

RailwayDB: adaptive storage of interaction graphs

Robert Soulé¹ · Buğra Gedik²

Received: 13 April 2015 / Revised: 17 September 2015 / Accepted: 6 October 2015 / Published online: 16 October 2015
© Springer-Verlag Berlin Heidelberg 2015

Abstract We are living in an ever more connected world, where data recording the interactions between people, software systems, and the physical world is becoming increasingly prevalent. These data often take the form of a temporally evolving graph, where entities are the vertices and the interactions between them are the edges. We call such graphs interaction graphs. Various domains, including telecommunications, transportation, and social media, depend on analytics performed on interaction graphs. The ability to efficiently support historical analysis over interaction graphs requires effective solutions for the problem of data layout on disk. This paper presents an adaptive disk layout called the *railway layout* for optimizing disk block storage for interaction graphs. The key idea is to divide blocks into one or more sub-blocks. Each sub-block contains the entire graph structure, but only a subset of the attributes. This improves query I/O, at the cost of increased storage overhead. We introduce optimal integer linear program (ILP) formulations for partitioning disk blocks into sub-blocks with overlapping and nonoverlapping attributes. Additionally, we present greedy heuristics that can scale better compared to the ILP alternatives, yet achieve close to optimal query I/O. We provide an implementation of the railway layout as part of RailwayDB—an open-source graph database we have developed. To demonstrate the benefits of the railway layout, we provide an extensive experimental evaluation, includ-

ing model-based as well as empirical results comparing our approach to baseline alternatives.

Keywords Interaction graphs · Adaptive storage · I/O optimization

1 Introduction

We are living in an ever more connected world, where the data generated by people, software systems, and the physical world are more accessible than before and are much larger in volume, variety, and velocity. In many application domains, such as telecommunications, transportation, and social media, live data recording the interactions between people, systems, and the environment is available for analysis. These data often take the form of a temporally evolving graph, where entities are the vertices and the interactions between them are the edges. We call such graphs *interaction graphs*.

Data analytics performed on interaction graphs can bring new business insights and improve decision making. For instance, the graph structure may represent the interactions in a social network, where finding communities in the graph can facilitate targeted advertising. In the telecommunications (telco) domain, call details records (CDRs) can be used to capture the call interactions between people, and locating closely connected groups of people can be used for generating promotions.

Interaction graphs are temporal in nature, and more importantly, they are append-only. This is in contrast to relationship graphs, which are updated via insertion and deletion operations. An example of a relationship graph is a social network capturing the follower–followee relationship among users. Examples of interaction graphs include CDR graphs cap-

✉ Robert Soulé
robert.soule@usi.ch

Buğra Gedik
bgedik@cs.bilkent.edu.tr

¹ Faculty of Informatics, Università della Svizzera italiana, Lugano, Switzerland

² Department of Computer Engineering, Bilkent University, Ankara, Turkey

turing calls between telco customers or mention graphs capturing interactions between users of a micro-blogging service, like Twitter.

Since interaction graphs can potentially grow forever, they present a storage challenge for system designers. Even on modern servers with large amounts of memory, one cannot assume that the entire graph will fit into the main memory. The append-only nature of interaction graphs make storing them on disk a necessity. Furthermore, the analysis of this historical interaction data forms an important part of the analytical landscape.

The ability to efficiently support historical analysis over interaction graphs requires effective solutions for the problem of *data layout* on disk. Most graph algorithms are characterized by locality of access [1], which is a direct result of the traversal-based nature of most of the graph algorithms. This is often taken advantage of by co-locating edges in close proximity within the same disk blocks [2]. This way, once a disk block is loaded into main memory buffers, several edges from it can be used for processing, reducing the disk I/O.

In interaction graphs, the locality of access is even more pronounced. First, the analysis to be performed on the interactions can be restricted to a temporal view of the graph, such as finding the influential users over a given week of interactions. This means that edges that are temporally close are accessed together. Second, traversals are again key to many graph algorithms, such as connected components, clustering coefficient, and PageRank. This means that edges that are close in terms of the path between their incident vertices and their timestamps should be located together with the same blocks. In our earlier work [3], we introduced an interaction graph database that works on this principle of access locality. It uses a disk organization that consists of a set of blocks, each containing a list of *temporal neighbor lists*. A temporal neighbor list contains a head vertex and a set of incident edges within a time range. The layout optimizer aims at bringing together, into the same disk block, temporal neighbor lists that are (i) close in terms of their temporal ranges, (ii) have many edges between them, and (iii) have few edges going into temporal neighbor lists outside the block.

Many real-world graph databases contain attributes. In the case of interaction graphs, the attributes can be considered as properties associated with the edges representing the interactions. Attributes can be stored in two ways, either separately (e.g., in a relational table) or locally with the temporal neighbor lists. If they are stored separately, then the graph database cannot take full advantage of locality optimizations performed for block organization. The database must go back and forth between the disk blocks to access the edge attributes. On the other hand, if attributes are stored locally within the disk blocks containing the graph structure, then there can be significant overhead due to disk I/O when only a few attributes are needed to answer a query.

To query an interaction graph, most algorithms traverse the graph structure to access the relevant attributes. Frequently, there are correlations among the attributes accessed by different queries. For example, queries q_1 and q_5 might access attributes a_1 and a_2 , while queries q_2 , q_3 and q_4 access attributes a_3 and a_4 . Because interaction graphs are temporal, the co-access correlations for the attributes can vary for different temporal regions. Moreover, the co-access correlations might be unknown at the insertion time, but be discovered later, when the workload is known.

It is widely recognized that query workload and disk layout have a significant impact on database performance [4–6]. For table-based relational databases, this fact has led database designers to develop alternative approaches for storage layout: Row-oriented storage [7] is more efficient when queries access many attributes from a small number of records, and column-oriented storage [8] is more efficient when queries access a small number of attributes from many records [6]. Unfortunately, although interaction graph databases, like relational databases, are the target of diverse query workloads, there is no clear correspondence to a row-oriented or column-oriented storage layout.

This paper presents an adaptive disk layout called the *railway layout* and associated algorithms for optimizing disk block storage for interaction graphs. The key idea is to divide blocks into one or more sub-blocks, where each sub-block contains a subset of the attributes (potentially overlapping), but the entire graph structure is replicated within each sub-block. This way, a query can be answered completely by only reading the sub-blocks that contain the attributes of interest, reducing the overall I/O. The core concept is equally applicable to relationship graphs, yet the motivation is much stronger for interaction graphs, as the continuously increasing nature of the interaction graphs rule out in-memory processing.

There are a number of challenges in achieving an effective adaptive layout. First, we need to find the partitioning of attributes that minimizes the query I/O. To address this, we model the problems of overlapping and nonoverlapping attribute partitioning as integer linear programs (ILPs) and provide optimal solutions that minimize the query I/O cost. Second, the query workload and thus the attribute access pattern can change over time. For this purpose, our railway layout supports customization of the attribute partitioning of sub-blocks on a per-block basis. Third, such flexibility necessitates online configuration of attribute partitioning as the query workload evolves, which in turn requires fast algorithms for performing the attribute partitioning. For this purpose, we develop greedy heuristic algorithms for both overlapping and nonoverlapping partitioning scenarios. These algorithms can scale to larger numbers of attributes, yet provide close to optimal query I/O performance. Finally, the railway layout trades off storage space to gain improved query I/O performance. The storage overhead is more pro-

nounced for the case of overlapping partitioning. To address this, we limit the amount of storage overhead that can be tolerated, and integrate this limit to both our ILP formulations, as well as our greedy heuristics.

We implemented railway layout as part of RailwayDB—an open-source graph database developed as a testbed for our research. To demonstrate the benefits of the railway layout, we provide an extensive experimental evaluation, including both model-based and empirical results comparing our approach to baseline alternatives. The results show that, for a storage increase of just 25 %, the optimal overlapping partitioning algorithm reduces the query I/O cost by 45 %. When allowed to double the storage usage, the overlapping partitioning algorithm can reduce the I/O cost by 73 %. The heuristic algorithm performs almost as well, reducing the I/O cost by 72 %, but cuts the running time needed to find a solution by orders of magnitude.

In summary, this paper makes the following contributions:

- We introduce the railway layout for adaptive organization of interaction graphs on disk.
- We introduce optimal ILP formulations for partitioning disk blocks into sub-blocks with overlapping and nonoverlapping attributes, given a query workload. Our formulation also support upper bounding the amount of storage overhead introduced as a result of the railway layout.
- To support online adaptation, we develop greedy heuristics that can scale better compared to the ILP alternatives, yet achieve close to optimal query I/O.
- We describe a practical implementation of the railway layout within our open-source RailwayDB.
- We provide an extensive experimental study comparing our approach to a few baseline alternatives.

The rest of the paper is organized as follows. Section 2 gives an overview of the railway layout in the context of an interaction database and motivates its design. Section 3 formalizes the optimal railway layout design problem. Section 4 gives integer linear programming formulation of the optimal layout for overlapping and nonoverlapping scenarios. Section 5 introduces our heuristic solutions for the same. Section 6 describes our design of an interaction graph database using the railway layout. Section 7 presents the implementation details. Section 8 reports an experimental evaluation of our system. Section 9 discusses related work and Sect. 10 concludes the paper.

2 Adaptive storage overview

The design of the railway disk layout builds on our prior work [3], which organized the disk layout for interaction graph databases to improve access locality. However, that

work did not consider the I/O cost due to the edge attributes, which is the major contributor to the total disk I/O during query processing. The railway layout addresses this issue by enabling the system to adapt the layout to changing workloads, with the goal of reducing the disk I/O during querying, in exchange for a slight increase in the disk space used to store the graph.

2.1 Interaction graphs

Several systems have previously used the term interaction graph in an informal way [3,9,10]. For this paper, we assume a discrete, ordered time domain, \mathcal{T} , and a domain of attributes, \mathcal{A} . We make no assumptions about the values of those attributes. An interaction graph is an ordered pair $\mathcal{I} = (V, E)$ such that V is a set of vertices and $E \subseteq \{(\tau, u, l, v) : u, v \in V \wedge \tau \in \mathcal{T} \wedge l \subseteq \mathcal{A}\}$ is a set of interactions. Each interaction (τ, u, l, v) has a timestamp τ , a source vertex u , a destination vertex v , and a set of attributes l . Interaction graphs support append-only write operations. They may be queried by specifying a time range $[t_s, t_e]$, and operations such as selection and projection which should be applied to the attributes of the edges that fall in the time range. Queries may also take a source vertex, s , as an additional parameter, in which case the query only examines attributes on edges incident to s .

2.2 Motivating example

To explain the design of the railway layout, we first introduce a small, motivating example. Figure 1 shows a graph for the

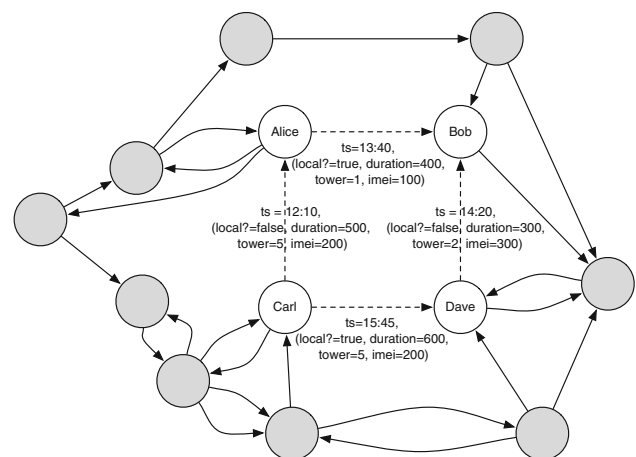


Fig. 1 A partial example interaction graph for call data records, capturing the telephone calls among a set of people. The running example focuses on a subgraph indicated by the nodes colored white. Each edge in the graph is associated with attributes for the interaction, including the time the call was placed, whether the call was local, the duration of the call, the cell phone tower, and the IMEI number identifying the device used

telephone call interactions among a set of people. Each node in the graph represents a person, and each edge in the graph represents a phone call from a caller to a callee. Each edge is associated with a set of attributes that maintain the details of the interaction, including the time the call was placed, whether the call was local or long distance, the duration of the call, the cell phone tower, and the IMEI number identifying the device used to place the call. Thus, the schema for the edges in the graph is as follows¹:

```
call(local?, duration, tower, imei)
```

Recall that interaction graphs are append-only and evolve over time. In other words, new timestamped edges are continuously added to the graph. For explanatory purposes, we focus on a subset of a graph at a particular time range. In the figure, the subset is indicated by the white nodes and the edges between them. In this subset, there were four call interactions. One of them was a local call from Alice to Bob, starting at 13:40. They spoke for 400 seconds. The call was made from cell phone tower 1, and the Alice's phone had an IMEI number of 100.

A telecommunications company performs various analytics by processing the graph. For example, in order to understand how they should price their service plans, they might want to capture the duration of all calls for each user. To plan for infrastructure provisioning, they might want to record a count of the number of calls that each cell phone tower handled.

In an interaction graph, queries are associated with a time range, $[t_{start}, t_{end}]$. To answer queries, the graph database system must traverse the subgraph that contains edges with timestamp t , such that $t_{start} \leq t \leq t_{end}$. As the system traverses the graph, it reads the relevant attributes to answer the query. Note that a query might access all or some of the attributes. As concrete examples, imagine we have two queries. Query q_1 asks for the average duration for calls from each tower, broken down by local or nonlocal nature of the calls. Query q_2 asks for the count of calls made by each type of device. In other words, we say that each query accesses a subset of the attributes:

```
q1 = {local?, duration, tower}, q2 = {imei}
```

2.3 Storage for locality

There are several ways in which one might store a graph on disk. The graph structure can be stored as a matrix representation or an adjacency list. Most graph databases choose an adjacency list representation because they reduce the storage overhead when graphs are sparse, and it is faster to iterate

¹ Since time is an implicit part of an interaction, we do not show it as part of the schema.

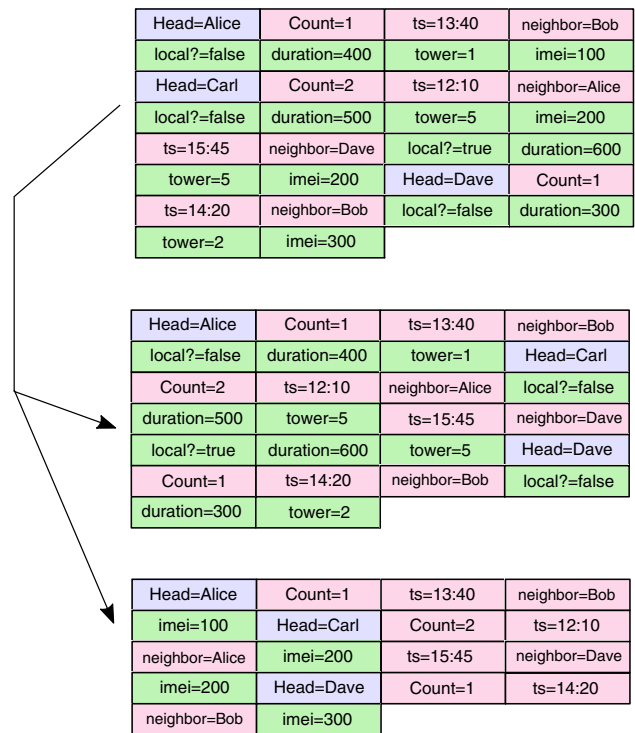


Fig. 2 Standard disk block storage for an interaction graph, and a partitioning into sub-blocks for the railway layout. Each sub-block maintains its own copy of the neighbor list, and a subset of the attributes

over the edges when traversing the graph. Attributes associated with each edge could be stored separately in a relational table, or along with the edges. Storing the attributes with the edges improves the locality, since the database can read the graph structure and associated attributes from the same disk block. To improve locality, typical disk layout schemes try to group as many adjacent nodes as possible in the same disk block.

Building on this basic design, our prior work [3] extended the notion of locality to include a temporal dimension for handling interaction graphs. Nodes are placed in the same block if they are close together both spatially and temporally. Based on the edge timestamps, the adjacency lists are divided into multiple pieces, and based on closeness of the nodes within the graph, these partial adjacency lists are combined into blocks. The locality of a block is determined by its *conductance* (i.e., the percentage of edges going out of a block) and its *cohesiveness* (i.e., a metric used to find highly connected components). Our earlier work describes a greedy algorithm for forming disk blocks with respect to this notion of locality.

Once the algorithm divides the graph into disk blocks, the graph data and attributes are stored in the layout scheme illustrated in the top of Fig. 2. Note that this is an adjacency list representation in which attributes are stored with the edges. Each disk block contains a sequence of vertices, identified

by a *head-node id*, followed by a *count* of the number of neighbors, and then the neighbor list itself. Each entry in the neighbor list is composed of a *timestamp*, an *id* for the destination vertex, and the properties for that edge. In the top of Fig. 2, all of the information from the example interaction subgraph is stored in a single disk block.

2.4 Railway layout

This paper introduces a new disk layout scheme, called the *railway layout* illustrated in the bottom of Fig. 2. With the railway layout scheme, blocks are partitioned into sub-blocks, such that each sub-block contains the adjacency list representation from the original block, but only a subset of the attributes. The subset of attributes assigned to each sub-block is determined by the query workload.

For example, given queries q_1 and q_2 , the railway layout would store the attributes *local?*, *duration*, and *tower* in one sub-block, and the attribute *imei* in a second sub-block. In the ideal case, a query can be answered completely by reading a single sub-block that contains only the relevant information and none of the irrelevant information, reducing the overall I/O cost. Of course, this layout comes at the expense of storage, as the graph structure information is duplicated in each sub-block. We argue that in general, I/O cost is more important than storage overhead, because a certain level of storage overhead can be accommodated by

adding additional disks. In Sects. 3 and 5, we will present our optimal and heuristic algorithms for discovering the sub-block partitions that keep the overhead below a user-specified threshold, while minimizing the disk I/O for the queries.

2.5 Adaptation

Because interaction graphs are append-only, and new edges are continuously added, there is a unique opportunity to adapt the disk layout with changing workloads over time. A database system utilizing the railway layout design can continually monitor the workload, and re-adjust the disk layout for historical data. This is illustrated in Fig. 3. In the figure, we have an interaction graph with four attributes, namely *a*, *b*, *c*, and *d*. Initially, without any workload optimization, all disk blocks have a single sub-block that contains the entire set of attributes. This is shown in the upper half of the figure. After some time, the database adapts to the workload. This is shown in the lower half of the figure. We see that blocks from different time ranges have adapted differently, as the workload they observe is different. For instance, blocks B_X and B_Y were partitioned into two sub-blocks as $\{a, b, c\}$ and $\{c, d\}$, whereas blocks B_Z and B_U were partitioned into three sub-blocks as $\{a\}$, $\{b, c\}$, and $\{c, d\}$, and the blocks B_V and B_W stayed intact as $\{a, b, c, d\}$. A partition index is kept to track the partitioning of blocks in different time regions of

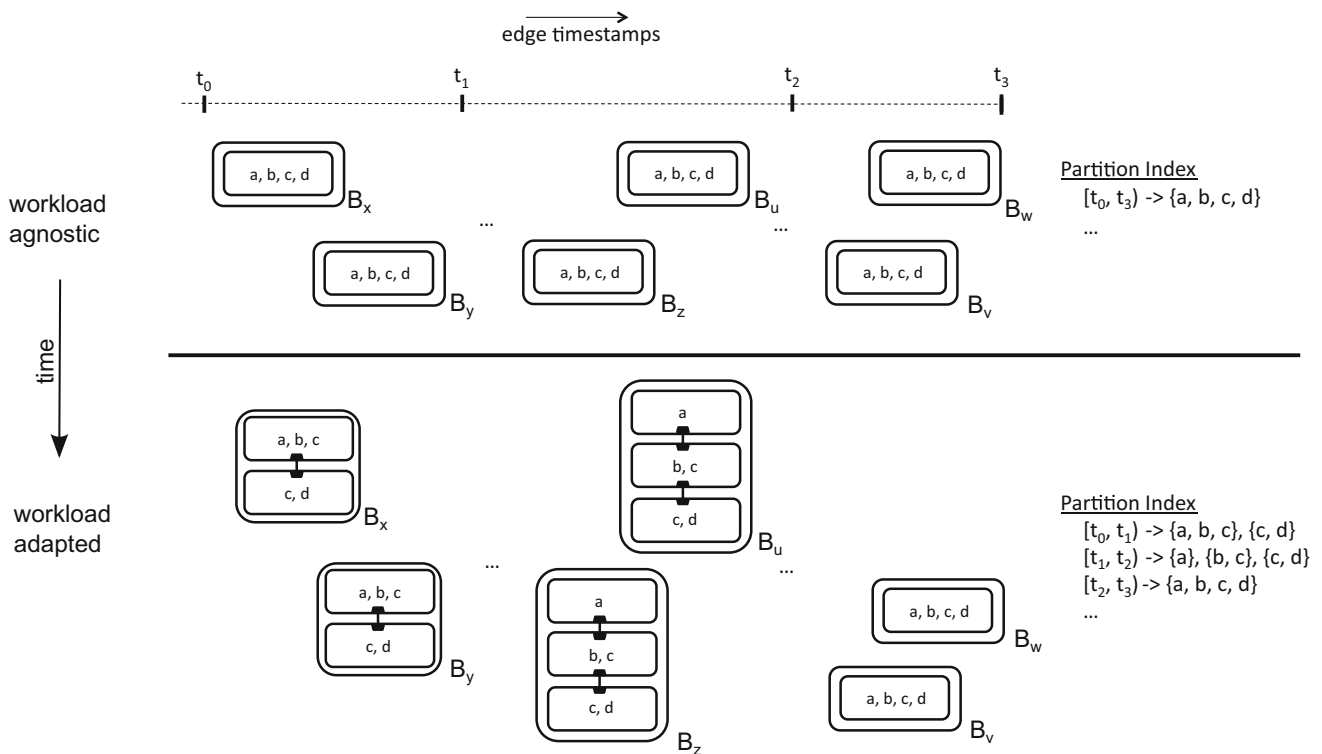


Fig. 3 A database system implementing the railway layout will adapt the disk storage over time

the interaction graph, which is shown on the right-hand side of the figure.

Sections 3, 4, and 5 of this paper focus on the problem of how to determine the best partitioning for a given workload, which is the key capability that enables the adaptation. Sections 6 and 7 describe how the system is implemented in practice.

3 Disk block partitioning problem

The optimal railway design concerns the partitioning of disk blocks into sub-blocks such that the query I/O is minimized, while the storage overhead induced is kept below a desired threshold. This optimization is guided by the query workload observed by the disk blocks within a given time range. Thus, the optimization problem is localized to a sequence of disk blocks that are in the same temporal range and the partitioning used for the sub-blocks created could be potentially different for disk blocks from different time ranges.

The partitioning of disk blocks into sub-blocks can be *nonoverlapping* or *overlapping*. In the nonoverlapping case, the attributes are partitioned among the sub-blocks with no overlap (i.e., a true partitioning). In the overlapping case, the subset of attributes contained within sub-blocks can overlap (i.e., an attribute can appear in multiple partitions). In both cases, the complete graph structure for the block is replicated within the sub-blocks, which results in a storage overhead.

In both overlapping and nonoverlapping partitioning, we trade increased storage overhead for reduced query I/O cost. In the overlapping case, the increase in the storage overhead is higher, as some of the attributes are replicated, in addition to the graph structure. On the other hand, enabling overlapping attributes is expected to reduce the query I/O (in the extreme case, there could be one sub-block per query). While the nonoverlapping partitioning scenario is a special case of the overlapping one, specialized algorithms can be used to solve the former problem.

In the rest of this section, we first introduce basic notation and then formulate the overall optimization problem. The modeling of the query I/O and storage overhead are presented next, which complete the formalization of the optimal railway design problem.

3.1 Basic notation

Let Q be the query workload, where each query $q \in Q$ accesses a set of attributes $q.A$ and traverses parts of the graph for the time range $q.T = [q.ts, q.te]$. Note that when we refer to a query, we mean *query kind*. That is, if q_1 is “all calls with a duration > 100” and q_2 is “all calls with a duration > 500”, then they have the same kind, as they only access the `duration` attribute. We denote the set of

all attributes as A . Given a block B , we denote its time range as $B.T$, which is the union of the time ranges of its temporal neighbor lists. Let $s(a)$ denote the size of an attribute a . We use $c_n(B)$ to denote the number of temporal neighbor lists within block B and $c_e(B)$ to denote the total number of edges in the temporal neighbor lists within the block. We overload the notation for block size and use $s(B)$ to denote the size of a block B . We have:

$$s(B) = c_e(B) \cdot \left(16 + \sum_{a \in B.A} s(a)\right) + c_n(B) \cdot 12 \quad (1)$$

Here, 16 corresponds to the cost of storing the edge id and the timestamp, and 12 corresponds to the cost of storing the head vertex (8 bytes) plus the number of entries (4 bytes) for a temporal neighbor list.

Our goal is to create a potentially overlapping partitioning of attributes for block B , resulting in a set of sub-blocks denoted by $\mathcal{P}(B)$. In other words, we have $\bigcup_{B' \in \mathcal{P}(B)} B'.A = A$. Here, \mathcal{P} denotes the partitioning function.

3.2 Optimization problem

We aim to find the partitioning function \mathcal{P} that minimizes the query I/O over B , while keeping the storage overhead below a limit, say $1 + \alpha$ times the original. The original corresponds to the case of a single block that contains all the attributes. Let us denote the query I/O as $L(\mathcal{P}, B)$ and the storage overhead as $H(\mathcal{P}, B)$. Then our goal is to find:

$$\mathcal{P} \leftarrow \operatorname{argmin}_{\{\mathcal{P}: H(\mathcal{P}(B)) < \alpha\}} L(\mathcal{P}, B) \quad (2)$$

3.3 Storage overhead formulation

The storage overhead is defined as the additional amount of disk space used to store the sub-blocks, normalized by the original space needed by a single block (no partitioning). The storage overhead can be formalized as follows, for the nonoverlapping case:

$$H(\mathcal{P}, B) = (|\mathcal{P}(B)| - 1) \cdot \left(1 - \frac{c_e(B) \cdot \sum_{a \in A} s(a)}{s(B)}\right) \quad (3)$$

Basically, for the nonoverlapping case, there is no overhead due to the attributes, as they are not repeated. However, there is overhead for the block structure that is repeated for each sub-block. There are $|\mathcal{P}(B)| - 1$ such extra sub-blocks, and for each, the contribution to the overhead due to storing the graph structure is given by $s(B) - c_e(B) \cdot \sum_{a \in A} s(a)$. Equation 3 has one nice feature, that is, it does not depend on the details of the attribute partitioning, other than the number of partitions. We make use of this feature, later for the ILP formulation of the problem.

For the general case of overlapping partitioning, we can formulate the storage overhead as follows:

$$H(\mathcal{P}, B) = \frac{\sum_{B' \in \mathcal{P}(B)} s(B')}{s(B)} - 1 \tag{4}$$

This formulation follows directly from the definition of storage overhead. While simple, it depends on the details of the partitioning, as $s(B')$ is the size of a sub-block B' , which in turn depends on the list of attributes within the sub-block.

3.4 Query I/O formulation

Let m be a function that maps a query q to the set of sub-blocks that are accessed to satisfy it for a relevant block B under a given partitioning \mathcal{P} .

For the case of nonoverlapping attributes, the m function lists all the sub-blocks whose attributes intersect with those from the query. Formally:

$$m(\mathcal{P}, B, q) = \{B' : B' \in \mathcal{P}(B) \wedge q.A \cap B'.A \neq \emptyset\} \tag{5}$$

For the case of overlapping attributes, we use a simple heuristic to define the set of sub-blocks used for answering the query. Algorithm 1 captures it. The idea is to start with an empty list of sub-blocks and greedily add new sub-blocks to it, until all query attributes are covered. At each iteration, the sub-block that brings the highest relative marginal gain is picked. The relative marginal gain is defined as the total size of the attributes from the sub-block that contribute to the query result, relative to the sub-block size. While computing the relative marginal gain, attributes that are already covered by sub-blocks that are selected earlier are not considered as contributing to the query result.

Algorithm 1: m -overlapping(\mathcal{P}, B, q)

Data: \mathcal{P} : partitioning function, B : block, q : query
 $S \leftarrow \emptyset; R \leftarrow \emptyset$ \triangleright Selected attributes; Resulting sub-blocks
while $S \subset q.A$ **do** \triangleright While unselected attributes remain
 $B' \leftarrow \operatorname{argmax}_{B' \in \mathcal{P}(B) \setminus R} \sum_{a \in B'.A \cap q.A \setminus S} \frac{c_a(B').s(a)}{s(B')}$ $S \leftarrow S \cup B'.A$
 $R \leftarrow R \cup B'$ \triangleright Extend the selected sub-blocks
return R \triangleright Final set of sub-blocks covering query attributes

Given that we have defined the function m that maps a query to the set of sub-blocks used to answer it, we can now formalize the total query I/O cost for a block under a given workload:

$$L(\mathcal{P}, B) = \sum_{q \in Q} w(q) \cdot \mathbf{1}(q.T \cap B.T \neq \emptyset) \cdot \sum_{B' \in m(\mathcal{P}, B, q)} s(B') \tag{6}$$

We simply sum the I/O cost contributions of the queries to compute the total I/O cost. A query contributes to the total I/O cost if and only if its time range intersects with that of the block ($\mathbf{1}(q.T \cap B.T \neq \emptyset)$). If it does, then we add the sizes of all the sub-blocks used to answer the query to the total I/O cost. Furthermore, we multiply the I/O cost contribution of a query with its frequency, denoted by $w(q)$ in the formula.

4 ILP solution

In this section, we formulate the optimal railway design problem as an integer linear program (ILP). The main challenge is to represent the objective function and the constraint as a linear combination of potentially integer variables.

For the ILP formulation, we define a number of binary (0 or 1) variables:

- $x_{a,p}$: 1 if attribute a is in partition p , 0 otherwise.
- $y_{p,q}$: 1 if partition p is used by query q , 0 otherwise.
- $z_{a,p,q}$: 1 if partition p is used by query q and attribute a is in partition p , 0 otherwise.
- u_p : 1 if partition p is assigned at least 1 attribute, 0 otherwise.

Each of these variables serve a purpose:

- x s define the attribute-to-partition assignments.
- y s help formulate the query I/O contribution of each partition due to the graph structure they contain (excluding their assigned attributes).
- z s help formulate the query I/O contribution of each partition, only considering the attributes they are assigned.
- u s help formulate the storage overhead requirement.

In total, we have $|A| \cdot (|A| + 1) \cdot (|Q| + 1)$ variables. Here, we assume that the maximum number of partitions is fixed. In fact, we cannot have more partitions than attributes, so the number of partitions is upper bounded by $k = |A|$, and thus $0 \leq p < k$. However, some of these partitions can be empty in the optimal solution, which means that the number of partitions found by the ILP solution is typically lower than the maximum possible. A simple post-processing step removes empty partitions and creates the final partitioning to be used.

Finally, we define a helper notation for representing whether a variable is accessed by a query or not: $q(a) \equiv \mathbf{1}(a \in q.A)$.

We are now ready to state the ILP formulation. We separate the cases of nonoverlapping and overlapping partitioning, as the former case can be formulated using a smaller number of constraints.

4.1 Nonoverlapping partitions

We start with the objective function, that is, the total query I/O, which is to be minimized:

$$\sum_{q \in Q} w(q) \cdot \left(\sum_{p=1}^k (16 \cdot c_e(B) + 12 \cdot c_n(B)) \cdot y_{p,q} + \sum_{a \in A} s(a) \cdot c_e(B) \cdot z_{a,p,q} \right) \tag{7}$$

In Eq. 7, we simply sum for each query and each partition, and add the I/O cost of reading in the structural information found in a sub-block, if the partition is used by the query. We then sum over each attribute as well and add the I/O cost of reading in the attributes. Note that $z_{a,p,q}$ could have been replaced with $x_{a,p} \cdot y_{p,q}$, but that would have made the objective function nonlinear.

We are now ready to state our constraints. Our first constraint is that each attribute must be assigned to a single partition. Formally:

$$\forall_{a \in A}, \sum_{p=1}^k x_{a,p} = 1 \tag{8}$$

Our second constraint is that if a query q contains an attribute a assigned to a partition p , then partition p is used by the query, i.e., $y_{p,q} = 1$. In essence, we want to state: $\forall_{\{p,q\} \in [1,\dots,k] \times Q}, y_{p,q} = \mathbf{1}(\sum_{a \in A} q(a) \cdot x_{a,p} > 0)$.

In order to formulate this constraint, we use the following ILP construction: Assume we have two variables, β_1 and β_2 , where $\beta_2 \in [0, 1]$ and $\beta_1 \geq 0$. We want to implement the following constraint: $\beta_2 = \mathbf{1}(\beta_1 > 0)$. This could be expressed as a linear constraint as follows, where K is a large constant guaranteed to be larger than β_1 for all practical purposes:

$$\begin{aligned} \beta_1 - \beta_2 &\geq 0 \\ K \cdot \beta_2 - \beta_1 &\geq 0 \end{aligned} \tag{9}$$

We apply this construction to our second constraint, where $\beta_1 = \sum_{a \in A} q(a) \cdot x_{a,p}$ and $\beta_2 = y_{p,q}$. This results in the following linear constraints:

$$\begin{aligned} \forall_{\{p,q\} \in [1,\dots,k] \times Q}, \sum_{a \in A} q(a) \cdot x_{a,p} - y_{p,q} &\geq 0 \\ \forall_{\{p,q\} \in [1,\dots,k] \times Q}, K \cdot y_{p,q} - \sum_{a \in A} q(a) \cdot x_{a,p} &\geq 0 \end{aligned} \tag{10}$$

Our third constraint is that if an attribute a is assigned to a partition p , and partition p is used by a query q , then the corresponding z variable must be set to 1. That is, we

want: $\forall_{\{a,p,q\} \in A \times [1,\dots,k] \times Q}, z_{a,p,q} = \mathbf{1}(x_{a,p} = y_{p,q} = 1)$. We express this as a linear constraint, as follows:

$$\forall_{\{a,p,q\} \in A \times [1,\dots,k] \times Q}, z_{a,p,q} - (x_{a,p} + y_{p,q}) \geq -1 \tag{11}$$

In Eq. 11, when both the x and y variables both 1, the z variable is simply forced to be 1. Otherwise, the z variable can be either 0 or 1, but since the z variables appear in the objective function as positive terms, the solver will set them to 0 to minimize the I/O cost (note that the z variables do not appear in any other constraint).

Our fourth constraint is that if a partition is nonempty, then its corresponding u variable must be set to 1. In other words, we want $\forall_{p \in [1,\dots,k]}, u_p = \mathbf{1}(\sum_{a \in A} x_{a,p} > 0)$. This is expressed as linear constraints, as follows:

$$\begin{aligned} \forall_{p \in [1,\dots,k]}, \sum_{a \in A} x_{a,p} - u_p &\geq 0 \\ \forall_{p \in [1,\dots,k]}, K \cdot u_p - \sum_{a \in A} x_{a,p} &\geq 0 \end{aligned} \tag{12}$$

Equation 12 uses the same construction as the second constraint, where $\beta_1 = \sum_{a \in A} x_{a,p}$ and $\beta_2 = u_p$.

Our fifth, and the last, constraint deals with the storage overhead. We want to make sure that the storage overhead does not go over α . Recall that for the nonoverlapping attributes case, the storage overhead depends on the number of partitions used (Eq. 3). That means that the only ILP variables it depends on are the u s. In particular, the number of partitions used is given by $\sum_{p=1}^k u_p$. This results in the following linear constraint:

$$\sum_{p=1}^k u_p \leq 1 + \frac{\alpha}{1 - \frac{c_e(B) \cdot \sum_{a \in A} s(a)}{s(B)}} \tag{13}$$

The final ILP formulation for the nonoverlapping partitioning is given in Fig. 4. We have a total of $|A|^2 \cdot |Q| + 2 \cdot |A| \cdot |Q| + 3 \cdot |A| + 1$ constraints and the objective function contains $|A| \cdot |Q| \cdot (1 + |A|)$ variables.

4.2 Overlapping partitions

We present an ILP formulation of the problem, as we did for the case of nonoverlapping partitions in Sect. 4.1. We use the same set of variables and the same objective function. However, the formulation of the constraints differ.

Our first constraint is that each attribute must be assigned to at least one partition. Formally:

$$\forall_{a \in A}, \sum_{p=1}^k x_{a,p} \geq 1 \tag{14}$$

$$\begin{aligned}
 & \text{minimize } \sum_{q \in Q} w(q) \cdot \left(\sum_{p=1}^k (16 \cdot c_e(B) + 12 \cdot c_n(B)) \cdot y_{p,q} \right. \\
 & \qquad \qquad \qquad \left. + \sum_{a \in A} s(a) \cdot c_e(B) \cdot z_{a,p,q} \right) \\
 & \text{subject to} \\
 & \qquad \forall a \in A, \quad \sum_{p=1}^k x_{a,p} = 1 \\
 & \qquad \forall \{p,q\} \in [1..k] \times Q, \quad \sum_{a \in A} q(a) \cdot x_{a,p} - y_{p,q} \geq 0 \\
 & \qquad \forall \{p,q\} \in [1..k] \times Q, \quad K \cdot y_{p,q} - \sum_{a \in A} q(a) \cdot x_{a,p} \geq 0 \\
 & \qquad \forall \{a,p,q\} \in A \times [1..k] \times Q, \quad z_{a,p,q} - (x_{a,p} + y_{p,q}) \geq -1 \\
 & \qquad \forall p \in [1..k], \quad \sum_{a \in A} x_{a,p} - u_p \geq 0 \\
 & \qquad \forall p \in [1..k], \quad K \cdot u_p - \sum_{a \in A} x_{a,p} \geq 0 \\
 & \qquad \sum_{p=1}^k u_p \leq 1 + \frac{\alpha}{1 - \frac{c_e(B) \cdot \sum_{a \in A} s(a)}{s(B)}}
 \end{aligned}$$

Fig. 4 ILP formulation for the nonoverlapping optimal railway design

As our second constraint, we require that for each attribute contained in a query, there needs to be a partition that is used by that query and that contains the attribute in question. Formally:

$$\forall \{a,q\} \in A \times Q, \quad \sum_{p=1}^k z_{a,p,q} \geq q(a) \tag{15}$$

As our third constraint, we require that if a query is using an attribute from a partition, then that partition must contain the attribute. That is, we need to link the z variables with the x variables as $\forall \{a,p,q\} \in A \times [1..k] \times Q, (z_{a,p,q} = 1) \implies (x_{a,p} = 1)$. This can be stated as linear constraints:

$$\forall \{a,p,q\} \in A \times [1..k] \times Q, \quad x_{a,p} - z_{a,p,q} \geq 0 \tag{16}$$

As our fourth constraint, we require that if a query is using at least one attribute from a partition, then that partition must be used by the query, i.e., we need to link the z variables with the y variables as $\forall \{p,q\} \in [1..k] \times Q, y_{p,q} = \mathbf{1}(\sum_{a \in A} z_{a,p,q} > 0)$. As before, we use the ILP construction from Eq. 9 for this, where $\beta_2 = y_{p,q}$ and $\beta_1 = \sum_{a \in A} z_{a,p,q}$. We get:

$$\begin{aligned}
 & \forall \{p,q\} \in [1..k] \times Q, \quad \sum_{a \in A} z_{a,p,q} - y_{p,q} \geq 0 \\
 & \forall \{p,q\} \in [1..k] \times Q, \quad K \cdot y_{p,q} - \sum_{a \in A} z_{a,p,q} \geq 0
 \end{aligned} \tag{17}$$

Our fifth constraint is that if an attribute a is assigned to a partition p , and partition p is used by a query q , then the

corresponding $z_{a,p,q}$ variable must be set to 1. This is same as the formulation for the nonoverlapping case from Eq. 11.

Our sixth constraint is that if a partition is nonempty, then its corresponding u variable must be set to 0. Again, this is same as the formulation for the nonoverlapping case from Eq. 12.

Our seventh, and the last, constraint deals with the storage overhead. However, the storage overhead formulation for the overlapping case is different from the one for the nonoverlapping case. This is because the overhead does not merely depend on the number of partitions, as attributes might have to be read multiple times from different partitions (due to the overlaps). As a result, we express the overhead using base variables as in the objective function. Formally:

$$\begin{aligned}
 & \sum_{p=1}^k \left((16 \cdot c_e(B) + 12 \cdot c_n(B)) \cdot u_p \right. \\
 & \qquad \left. + \sum_{a \in A} s(a) \cdot c_e(B) \cdot x_{a,p} \right) \leq s(B) \cdot (1 + \alpha)
 \end{aligned} \tag{18}$$

The final ILP formulation for the overlapping partitioning is given in Fig. 5. We have a total of $2 \cdot |A|^2 \cdot |Q| + 3 \cdot$

$$\begin{aligned}
 & \text{minimize } \sum_{q \in Q} w(q) \cdot \left(\sum_{p=1}^k (16 \cdot c_e(B) + 12 \cdot c_n(B)) \cdot y_{p,q} \right. \\
 & \qquad \qquad \qquad \left. + \sum_{a \in A} s(a) \cdot c_e(B) \cdot z_{a,p,q} \right) \\
 & \text{subject to} \\
 & \qquad \forall a \in A, \quad \sum_{p=1}^k x_{a,p} \geq 1 \\
 & \qquad \forall \{a,q\} \in A \times Q, \quad \sum_{p=1}^k z_{a,p,q} \geq q(a) \\
 & \qquad \forall \{a,p,q\} \in A \times [1..k] \times Q, \quad x_{a,p} - z_{a,p,q} \geq 0 \\
 & \qquad \forall \{p,q\} \in [1..k] \times Q, \quad \sum_{a \in A} z_{a,p,q} - y_{p,q} \geq 0 \\
 & \qquad \forall \{p,q\} \in [1..k] \times Q, \quad K \cdot y_{p,q} - \sum_{a \in A} z_{a,p,q} \geq 0 \\
 & \qquad \forall \{a,p,q\} \in A \times [1..k] \times Q, \quad z_{a,p,q} - (x_{a,p} + y_{p,q}) \geq -1 \\
 & \qquad \forall p \in [1..k], \quad \sum_{a \in A} x_{a,p} - u_p \geq 0 \\
 & \qquad \forall p \in [1..k], \quad K \cdot u_p - \sum_{a \in A} x_{a,p} \geq 0 \\
 & \qquad \sum_{p=1}^k \left((16 \cdot c_e(B) + 12 \cdot c_n(B)) \cdot u_p \right. \\
 & \qquad \qquad \left. + \sum_{a \in A} s(a) \cdot c_e(B) \cdot x_{a,p} \right) \leq s(B) \cdot (1 + \alpha)
 \end{aligned}$$

Fig. 5 ILP formulation for the overlapping partitioning

$|A| \cdot |Q| + 3 \cdot |A| + 1$ constraints and the objective function contains $|A| \cdot |Q| \cdot (1 + |A|)$ variables.

5 Heuristic solution

The ILP formulation described in Sect. 4 finds an optimal solution to the problem of partitioning disk blocks into sub-blocks such that the query I/O is minimized. Unfortunately, solving these types of constraint problems at scale can become a performance bottleneck, since integer programming is NP-Hard. In a graph database using the railway layout, the layout optimization of a block should be fast enough so that it could be piggybacked on disk I/O when significant workload change that necessitates a new layout is detected. We therefore introduce heuristic algorithms for both overlapping and nonoverlapping partitioning scenarios. Experiments in Sect. 8 demonstrate that these heuristic algorithms show significantly improved running times over the optimal approaches, while still appreciably reducing the query I/O cost.

5.1 Nonoverlapping attributes

For the nonoverlapping attributes scenario, we use a heuristic algorithm that greedily assigns attributes to partitions. The pseudo-code of it is given in Algorithm 2. One complication is that the number of partitions is not known a priori. Yet, we know that the number of partitions is bounded by the number of attributes. Furthermore, the number of partitions cannot be larger than one plus the number of attributes used by the queries, as in the worst case each attribute will be in a partition of its own and the unused attributes will be in a separate partition. As such, we start with a single partition and try increasing the number of partitions, until we hit the maximum number of partitions or the storage overhead goes beyond the threshold α . Among all partition counts tried, the one that provides the lowest query cost is selected as the final partitioning. Note that, for the nonoverlapping scenario, the storage overhead is an increasing function of the number of partitions. As such, once we exceed the storage overhead threshold, we can safely stop trying larger numbers of partitions.

For a fixed number of partitions, the algorithm operates by incrementally assigning attributes to partitions. We consider the attributes in decreasing order of their frequency. This is because the reverse, that is, assigning highly frequent attributes later, may result in making assignments that are hard to balance out later. Initially, all partitions are empty. We pick the next unassigned attribute and evaluate assigning it to one of the available partitions. The assignment that results in the lowest query cost is selected as the best assign-

Algorithm 2: Algorithm for partitioning blocks into sub-blocks with nonoverlapping attributes.

```

Data:  $B$ : block,  $Q$ : set of queries
 $c^* \leftarrow \infty$  ▷ Lowest cost over all # of partitions
 $l \leftarrow \min(|A|, 1 + |\cup_{q \in Q} q \cdot A|)$ 
for  $k = 1$  to  $l$  do ▷ For each possible # of partitions
   $R[i] \leftarrow \emptyset, \forall i \in [1, \dots, k]$  ▷ Initialize partitions
  for  $a \in A$ , in decr. order of  $f(a)$  do ▷ For each attribute
     $c \leftarrow \infty$  ▷ Lowest cost over all assignments
     $j \leftarrow -1$  ▷ Best partition assignment
    for  $i \in [1, \dots, k]$  do ▷ For each partition assignment
       $R[i] \leftarrow R[i] \cup \{a\}$  ▷ Assign attribute
      if  $L(R, B, Q) < c$  then ▷ If query cost is lower
         $c \leftarrow L(R, B, Q)$  ▷ Update the lowest cost
         $j \leftarrow i$  ▷ Update the best partition
       $R[i] \leftarrow R[i] \setminus \{a\}$  ▷ Un-assign attribute
     $R[j] \leftarrow R[j] \cup \{a\}$  ▷ Assign to best partition
  if  $H(R, B, Q) > \alpha$  then break ▷ If solution infeasible
  if  $L(R, B, Q) < c^*$  then ▷ If solution has lower cost
     $c^* \leftarrow L(R, B, Q)$  ▷ Update the lowest cost
     $\mathcal{P}(B) \leftarrow R$  ▷ Update the best partitioning
return  $\mathcal{P}(B)$  ▷ Final set of sub-blocks

```

ment and is applied. When computing the query cost, we only consider the attributes assigned so far.

5.1.1 Computational complexity

The computational complexity of the algorithm is $\mathcal{O}(k^2 \cdot |A| \cdot |Q|)$, where k is the maximum number of partitions tried. The $|Q|$ term is the number of unique queries and comes from the cost of computing the query I/O (this can be computed incrementally, even though this is not shown in the pseudo-code). While in the worst case we have $k = |A|$, resulting in a computational complexity of $\mathcal{O}(|A|^3 \cdot |Q|)$, in practice k is much lower due to the upper bound α on the storage overhead.

5.2 Overlapping Partitions

For the overlapping attributes scenario, we use a heuristic algorithm that starts with each query in its own partition and greedily merges partitions until the storage overhead is below the limit. The pseudo-code of it is given in Algorithm 3.

We start the algorithm in a state where for each unique query there is a separate sub-block that contains the attributes from that query. If there are attributes not covered by the queries, they are assigned to a special sub-block. This is the “ideal” partitioning, because the I/O cost would be minimized for the workload at hand. However, in most practical settings, this partitioning will have excessive storage overhead. Thus, we iteratively combine the pair of partitions that have the lowest cost. This is repeated until the storage overhead is below the threshold α . The end result is the final overlapping partitioning.

We define the cost of a merge based on the query I/O and storage cost. In particular, we measure the increase in the query I/O due to the merge, per reduction in the storage

Algorithm 3: Algorithm for partitioning blocks into sub-blocks with overlapping attributes.

```

Data:  $B$ : block,  $Q$ : set of queries
 $\mathcal{P}(B) \leftarrow \{q.A : q \in Q\}$ 
 $A' \leftarrow A \setminus \bigcup_{q \in Q} q.A$ 
if  $A' \neq \emptyset$  then
     $\mathcal{P}(B) \leftarrow \mathcal{P}(B) \cup \{A'\}$ 
while  $H(\mathcal{P}, B) > \alpha$  do
     $c^* \leftarrow \infty$ 
     $(b_x, b_y) \leftarrow (\emptyset, \emptyset)$ 
    for  $\{b_i, b_j\} \in \mathcal{P}(B)$  do
         $\mathcal{P}'(B) \leftarrow \mathcal{P}(B) \setminus \{b_i, b_j\} \cup \{b_i \cup b_j\}$ 
         $c \leftarrow \frac{L(\mathcal{P}', B, Q) - L(\mathcal{P}, B, Q)}{H(\mathcal{P}, B) - H(\mathcal{P}', B)}$ 
        if  $c < c^*$  then
             $c^* \leftarrow c$ 
             $(b_x, b_y) \leftarrow (b_i, b_j)$ 
         $\mathcal{P}(B) \leftarrow \mathcal{P}(B) \setminus \{b_x, b_y\} \cup \{b_x \cup b_y\}$ 
return  $\mathcal{P}(B)$ 
    
```

▷ Each query gets its own sub-block
 ▷ Attributes not covered by the queries
 ▷ There are uncovered attributes
 ▷ Add missing attributes
 ▷ Until storage overhead is below α
 ▷ Lowest cost over all sub-block pairs
 ▷ Sub-block pair with the lowest cost
 ▷ For each pair of blocks
 ▷ Cost of merge
 ▷ Cost is lower
 ▷ Update the lowest cost
 ▷ Update the best pair
 ▷ Final set of sub-blocks

space used. We want to minimize this metric. More formally, assuming \mathcal{P} is the partitioning before the merge and \mathcal{P}' is the partitioning after the merge, the utility can be formulated as:

$$\frac{L(\mathcal{P}', B, Q) - L(\mathcal{P}, B, Q)}{H(\mathcal{P}, B) - H(\mathcal{P}', B)}$$

5.2.1 Computational complexity

The computational complexity of the algorithm is $\mathcal{O}(|A| \cdot |Q|^3)$. At each iteration, the algorithm reduces the number of partitions by one and initially there are $|Q|$ partitions. As such, in the worst case, there will be $|Q|$ iterations. The number of pairs considered is bounded by $|Q|^2$. The utility metric can be computed incrementally, but requires iterating over the query attributes, bringing in the $|A|$ term.

6 Database system design

We have implemented an interaction graph database, named RailwayDB, which encompasses the adaptation capabilities outlined in the earlier sections.

Figure 6 presents the system architecture of RailwayDB. The design extends our earlier work [3], with major new components for handling adaptation, which is the main focus of this work. These new components include the *Optimizer*, *Stat Collector*, *Re-partitioner*, and the *Partition Manager*. Here, we give a brief overview of the RailwayDB components and their interactions, with a particular focus on the components that are used for adapting the storage layout.

The RailwayDB system has three entry points. The first entry point is via streaming inserts. As new interactions happen, they are added into the system via streaming inserts. This is shown on the top left side of Fig. 6. The second entry point is the queries. *Interval queries* can be asked to locate vertices that are involved in interactions within a given time interval and *focused interval queries* can be asked to locate all the

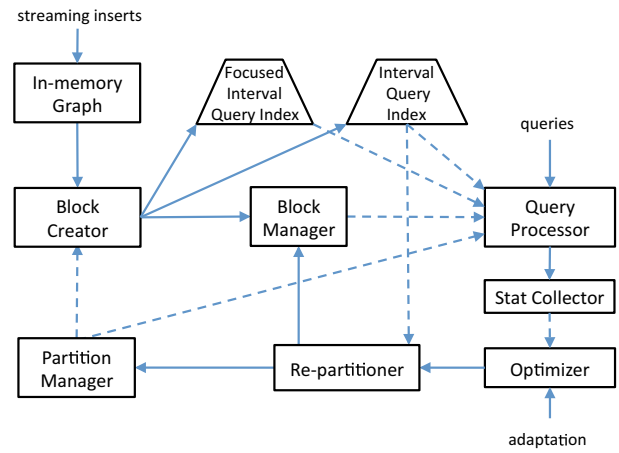


Fig. 6 Architecture of RailwayDB. Solid lines represent ‘writes’ relationships, whereas dashed lines represent ‘reads from’ relationships

interactions of a given vertex within a given time interval. Finally, adaptation can be triggered by asking the optimizer to re-partition the disk layout.

6.1 Streaming inserts

Incoming interactions are first buffered within the *In-memory Graph* component. These edges are also inserted into a FIFO queue. The FIFO queue has a limited temporal extent, and as the interactions expire from the queue, they are removed from the in-memory graph and sent to the *Block Creator* component. This component forms blocks of temporal neighbor lists in batches, with the aim of maximizing locality. The details of this process are given in [3]. The *Block Creator* uses the information it gets from the *Partition Manager* to decide on the partitions it will use for creating sub-blocks. The *Partition Manager* manages the partition index, which provides the partitioning of edge attributes for a given timestamp. This index is queried using the mid-point of the temporal range of a block as the timestamp, in order to locate the partitioning to be used for creating the sub-blocks. The partition index uses an LRU cache to keep the commonly used mappings in memory. Once the sub-blocks are formed, they are forwarded to the *Block Manager*, which writes them to the disk. The *Block Manager* also uses an LRU buffer to keep commonly used blocks in memory. The *Block Creator* also updates the *Focused Query Interval Index* and the *Interval Query Index* with information about the newly created block. We will detail the use of these indexes shortly.

6.2 Queries

Interval queries are supported by the *Interval Query Index*, which indexes the temporal neighbor list time ranges of the

head vertices within each block. Given a time range, this index can locate all head vertices (with their block ids) that are involved in interactions with timestamps in the given time range. The focused interval queries are supported by the *Focused Interval Query Index*, which indexes the pairs of head vertex and temporal neighbor list end timestamps within each block. Given a time range and a vertex, this index can locate all the block ids that contain interactions involving the given vertex with timestamps in the given time range. These two indexes are agnostic to the partitioning of blocks into sub-blocks and always work with a special sub-block called the *master* block. The master block contains all the ids for the sub-blocks corresponding to the partitions. This enables the *Block Manager* to locate all the sub-blocks given the master. The *Query Processor* uses the *Partition Manager* to decide which sub-blocks to retrieve from the *Block Manager*. Depending on the partitions that are needed to answer the query, the list of sub-blocks to be retrieved can change. And this list could be different for different time points within the query time range.

The *Query Processor* also updates the *Stat Collector* as it processes the queries. The *Stat Collector* maintains statistics about which set of attributes are queried how many times. These statistics are maintained separately for different time ranges. In our implementation, we use fixed-size, nonoverlapping time ranges for the purpose of statistics collection.

6.3 Layout Adaptation

The *Optimizer* component uses the statistics kept by the *Stat Collector* to decide on the new partitionings for different time intervals. Various exact solutions and heuristics we described in Sects. 4 and 5 are run as part of the *Optimizer*. The results are fed into the *Re-partitioner*, which is responsible for orchestrating the changes on the disk blocks. For a given set of partitions associated with a time interval, the *Re-partitioner* uses the interval query index to find the list of blocks that are impacted by the change. For each block, it loads the master block and uses the sub-block ids contained there to selectively retrieve and update the sub-blocks. Depending on the changes between the old and the new partitioning, one or more sub-blocks may be added or removed. The retrieval and storage of blocks are performed by interacting with the *Block Manager*. The *Optimizer* then updates the mapping stored in the partition index by calling the *Partition Manager*.

7 System implementation

We have implemented a prototype of the RailwayDB, following the design described in the previous section. The prototype is implemented in C++11 using the LLVM 3.5

compiler. It uses an LSM-tree (LevelDB [11]) for storing the vertex and edge data, as well as the focused interval query index and the partition index. An R-tree (libspatialindex [12]) is used for storing the interval query index.

To solve the ILP formulation of the partitioning problem, the optimizer relies on the C libraries from an integer linear program solver (Gurobi [13]).

In addition to the database, we have implemented infrastructure for evaluating the system, including a workload simulator, which allows us to vary a number of parameters that impact performance, and experiments that use real-world workload from Twitter.

All source codes for the database, simulator, and experiments are publicly available².

8 Evaluation

In this section, we describe two sets of experiments that evaluate our adaptive storage scheme. The first set of experiments is based on the analytic cost model described in Sect. 3 and uses a workload simulator that allows us to evaluate the performance under a number of different workload parameters. The second set of experiments measures the system performance of our database implementation using a real-world data set drawn from Twitter messages. Overall, the results demonstrate that the railway layout scheme significantly reduces query I/O for interaction graphs.

8.1 Environment

We ran all experiments on a machine with a 2.3 GHz Intel i7 processor that has 32 KB L1 data, 32 KB L1 instruction, 256 KB L2 (per core), 6 MB L3 (shared) cache, and 16 GB of main memory. The processor has four cores, but our implementation only uses a single core. The operating system was OS X 10.9.4.

8.2 Model-based experiments

Our first set of experiments evaluates the partitioning algorithms in isolation using the cost model from Sect. 3 and a workload simulator. The simulator allows us to measure the impact of different workload parameters on three performance metrics: (i) the reduction in query I/O due to using the railway layout, (ii) the expected increase in storage cost resulting from the railway layout, and (iii) the scalability of the partitioning algorithms.

² <https://github.com/usi-systems/railwaydb.git>

Table 1 Workload generation parameter defaults

Parameter	Default
# of attributes	10
attribute sizes	Zipf ($z = 0.5$, {4, 1, 8, 2, 16, 32, 64})
query length	Normal ($\mu = 3$, $\sigma = 2.0$)
# of query kinds	5
query kind freq.	Zipf ($z = 0.5$, $n = 5$)
storage ohd. threshold	$\alpha = 1.0$

8.2.1 Workload simulator

The parameters and default settings of the workload simulator are shown in Table 1. The default number of attributes in the graph database schema is taken as 10, even though we experiment with a range of values for it. The size of the attributes come from the list of sizes given in Table 1 and are picked randomly from a Zipf distribution with $z = 0.5$. The average number of unique query kinds we have in the workload for a particular time point is taken as 5, which is another parameter we vary throughout our experiments. The frequencies of different queries follow a Zipf distribution with $z = 0.5$.

8.2.2 Experiment Setup

We measured these three respective values, *query I/O cost*, *storage overhead*, and *running time*, for each partitioning algorithm, as we varied three parameters to the default workload in Table 1:

- *Number of attributes* is the total number of attributes in the iteration graph schema. We first increased the attribute count by multiples of two from 2 to 16. Then, to measure large attribute sets, we increased the attribute count by powers of two from 16 to 128.
- *Number of query kinds* is the number of unique queries in the workload. Queries are of a different kind if the difference of their attribute sets is nonnull. Queries that ask for the same set of attributes, but differ in start node or time interval, are considered to be the same query kind. We increased the number of query kinds by multiples of two from 2 to 16. Beyond 16, the optimal solvers were no longer able to find solutions in a reasonable amount of time.
- *Storage overhead threshold* is the user-specified parameter that dictates how much storage overhead will be tolerated for a solution. We increased the storage overhead threshold by increments of 0.25 from 0 to 2.0.

As baseline comparisons, we also measured the results for two naïve partitioning schemes: SinglePartition places all attributes into a single partition, and PartitionPerAttribute creates a separate partition for each attribute. The SinglePartition scheme represents the standard disk layout [3], and the PartitionPerAttribute approach represents an extreme partitioning (although not an optimal one, as it potentially increases both the query I/O and storage costs). Furthermore, with the SinglePartition approach, the storage overhead is minimized, and thus no other approach can have smaller storage overhead. With the PartitionPerAttribute approach, the query I/O is optimized for single-attribute queries, and thus no other approach can have a lower I/O when only one attribute is queried. These two approaches also correspond to the classic record-oriented (SinglePartition) and column-oriented (PartitionPerAttribute) storage in relational databases.

For each configuration, we ran the experiment 10 times. Each partitioning algorithm used the same workload for each run, but each run was on a different random workload using the same configuration parameters. We report the average (arithmetic mean) and standard deviation.

For all experiments, other than the experiment in which we explicitly altered the value, we used a default storage overhead threshold value of 1.0. We believe this is a reasonable number, as it corresponds with doubling the available storage space. Note that this number is an upper bound on the storage overhead. An optimal partitioning need not use all of the extra space.

8.2.3 Query I/O

Figure 7 shows the results from the query I/O cost measurements. In all three experiments, we see the benefit of the railway layout. The SinglePartition and PartitionPerAttribute layouts represent baseline measurements for a traditional layout and pathological partitioning scheme. All versions of the railway layout result in better query I/O than the baseline measurements, except when the storage threshold is set to not allow any overhead (as we would expect).

In the left graph, we see that the benefits of the railway layout become more pronounced as we increase the number of attributes. At the low end of the graph, with a schema of only two attributes, the optimal overlapping partitioning algorithm results in a 10% reduction in query I/O cost over the SinglePartition scheme. With 16 attributes, there is a 77% reduction in I/O cost. Note that the heuristic overlapping is just as good, also giving a 77% reduction in I/O cost. For large attribute sets, the results are even stronger. With 128 attributes, the heuristic and optimal schemes exhibit a 96% reduction in I/O cost.

In the middle figure, we see that the benefits of the railway layout remain relatively constant as we increase the number

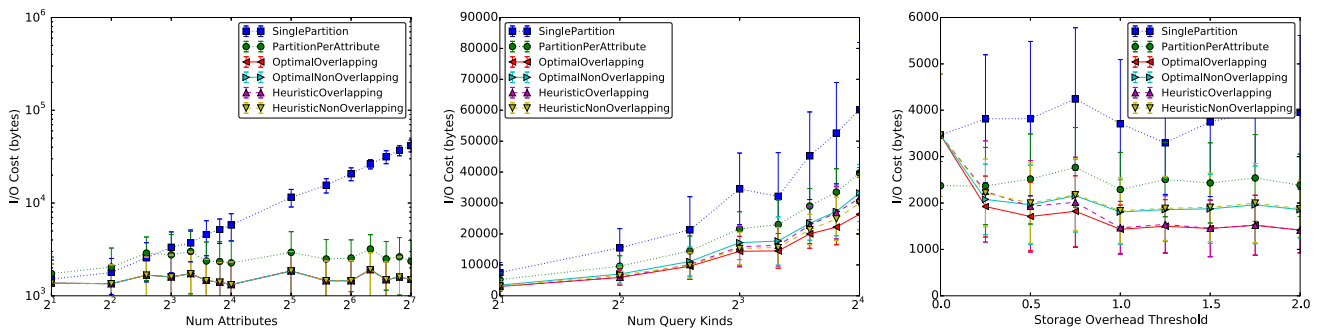


Fig. 7 Query I/O cost for different partitioning algorithms for increasing number of attributes, number of query kinds, and for increasing storage overhead threshold

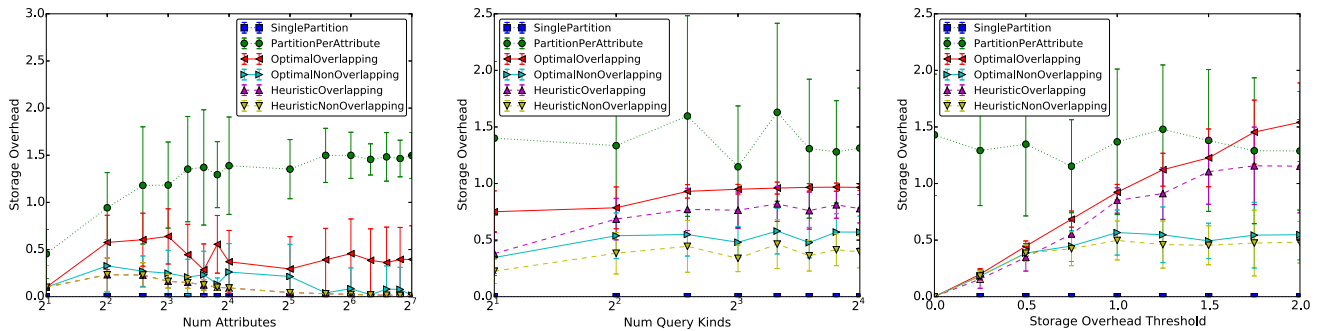


Fig. 8 Storage overhead for different partitioning algorithms for increasing number of attributes, number of query kinds, and for increasing storage overhead threshold

of query kinds. In the case of two query kinds, we see a 60 % difference between the optimal overlapping and single partitioning schemes, while in the case of 16 query kinds, we see a 56 % difference. While increasing the number of query kinds did not have a big impact on query I/O, it did have a large impact on running time, as we will see.

The railway layout makes a tradeoff between query I/O cost and storage cost. We see in the right graph of Fig. 7 that when the user explicitly disallows any increase in storage (i.e., sets the threshold to 0), then the railway layout does not help. However, with even just a slight 25 % increase in storage, all railway layouts reduce query I/O, demonstrating reductions of 45 %.

8.2.4 Storage overhead

The experiments in Fig. 8 quantify the storage overhead that one can expect with using the railway layout. In the left graph, we see that the optimal overlapping and heuristic overlapping approach the user-specified limit of doubling the storage space. As expected, the algorithms will make use of extra storage in order to reduce the query I/O cost. The nonoverlapping schemes are limited in the amount of storage overhead that they use, since they cannot duplicate attributes in separate partitions. So, the extra storage overhead is attributed to duplicating the graph structure.

The middle graph shows a similar result. The overlapping partitioning algorithms approach the user-specified threshold, while the nonoverlapping schemes are bounded.

The right graph in Fig. 8 is interesting. It shows that as the user increases the threshold to a value of 2.0 (i.e., tripling the available storage) both the optimal and heuristic overlapping schemes will try to take advantage of the extra space to reduce query I/O.

8.2.5 Scalability

The experiments in Fig. 9 show the running times for our four algorithms. For attribute sets smaller than 16 attributes, the running times for all schemes are comparable. However, for larger attribute sets, the heuristic approaches demonstrate significantly faster running times. With 32 attributes, the HeuristicNonOverlapping is 94 % faster than the OptimalNonOverlapping scheme. When the schema had 128 attributes, the OptimalOverlapping approach took 18.95 s to find a solution. In contrast, both heuristic solutions took milliseconds to solve.

The number of query kinds had a large impact on solving time. With 16 attributes, the OptimalOverlapping scheme took 17.22 min to find a solution. The OptimalNonOverlapping took 4.99 s. However, the heuristic greedy algorithms were still quite fast. The HeuristicOverlapping took just

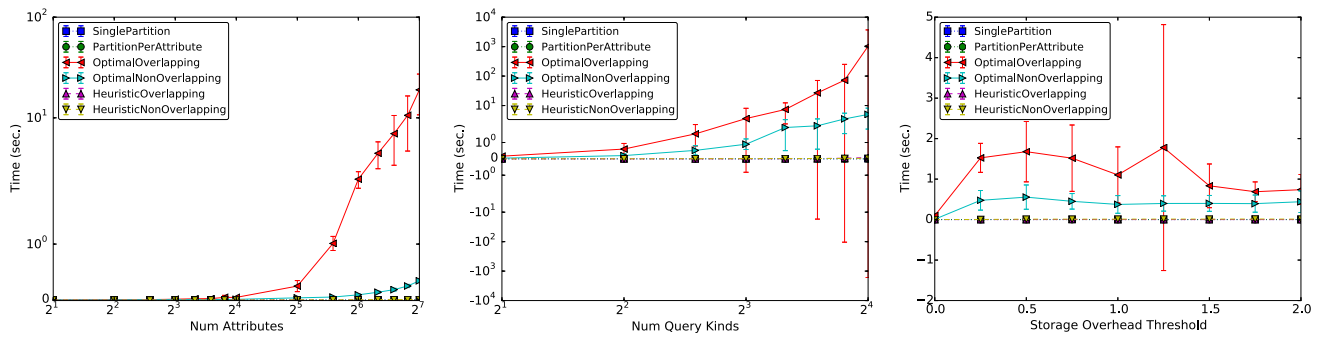


Fig. 9 Running time of different partitioning algorithms for increasing number of attributes, number of query kinds, and for increasing storage overhead threshold

79ms and the HeuristicNonOverlapping took 30ms. After leaving the experiment running for more than 12 h, we were not able to complete the optimal overlapping measurement for the case of 32 query kinds. This experiment demonstrates the benefit of our heuristic greedy algorithms.

However, as shown in the right graph, the storage overhead threshold did not have a significant impact on the running time. This is as expected, since the optimal solvers scale with the number of variables in the constraint problem, and the number of variables does not increase as we alter the storage overhead threshold.

8.2.6 Summary

Overall, our experiments demonstrate the benefits of the railway layout. For a storage increase of just 25 %, the optimal partitioning algorithm reduces the query I/O cost by 45 %. When allowed to double the storage usage, the overlapping partitioning algorithm can reduce the I/O cost by 73 %. The heuristic algorithm performs almost as well, reducing the I/O cost by 72 %, while also reducing the running time needed to find a solution by orders of magnitude.

8.3 System-based experiments

Our second set of experiments evaluates the impact of disk adaptation using the RailwayDB prototype described in Sect. 7 and a real-world data set. Specifically, we measured the actual query I/O and query processing time of the RailwayDB under three different scenarios: (i) varying disk block sizes, (ii) a varying number of query kinds, and (iii) varying the traversal size of the queries.

8.3.1 Twitter data set

We populated our interaction graph with data drawn from Twitter messages from the time interval May 14 to June 05, 2013. The data set contains messages from 500K most prolific Twitter users from Turkey. Interestingly, the time interval

Table 2 Average sizes of attribute data from the Twitter data set

Attribute	Avg. size (bytes)	Attribute	Avg. size (bytes)
Time	12	isTruncated	9
TweetId	22	mentionedUsers	12.9
UserId	12.9	hashTags	6.1
RetweetId	9.9	text	93.9
ReplyToStatus	5	dir	5

during which the data are collected coincides with the 2013 protests in Turkey [14] that generated massive amount of discussion and interaction in the social media. We convert the twitter data into an interaction graph, as follows: If a tweet from user x mentions another user y , then an interaction between x and y is established.

Each tweet has 10 different attributes, which may be of variable size. The names of the attributes and the average (mean) size of each value appear in Table 2.

8.3.2 Queries

To provide a workload, we generated 100 randomized queries on the Twitter data set. Each query includes the following information: (i) a start vertex, v , in the interaction graph, (ii) a start time for the time interval of the query, (iii) an end time for the same, and (iv) a set of attributes to retrieve from the data on the outgoing edges of v .

For each query in the workload, we chose a random vertex using a uniform probability distribution from the entire set of vertices to act as the start vertex. To perform our experiments, we limited the query time ranges to a day’s worth of interactions. This has enabled us to run many queries and report averages as well as standard deviations. The query length (# of attributes used) was a randomly chosen value in the range of [1,10] using a normal distribution, with a mean of three attributes and a standard deviation of 2. The attributes were

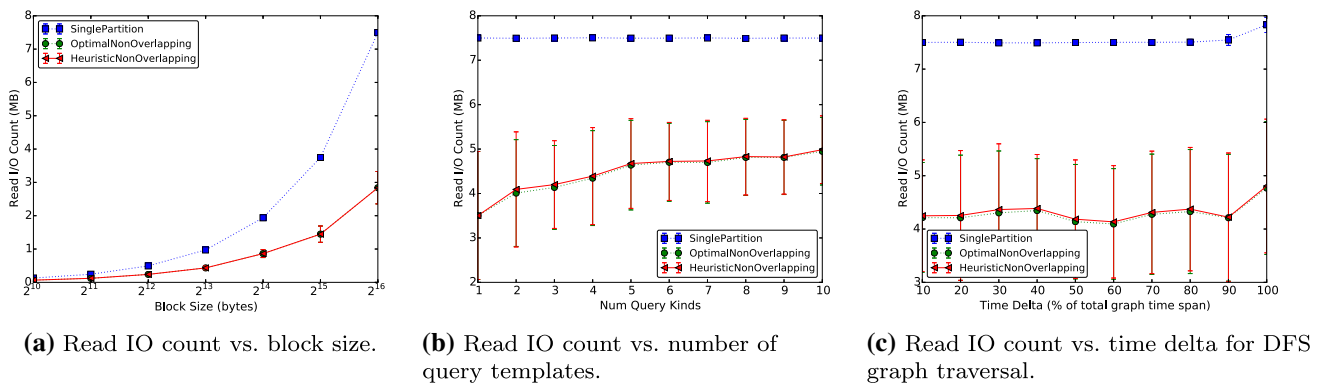


Fig. 10 Experiments on real Twitter data show the benefits in terms of IO cost for the RailwayDB adaptive storage system. **a** Read IO count versus block size, **b** Read IO count versus number of query templates, **c** Read IO count versus time delta for DFS graph traversal

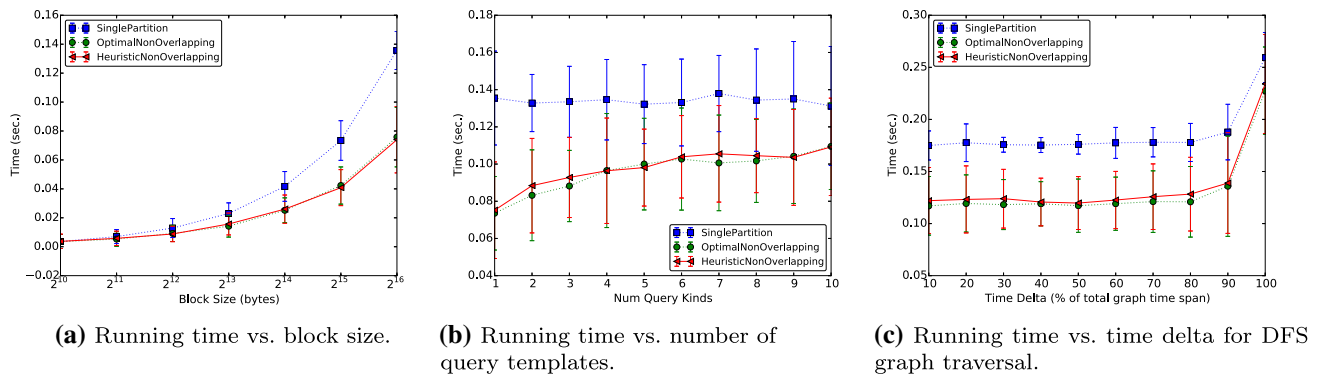


Fig. 11 RailwayDB adaptive storage system demonstrates significant benefits in terms of query processing time for the Twitter data set, **a** Running time versus block size, **b** Running time versus number of query templates, **c** Running time versus time delta for DFS graph traversal

chosen randomly using a Zipf distribution, with a parameter of 0.5, from the set of attributes in Table 2.

8.3.3 Experiment Setup

We measured the *query I/O cost* and *query processing time* using three partitioning algorithms, SinglePartition, OptimalNonOverlapping, and HeuristicNonOverlapping, as we varied three parameters:

- *Block size* is the size of each disk block. We increased the block size by multiples of two, from 1KB to 32KB. For all other experiments, the default block size used was 32KB.
- *Number of query kinds* is the number of unique query types in the workload. Two queries that ask for the same set of attributes, but possibly differ in start node or time interval, are considered to be the same query kind. We increased the number of query kinds from 1 to 10. For all other experiments, the number of query kinds was 3.
- *Traversal size of queries* represents the time ranges of the queries. We increased the time range as a percentage from 10 to 100, the latter representing an entire day's worth of

interactions. For all other experiments, the default time range was 100%.

For each configuration, we ran the experiment 100 times for each partitioning scheme. For each run, we generated a different set of random queries with the same configuration parameters. We report the average (arithmetic mean) and standard deviation. The results of the experiments appear in Figs. 10 and 11.

8.3.4 Block size

Figures 10a and 11a show the query I/O cost and query processing time measurements for the block size experiments. When the block size is small, the total amount of I/O that the database must perform is smaller, and the time to process the queries is also smaller. Thus, the impact of the railway partitioning algorithms is harder to visualize in the graph.

However, there is still a significant difference in the query I/O cost measurements, even when the block size is 1,024 bytes. For example, the OptimalNonOverlapping scheme reads 66,257, while the SinglePartition scheme reads 129,506

bytes, corresponding to a 48 % reduction in I/O. As the block size increases, the affects are more pronounced. When the block size is 65,536 bytes, the OptimalNonOverlapping and HeuristicNonOverlapping both reduce I/O by 62 %.

The query processing time results follow the same trend as the query I/O cost measurements. When the block size is 1,024 bytes, the query processing time is reduced by 2 %. However, when the block size is 65,536 bytes, both the optimal and heuristic approaches reduce the query processing time by 44 %.

8.3.5 Number of query kinds

Figures 10b and 11b show the results for the number of query kinds experiment. The number of query kinds does not impact the SinglePartition scheme, since the database has to read all attributes, regardless of the query workload.

For both query I/O and running time, there is very little difference between the OptimalNonOverlapping and HeuristicNonOverlapping schemes. Both demonstrate a significant reduction in query I/O and query processing time, in comparison with the SinglePartition scheme. When there is a single query kind, both the optimal and heuristic schemes reduce I/O by 53 % and reduce running time by 45 and 44 %, respectively. When there are 10 query kinds, both schemes reduce I/O by 33 % and reduce running time by 16 %.

The graphs show the expected behavior; increasing the number of templates reduces the benefits of the partitioning scheme. However, even with a relatively diverse number of query kinds, the railway scheme significantly reduces I/O and query processing time.

8.3.6 Traversal size of queries

To measure the impact of query traversal size, we modified the workload described above to perform a depth-first search (DFS) bounded by a time interval, rather than query a single node. To increase the traversal size, we increased the time interval for the query as a percentage of the total graph time interval.

Figures 10c and 11c show the results for the traversal size experiment. The benefits of the RailwayDB partitioning schemes remain relatively constant for all measurements, and both the optimal and heuristic approaches demonstrate reductions in read I/O and query processing time. When the time delta is 10 %, both schemes reduce I/O by 43 %. When the time delta is 100 %, the OptimalNonOverlapping reduces I/O by 39 % and the HeuristicNonOverlapping reduces I/O by 38 %. Similarly, when the time delta is 10 % the OptimalNonOverlapping reduces I/O by query processing time by 33 % and the HeuristicNonOverlapping reduces time by

30 %. When the time delta is 100 % the OptimalNonOverlapping reduces I/O by query processing time by 12 % and the HeuristicNonOverlapping reduces time by 10 %.

8.3.7 Summary

All of the experiments demonstrate that the adaptive railway layout significantly reduces the query I/O cost and query processing time when compared to a standard disk layout (i.e., SinglePartition). The effectiveness of the railway layout improves if there are large disk blocks, a smaller number of query kinds, and if queries require more I/O.

9 Related Work

There has recently been increased research interest in large-scale graph analysis and programming models. These include synchronous vertex programming pioneered by Pregel [15], such as Apache Giraph [16], asynchronous vertex programming pioneered by GraphLab [17, 18], and generalized iterated matrix-vector multiplication pioneered by PEGASUS [19]. These systems largely focus on the problem on analytical processing, while our work focuses on data management. Moreover, the graphs these systems provide do not have a temporal dimension.

The railway layout and algorithms build on our prior work [3], which added a temporal dimension to the notion of locality for organizing the disk layout of interaction graph databases. Graph database nodes are placed in the same disk block if they are close together both spatially and temporally. The railway layout extends this design to partition disk blocks into sub-blocks that reduces the query I/O cost. Because interaction graphs are append-only, the railway design enables the disk layout to adapt with changing workloads.

Our adaptation scheme is similar to work on adaptive layouts for relational database. The *H₂O* [4] system can adapt its data layout into three types, row-major, column-major, or groups of columns, depending on the workload. HYRISE [5] provides a similar adaptive layout scheme for an in-memory relational database. Both systems use heuristic, iterative solutions to determine partitioning. The railway layout scheme differs, in that it targets interaction graphs, and we present optimal solutions, in addition to heuristic solutions.

A related area of work in relational databases is adaptive indexing. The goal of adaptive indexing is to create and adjust indexes in response to the queries processed by the system. In this respect, like our work, adaptive indexing is query workload aware. Adaptive indexing approaches include *database cracking* [20, 21], *adaptive merging* [22, 23], or a combination of the two [24]. In database cracking, indexes are created over columns used by the query predicates, by form-

ing range partitioned columns and associated indexes over the range boundaries. As new query predicates are seen for the columns, new partitions are added, getting closer to a sorted column. The idea behind adaptive merging is to initially organize a column as a list of internally sorted partitions, and as new queries arrive, incrementally merge parts of these partitions satisfying query predicates into a single sorted partition. A similar area is fine-grained indexing [25], where indexes are formed only on the parts of the table that are being used by the queries. Different than all these works, RailwayDB partitions the set of attributes (aka columns) and can independently adjust this partitioning for different time ranges over the temporally ordered interactions in the database.

The rise in popularity of social networks, and the recognition that workloads for social network data differ from traditional workloads, has led to increased scrutiny on the problem of disk layout for graph databases. Bondhu [1], the layout manager for the Neo4j graph database [26], aims to minimize the number of seek operations for small user block sizes by fetching multiple friends' data at the same time and by clustering related data into the same block. Bondhu differs from our work in that the cost model does not include a notion of time, nor does it allow for adaptive layouts.

Instead of storing graph data with an adjacency list representation, GBase [27] uses a sparse matrix format. The matrix representation allows GBase to use compression schemes to store homogenous regions of graphs, significantly reducing the storage overhead for large graphs. On top of this storage layout, GBase provides a parallel indexing mechanism that accelerates queries. While the high-level motivations of GBase (i.e., improving query response time for graph database queries) are similar to our work, they are largely focused on the storage overhead. In contrast, we focus on reducing the query I/O cost.

DeltaGraph [28], like our work, includes a temporal component to the layout design, to efficiently support queries over historical graph data. DeltaGraph differs from our work in that they are targeting distributed graph databases that partition data across a set of machines. Consequently, they propose a quite different cost model. Moreover, our railway design lays the foundation for an adaptive disk layout mechanism, which can change over time. Since the DeltaGraph mechanism is static, we expect that the two designs are complementary.

Finally, there is prior work on temporal RDF databases, which aims to improve the response time of SPARQL queries. Notably, Bornea et al. [29] describe a way of mapping an RDF store to a relational database, in order to leverage the overwhelming amount of work on relational database query optimization. Their work is similar to ours in that they use a ILP formulation of a constraint problem in order to optimally determine data placement.

10 Conclusion

Many of today's most popular applications rely on data analytics performed on interaction graphs. The ability to efficiently support historical analysis over interaction graphs requires effective solutions for the problem of data layout on disk. In this paper, we have presented a novel disk layout design for graphs called the railway layout. The design is analogous to hybrid column and row stores in relational databases. Our simulations and experiments show that the railway layout significantly reduces query I/O cost for randomized workloads. We have identified the key challenge for systems to implement the railway layout, which is how to partition blocks into sub-blocks. To solve that problem, we first presented optimal solutions for overlapping and nonoverlapping partitioning using an ILP formulation. To improve the scalability of the partitioner, and enable future work in online adaptation of the disk layout, we have also presented heuristic greedy algorithms that find results close to the optimal solutions, but exhibit faster running times on large graph schemas and workloads. To compare the four partitioning algorithms, we have presented a number of experiments that evaluate the effectiveness and tradeoffs of the various approaches. Overall, the railway layout design appreciably improves the performance of data analytics on interaction graphs.

References

1. Hoque, I., Gupta, I.: Disk layout techniques for online social network data. *Internet Comput.* **16**(3), 24–36 (2012)
2. Steinhaus, R.: G-Store: a storage manager for graph data. Master's thesis, University of Oxford (2011)
3. Gedik, B., Bordawekar, R.: Disk-based management of interaction graphs. *IEEE Trans. Knowl. Data. Eng. (TKDE)* **26**(11), 650–665 (2014)
4. Alagiannis, I., Idreos, S., Ailamaki, A.: H2O: a hands-free adaptive store. In: *ACM International Conference on Management of Data (SIGMOD)*, pp. 1103–1114 (2014)
5. Grund, M., Zeier, A., Krüger, J., Plattner, H., Madden, S.: HYRISE: a main memory hybrid storage engine. *VLDB J. (VLDBJ)* **4**(2), 105–116 (2010)
6. Stonebraker, M., Abadi, D.J., Batkin, A., Chen, X., Cherniack, M., Ferreira, M., Lau, E., Lin, A., Madden, S., O'Neil, E., O'Neil, P., Rasin, A., Tran, N., Zdonik, S.: C-store: a column-oriented DBMS. In: *International Conference on Very Large Data Bases (VLDB)*, pp. 553–564 (2005)
7. Hellerstein, J.M., Stonebraker, M., Hamilton, J.: Architecture of a database system. *Found. Trends Databases* **1**(2), 141–259 (2007). doi:[10.1561/19000000002](https://doi.org/10.1561/19000000002)
8. Abadi, D., Boncz, P., Harizopoulos, S., Idreos, S., Madden, S.: The design and implementation of modern column-oriented database systems. *Found. Trends Databases* **5**(3), 197–280 (2013)
9. Wilson, C., Boe, B., Sala, A., Puttaswamy, K.P., Zhao, B.Y.: User interactions in social networks and their implications. In: *European Conference on Computer Systems (EUROSYS)*, pp. 205–218 (2009)

10. Wilson, C., Sala, A., Puttaswamy, K.P.N., Zhao, B.Y.: Beyond social graphs: user interactions in online social networks and their implications. *ACM Trans. Web (TWEB)* **6**(4), 17:1–17:31 (2012)
11. Ghemawat, S.: LevelDB - a fast and lightweight key/value database library by Google. Available at <https://github.com/google/leveldb>, retrieved March, 2015
12. Hadjieleftheriou, M., Hoel, E.G., Tsotras, V.J.: SaIL: a spatial index library for efficient application integration. *GeoInformatica* **9**(4), 367–389 (2005)
13. Gurobi Optimization Inc.: The Gurobi Optimizer. Available at <http://www.gurobi.com>
14. Turkey Protests Spread from Istanbul to Ankara, Euronews. <http://www.euronews.com/2013/05/31/turkey-protests-spread-from-istanbul-to-ankara/>
15. Malewicz, G., Austern, M.H., Bik, A.J., Dehnert, J.C., Horn, I., Leiser, N., Czajkowski, G.: Pregel: a system for large-scale graph processing. In: *ACM International Conference on Management of Data (SIGMOD)*, pp. 135–146 (2010)
16. Apache giraph. <http://giraph.apache.org>
17. Gonzalez, J.E., Low, Y., Gu, H., Bickson, D., Guestrin, C.: PowerGraph: distributed graph-parallel computation on natural graphs. In: *USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, pp. 17–30 (2012)
18. Kyrola, A., Bleslo, G., Guestrin, C.: GraphChi: large-scale graph computation on just a PC. In: *USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, pp. 31–46 (2012)
19. Kang, U., Tsourakakis, C.E., Faloutsos, C.: PEGASUS: A petascale graph mining system implementation and observations. In: *IEEE International Conference on Data Mining (ICDM)*, pp. 229–238 (2009)
20. Idreos, S., Kersten, M.L., Manegold, S.: Self-organizing tuple reconstruction in column-stores. In: *ACM International Conference on Management of Data (SIGMOD)*, pp. 297–308 (2009)
21. Idreos, S., Kersten, M.L., Manegold, S.: Database cracking. In: *Conference on Innovative Data Systems Research (CIDR)*, pp. 68–78 (2007)
22. Graefe, G., Kuno, H.: Self-selecting, self-tuning, incrementally optimized indexes. In: *International Conference on Extending Database Technology (EDBT)*, pp. 371–381 (2010)
23. Graefe, G., Kuno, H.A.: Adaptive indexing for relational keys. In: *Workshop on Self-Managing Database Systems (SMDB)*, pp. 69–74 (2010)
24. Idreos, S., Manegold, S., Kuno, H., Graefe, G.: Merging what's cracked, cracking what's merged: adaptive indexing in main-memory column-stores. In: *International Conference on Very Large Data Bases (VLDB)*, pp. 586–597 (2014)
25. Wu, E., Madden, S.: Partitioning techniques for fine-grained indexing. In: *IEEE International Conference on Data Engineering (ICDE)*, pp. 1127–1138 (2011)
26. Neo4j Graph Database. <http://www.neo4j.org>
27. Kang, U., Tong, H., Sun, J., Lin, C.Y., Faloutsos, C.: GBASE: A scalable and general graph management system. In: *ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD)*, pp. 1091–1099 (2011)
28. Khurana, U., Deshpande, A.: Efficient snapshot retrieval over historical graph data. In: *IEEE International Conference on Data Engineering (ICDE)*, pp. 997–1008 (2013)
29. Bornea, M.A., Dolby, J., Kementsietsidis, A., Srinivas, K., Dantresangle, P., Udrea, O., Bhattacharjee, B.: Building an efficient rdf store over a relational database. In: *ACM International Conference on Management of Data (SIGMOD)*, pp. 121–132 (2013)