–1813 (2012)
18. Cheng, Y., Rusu, F.: Parallel in-situ data processing with speculative loading. In: *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pp. 1287–1298 (2014)
19. Chou, J.C.-Y., Howison, M., Austin, B., Wu, K., Qiang, J., Bethel, E.W., Shoshani, A., Rübél, O., Prabhat, Ryne, R.D.: Parallel index and query for large scale data analysis. In: *Proceedings of the ACM/IEEE International Conference for High Performance Computing, Networking, Storage and Analysis*, pp. 30:1–30:11 (2011)
20. Clopper, C.J., Pearson, E.S.: The use of confidence or fiducial limits illustrated in the case of the binomial. *Biometrika* **26**(4), 404–413 (1934)
21. DeWitt, D.J., Halverson, A., Nehme, R.V., Shankar, S., Aguilar-Saborit, J., Avanes, A., Flaszka, M., Gramling, J.: Split query processing in polybase. In: *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pp. 1255–1266 (2013)
22. Finkelstein, S.J., Schkolnick, M., Tiberio, P.: Physical database design for relational databases. *ACM Trans. Database Syst. (TODS)* **13**(1), 91–128 (1988)
23. Furtado, C., Lima, A.A.B., Pacitti, E., Valduriez, P., Mattoso, M.: Physical and virtual partitioning in OLAP database clusters. In: *Proceedings of the Symposium on Computer Architecture and High Performance Computing (SBAC-PAD)*, pp. 143–150 (2005)
24. Gankidi, V.R., Teletia, N., Patel, J.M., Halverson, A., DeWitt, D.J.: Indexing HDFS data in PDW: splitting the data from the index. *Proc. VLDB Endow.* **7**(13), 1520–1528 (2014)
25. Graefe, G., Kuno, H.: Self-selecting, self-tuning, incrementally optimized indexes. In: *Proceedings of the International Conference on Extending Database Technology (EDBT)*, pp. 371–381 (2010)
26. Graefe, G., McKenna, W.J.: The volcano optimizer generator: extensibility and efficient search. In: *Proceedings of the IEEE International Conference on Data Engineering (ICDE)*, pp. 209–218 (1993)
27. Grund, M., Krüger, J., Plattner, H., Zeier, A., Cudre-Mauroux, P., Madden, S.: HYRISE: a main memory hybrid storage engine. *Proc. VLDB Endow.* **4**(2), 105–116 (2010)
28. Halim, F., Idreos, S., Karras, P., Yap, R.H.C.: Stochastic database cracking: towards robust adaptive indexing in main-memory column-stores. *Proc. VLDB Endow.* **5**(6), 502–513 (2012)
29. Härder, T.: Selecting an optimal set of secondary indices. In: *Proceedings of the European Cooperation in Informatics (ECI)*, pp. 146–160 (1976)
30. Hu, G., Ma, J., Huang, B.: High throughput implementation of MD5 algorithm on GPU. In: *Proceedings of the International Conference on Ubiquitous Information Technologies & Applications (ICUT)*, pp. 1–5 (2009)
31. Idreos, S., Alagiannis, I., Johnson, R., Ailamaki, A.: Here are my data files. Here are my queries. Where are my results? In: *Proceedings of the Biennial Conference on Innovative Data Systems Research (CIDR)*, pp. 57–68 (2011)
32. Idreos, S., Kersten, M.L., Manegold, S.: Database cracking. In: *Proceedings of the Biennial Conference on Innovative Data Systems Research (CIDR)* (2007)
33. Idreos, S., Kersten, M.L., Manegold, S.: Self-organizing tuple reconstruction in column-stores. In: *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pp. 297–308 (2009)
34. Idreos, S., Manegold, S., Kuno, H., Graefe, G.: Merging what's cracked, cracking what's merged: adaptive indexing in main-memory column-stores. *Proc. VLDB Endow.* **4**(9), 586–597 (2011)
35. Idreos, S., Zoumpatianos, K., Athanassoulis, M., Dayan, N., Hentschel, B., Kester, M.S., Guo, D., Maas, L., Qin, W., Abdul, W., Sun, Y.: The periodic table of data structures. *Bull. IEEE Comput. Soc. Tech. Comm. Data Eng.* **41**(3), 64–75 (2018)
36. Ivanova, M., Kersten, M.L., Manegold, S.: Data vaults: a symbiosis between database technology and scientific file repositories. In: *Proceedings of the International Conference on Scientific and Statistical Database Management (SSDBM)*, pp. 485–494 (2012)
37. Jindal, A., Dittrich, J.: Relax and let the database do the partitioning online. In: *Proceedings of the International Conference on Very Large Data Bases (VLDB)*, pp. 65–80 (2011)
38. Kargin, Y., Kersten, M.L., Manegold, S., Pirk, H.: The DBMS—your big data sommelier. In: *Proceedings of the IEEE International Conference on Data Engineering (ICDE)*, pp. 1119–1130 (2015)
39. Karlin, A.R., Manasse, M.S., McGeoch, L.A., Owicki, S.S.: Competitive randomized algorithms for non-uniform problems. In: *Proceedings of the Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pp. 301–309 (1990)
40. Karpathiotakis, M., Alagiannis, I., Ailamaki, A.: Fast queries over heterogeneous data through engine customization. *Proc. VLDB Endow.* **9**(12), 972–983 (2016)
41. Karpathiotakis, M., Alagiannis, I., Heinis, T., Branco, M., Ailamaki, A.: Just-in-time data virtualization: lightweight data management with ViDa. In: *Proceedings of the Biennial Conference on Innovative Data Systems Research (CIDR)* (2015)
42. Karpathiotakis, M., Branco, M., Alagiannis, I., Ailamaki, A.: Adaptive query processing on RAW data. *Proc. VLDB Endow.* **7**(12), 1119–1130 (2014)
43. Kerrisk, M.: *The Linux programming interface: a Linux and UNIX system programming handbook*. No Starch Press, San Francisco (2010)
44. Kester, M.S., Athanassoulis, M., Idreos, S.: Access path selection in main-memory optimized data systems: should I scan or should I probe? In: *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pp. 715–730 (2017)
45. Kornacker, M., Behm, A., Bittorf, V., Bobrovitsky, T., Ching, C., Choi, A., Erickson, J., Grund, M., Hecht, D., Jacobs, M., Joshi, I., Kuff, L., Kumar, D., Leblang, A., Li, N., Pandis, I., Robinson, H., Rorke, D., Rus, S., Russell, J., Tsiogiannis, D., Wanderman-

- Milne, S., Yoder, M.: Impala: a modern, open-source SQL engine for Hadoop. In: Proceedings of the Biennial Conference on Innovative Data Systems Research (CIDR) (2015)
46. Lightstone, S., Teorey, T.J., Nadeau, T.P.: Physical Database Design: The Database Professional's Guide to Exploiting Indexes, Views, Storage, and More. Morgan Kaufmann, Burlington (2007)
 47. López-Blázquez, F., Mino, B.S.: Binomial approximation to hypergeometric probabilities. *J. Stat. Plan. Inference* **87**(1), 21–29 (2000)
 48. McCrary, S.: Implementing algorithms to measure common statistics. *VLDB J.* **8**, 1–17 (2015)
 49. Melnik, S., Gubarev, A., Long, J.J., Romer, G., Shivakumar, S., Tolton, M., Vassilakis, T.: Dremel: interactive analysis of web-scale datasets. *Proc. VLDB Endow.* **3**(1), 330–339 (2010)
 50. Moerkotte, G.: Small materialized aggregates: a light weight index structure for data warehousing. In: Proceedings of the International Conference on Very Large Data Bases (VLDB), pp. 476–487 (1998)
 51. Mühlbauer, T., Rödiger, W., Seilbeck, R., Reiser, A., Kemper, A., Neumann, T.: Instant loading for main memory databases. *Proc. VLDB Endow.* **6**(14), 1702–1713 (2013)
 52. O'Neil, P.E.: Model 204 architecture and performance. In: Proceedings of the International Workshop on High Performance Transaction Systems (HPTS), pp. 40–59 (1987)
 53. Papadomanolakis, S., Ailamaki, A.: AutoPart: Automating schema design for large scientific databases using data partitioning. In: Proceedings of the International Conference on Scientific and Statistical Database Management (SSDBM), pp. 383 (2004)
 54. Pearson, K.: Contributions to the mathematical theory of evolution. II. Skew variation in homogeneous material. *Philos. Trans. R. Soc. Lond.* **186**(Part I), 343–424 (1895)
 55. Peterson, W.W., Brown, D.T.: Cyclic codes for error detection. *Proc. IRE* **49**(1), 228–235 (1961)
 56. Petraki, E., Idreos, S., Manegold, S.: Holistic indexing in main-memory column-stores. In: Proceedings of the ACM SIGMOD International Conference on Management of Data (2015)
 57. Richter, S., Quiané-Ruiz, J.-A., Schuh, S., Dittrich, J.: Towards zero-overhead static and adaptive indexing in Hadoop. *VLDB J.* **23**(3), 469–494 (2013)
 58. Rivest, R.L.: The MD5 message-digest algorithm. *RFC* **1321**, 1–21 (1992)
 59. Schnaitter, K., Abiteboul, S., Milo, T., Polyzotis, N.: COLT: continuous on-line database tuning. In: Proceedings of the ACM SIGMOD International Conference on Management of Data, pp. 793–795 (2006)
 60. Schuhknecht, F.M., Jindal, A., Dittrich, J.: The uncracked pieces in database cracking. *Proc. VLDB Endow.* **7**(2), 97–108 (2013)
 61. Sidirourgos, L., Kersten, M.L.: Column imprints: a secondary index structure. In: Proceedings of the ACM SIGMOD International Conference on Management of Data, pp. 893–904 (2013)
 62. Sinha, R.R., Mitra, S., Winslett, M.: Bitmap indexes for large scientific data sets: a case study. In: Proceedings of the IEEE International Symposium on Parallel and Distributed Processing (IPDPS) (2006)
 63. Sun, L., Franklin, M.J., Krishnan, S., Xin, R.S.: Fine-grained partitioning for aggressive data skipping. In: Proceedings of the ACM SIGMOD International Conference on Management of Data, pp. 1115–1126 (2014)
 64. Thusoo, A., Sarma, J.S., Jain, N., Shao, Z., Chakka, P., Anthony, S., Liu, H., Wyckoff, P., Murthy, R.: Hive—a warehousing solution over a map-reduce framework. *Proc. VLDB Endow.* **2**(2), 1626–1629 (2009)
 65. Wang, X., Yu, H.: How to break MD5 and other hash functions. In: Proceedings of the Annual International Conference on the Theory and Applications of Cryptographic Techniques (EUROCRYPT), pp. 19–35 (2005)
 66. Wu, E., Madden, S.: Partitioning techniques for fine-grained indexing. In: Proceedings of the IEEE International Conference on Data Engineering (ICDE), pp. 1127–1138 (2011)
 67. Wu, K., Ahern, S., Bethel, E.W., Chen, J., Childs, H., Cormier-Michel, E., Geddes, C., Gu, J., Hagen, H., Hamann, B., Koegler, W., Lauret, J., Meredith, J., Messmer, P., Otoo, E.J., Perevotzhikov, V., Poskanzer, A., Rübel, O., Shoshani, A., Sim, A., Stockinger, K., Weber, G., Zhang, W.-M.: FastBit: interactively searching massive data. *J. Phys.: Conf. Ser.* **180**(1), 012053 (2009)
 68. Zilio, D.C., Rao, J., Lightstone, S., Lohman, G.M., Storm, A., Garcia-Arellano, C., Fadden, S.: DB2 design advisor: integrated automatic physical database design. In: Proceedings of the International Conference on Very Large Data Bases (VLDB), pp. 1087–1097 (2004)

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.