



**SPECIAL ISSUE PAPER**

# Specializing parallel data structures for Datalog

Herbert Jordan<sup>1</sup>  | Pavle Subotić<sup>2</sup> | David Zhao<sup>3</sup> | Bernhard Scholz<sup>3</sup> 

<sup>1</sup>Universität Innsbruck, Innsbruck, Austria

<sup>2</sup>Amazon, London, United Kingdom

<sup>3</sup>The University of Sydney, Sydney, Australia

**Correspondence**

Bernhard Scholz, The University of Sydney,  
Sydney, Australia.

Email: Bernhard.Scholz@sydney.edu.au

**Funding information**

ARC Discovery Project, Grant/Award Number:  
DP180104030; Austrian Research Promotion  
Agency (FFG), Grant/Award Number: 868018

## Summary

We see a resurgence of Datalog in a variety of applications, including program analysis, networking, data integration, cloud computing, and security. The large-scale and complexity of these applications need the efficient management of data in relations. Hence, Datalog implementations require new data structures for managing relations that (1) are parallel, (2) are highly specialized for Datalog evaluation, and (3) can accommodate different workloads depending on the applications concerning memory consumption and computational efficiency. In this article, we present a data structure framework for relations that is specialized for shared-memory parallel Datalog implementations such as the soufflé Datalog compiler. The data structure framework permits a portfolio of different data structures depending on the workload. We also introduce two concrete parallel data structures for relations, designed for various workloads. Our benchmarks demonstrate a speed-up of up to 6× by using a portfolio of data structures compared with using a B-tree alone, showing the advantage of our data structure framework.

## KEYWORDS

B-tree, Datalog, parallel data structure, Trie

## 1 | INTRODUCTION

Datalog is a popular domain-specific language (DSL) for specifying a range of important problems, including static program analysis,<sup>1</sup> network analysis,<sup>2</sup> graph databases,<sup>3</sup> block-chain analysis,<sup>4</sup> among many others. Recently, Datalog has been shown to scale to industrial-sized projects by use of effective shared-memory-based parallelization techniques, and in particular, concurrent data structures<sup>5</sup> that provide the underlying mechanism to perform Datalog evaluation.

A Datalog program consists of a set of rules defined on sets of relations, equivalent to tables in conventional relational database systems. However, not all data within the relations is provided explicitly. Deductive rules define how tuples in relations are formed from the content of other previously computed relations. Datalog engines organize and evaluate the rules to compute output tuples. Efficient state-of-the-art Datalog evaluation strategies, known as *semi-naïve evaluation* strategies,<sup>6</sup> reduce the evaluation of deductive rules to a sequence of relational algebra operations, which require an internal representation for logical relations. Thus, at the core of each Datalog engine<sup>6–9</sup> is a data structure for representing relations that facilitates relational algebra operations including join, selection, and projection operations.

The data structures used in Datalog engines vary, including (in-memory) B-trees, tries, ordinary arrays, and BDDs. The use of specialized data structures has been argued in niche applications. For example, *bddbdb*<sup>7</sup> employs binary decision diagrams to store relations in points-to based static analysis. The *Socialite*<sup>3</sup> engine employs hash-based sets for graph database problems. The *NoD/μZ*<sup>10</sup> engine employs a data structure optimized for Network analysis.<sup>2</sup> These engines argue that their respective data structures achieve the best average case performance for their respective domains. However, despite these strong performance results in certain situations, a single data structure ultimately fails to meet peak performance across diverse sets of benchmarks.

In this paper, we argue for the use of specialized data structures in Datalog engines. We discuss the design of efficient data structures that are specialized for an aspect of Datalog, eg, evaluation algorithm, ruleset characteristics, and dataset characteristics. For example, since Datalog conceptually operates on sets of tuples and assumes monotonicity during the evaluation, a given tuple will only be inserted once in a given relation. Furthermore, the semi-naïve rule evaluation is performed in two phases,<sup>5</sup> namely a read phase and a write phase.<sup>6</sup> This information can be exploited in parallel execution to avoid unnecessary synchronization as writes cannot occur interleaved with these reads.

The approach above is reflected in the design of our Datalog engine Soufflé. Unlike other engines, Soufflé does not depend on a single data structure and does not attempt to provide performance by finding and employing a *universally best data structure*. Instead, Soufflé provides a data structure framework that allows the use of a *portfolio of data structures* during evaluation. In this way, the most effective data structure can be selected by users for a given relation or automatically inferred based on various ruleset/dataset characteristics. Apart from improved overall performance, Soufflé provides researchers a framework to quickly and seamlessly prototype new experimental Datalog-Enabled Relational (DER) data structures for relations without the need for developing yet another Datalog engine. The framework has resulted in the development and integration of several novel specialized data structures and has allowed Soufflé to exhibit strong performance across a wide range of benchmarks, including industrial grade benchmarks sourced from program analysis, network reachability, and block-chain analysis use cases.

Within the Soufflé data structure framework, we have designed two pivotal novel data structures that highlight the need for a portfolio of specialized data structures in Datalog engines.

The first data structure is a specialized concurrent data structure for Datalog that is based on in-memory B-trees and performs well for nondense data. Our data structure is fundamental for the efficient parallel execution of relational algebra operations on contemporary shared-memory multicore machines. Unlike B-Tree implementations in engines such as Logicblox v3,<sup>11</sup> our B-tree is specialized to work with the seminaïve Datalog evaluation and thus exploits the properties of the two phases of seminaïve and thus provides an efficient locking mechanism, as well as a “hint” mechanism to reduce the overheads of tree traversal.

The second data structure is a specialized concurrent data structure for Soufflé called the *Brie*, which has been designed for *high-density* relations with a *large data volume* only. Datalog relations with these characteristics show up in applications such as points-to program analysis. The new data structure adopts design ideas from both B-tree<sup>12</sup> and trie data structures.<sup>13,14</sup> However, the BRIE has several key advantages over both B-trees and tries for dense data only. Compared to B-trees, the BRIE achieves a much higher *coding density* per element by taking advantage of the similarities of maintained keys and thus reducing the amount of duplicated data stored. Compared to tries, the BRIE achieves superior data locality, thus improving cache coherence for dense input data sets.

We have implemented both data structures in the Soufflé data structure framework and have evaluated their performance on parallel microbenchmarks and industrial Datalog benchmarks. We examine the cases where each data structure is effective and their combined usage to achieve high performance.

The contributions of this article are summarized as follows.\*

- We present the properties of a generic data structure specification for Datalog relations. We enable the development of DERs via an interface so that data structures for applications can be adjusted depending on the workload. The data structure interface enables research in novel data structures for Datalog without the need to implement a whole Datalog engine.
- We present an in-memory B-tree that utilizes the data structure framework. The in-memory B-tree has been specialized for Datalog, exploiting the two phases (read/write) of seminaïve evaluation. It has a low synchronization overhead and supports a wide range of applications that have a sparse to moderate data density.
- We present a BRIE data structure for data with high volume and high-density characteristics. The BRIE data structure uses some ideas of tries, with compressed representations of subtrees, and path compression while maintaining some of the cache-friendliness of B-trees.

We also conduct experiments with various industrial and research benchmarks showing the effectiveness of the novel data structures complying to the Soufflé's relational interfaces.

The article is structured as follows. In Section 2, we discuss the need of specialized Datalog relations for parallelized Datalog evaluation. In Section 3, we discuss the properties of a generic data structure interface for Datalog evaluation. In Section 4, we present the specialized in-memory B-tree data structure, and in Section 5, we present the Brie data structure. In Section 6, we evaluate both data structures. In Section 7, we survey related work, and in Section 8, we draw our conclusions.

## 2 | PARALLEL DATALOG EVALUATION

Datalog is a logic programming language, where logical specifications define the computation of input and output relations. Given a set of input relations, the output relations are computed by logic rules. This computation may also involve many intermediate relations in the case of complex Datalog programs.

For example, let *connection* be a binary input relation, representing a set of edges in a graph. The two logical rules

$$tc(X, Y) :- \text{connection}(X, Y) \quad (1)$$

$$tc(X, Z) :- tc(X, Y), \text{connection}(Y, Z) \quad (2)$$

\* This journal article is based on previous conference articles<sup>5,15</sup> written by the coauthors.

```

1  using Tuple = array<size_t,2>;
2  using Relation = set<Tuple>;
3
4  Relation evaluate(const Relation &connection) {
5      Relation tc = connection, deltaTc = connection;
6      while(!deltaTc.empty()){
7          Relation newTc;
8          #pragma omp parallel for
9          for (const auto &t1: deltaTc){
10             auto l = connection.lower_bound({t1[1],0});
11             auto u = connection.lower_bound({t1[1]+1,0});
12             for (auto it=l; it!=u; ++it){
13                 auto& t2 = *it;
14                 Tuple t3({t1[0],t2[1]});
15                 if (tc.find(t3) == tc.end())
16                     newTc.insert(t3);
17             }
18         }
19         tc.insert(newTc.begin(), newTc.end());
20         deltaTc.swap(newTc);
21     }
22     return tc;
23 }

```

**FIGURE 1** Synthesised C++/STL code for path example

implicitly define the content of the output relation  $tc$  computing the transitive closure of the *connection* relation. The first rule states that if there is a direct connection from node  $X$  to node  $Y$ , then  $(X, Y)$  is included in the transitive closure. The second rule is a recursive rule, adding transitive paths inductively. If there is a transitive path from node  $X$  to node  $Y$  and there is a direct connection from  $Y$  to  $Z$ , then there is a transitive path from  $X$  to  $Z$ . While this is a simple example of Datalog, real-world use cases may comprise hundreds of relations, connected via hundreds of potentially mutually recursive rules.

The standard approach for Datalog evaluation is a bottom-up approach. A naïve bottom-up method proceeds in multiple iterations, each of which evaluates all rules in the program and generates new tuples based on the tuples in the previous iteration. This evaluation proceeds until no new facts can be generated, where we say that a *fixpoint* has been reached.

However, the naïve evaluation will recompute tuples that have already been generated in prior iterations, thus leading to wasted computational effort. Therefore, seminaïve evaluation is an optimization which assumes that each iteration should depend on tuples generated in the immediate previous iteration. To this end, seminaïve keeps track of auxiliary relations: a *new* relation and a *delta* relation for each original relation, which keep track of new tuples generated in the current iteration. For our example, the recursive rule of the program would be modified as follows:

$$\begin{aligned}
 tc^0 &:= \text{connection} \\
 \Delta_{tc}^0 &:= \text{connection} \\
 new_{tc}^{i+1}(X, Z) &:= \Delta_{tc}^i(X, Y), \text{connection}(Y, Z). \\
 \Delta_{tc}^{i+1} &:= new_{tc}^{i+1} \setminus tc^i \\
 tc^{i+1} &:= tc^i \cup \Delta_{tc}^{i+1}.
 \end{aligned}$$

In this example, the new knowledge in iteration  $i + 1$  depends directly on tuples computed in iteration  $i$ . Once the new knowledge is computed, two set operations are required to merge the new knowledge into the existing relation  $tc$ .

The Soufflé<sup>8</sup> Datalog engine compiles the transitive closure example to parallel C++ code as listed in Figure 1. The code outlines the anatomy of the underlying least fixed-point calculations of the seminaïve evaluation<sup>6</sup> using STL's set data structures for relations. The synthesized function `evaluate()` has the input relation *connection* as an argument and computes the output relation  $tc$ . The nonrecursive rule (1) causes the initialization of relation  $tc$  with the tuples of relation *connection* (line 5).

The fixed-point calculation for the recursive rule (2) is performed in the while-loop from line 6 to line 21. In each iteration of the while-loop, new tuples are produced for relation  $tc$ . Thereby, the auxiliary relations *newTc* and *deltaTc* are used to avoid recomputations of already generated tuples. If no further tuples can be found, the fixed-point algorithm will stop, and relation  $tc$  will be the transitive closure of *connection*.

The seminaïve algorithm thus unrolls recursive rules into a sequence of nonrecursive rules. In the outlined implementation, these nonrecursive rules are computed using a *nested-loop join*, ie, a nested loop structure that filters tuples from relations and results in a potential insertion in the innermost loop. The filtering mechanism could be implemented by traversing over all tuples in a relation and checking an equality predicate to determine if a tuple should be included in an intermediate result. However, this full relation scans would be computationally inefficient. Instead, we maintain tuples lexicographically ordered and exploit this property to efficiently narrow down iteration ranges to tuples of matching attributes. This improves the search from a linear to logarithmic time complexity.

For example, the nested loop-join in lines 9-12 finds matching connections  $t2 \equiv (Y, Z)$  for each newly discovered transitive connection  $t1 \equiv (X, Y)$  in relation *deltaTc*. The adjacent connection  $t2$  can be found efficiently because the STL set uses a lexicographical order over the edge set. The lexicographical order for two edges  $(u, v)$  and  $(u', v')$  is defined as  $(u, v) < (u', v')$  if and only if  $u < u' \vee (u = u' \wedge v < v')$ . If the node-set is totally ordered (ie, for each node, there exists a unique number), then the lexicographical order is a total order of the connection set. With the lexicographical order, STL's red-black tree efficiently performs the range traversal of all matching/adjacent connections, and no scan over the whole relation *connection* is necessary. If a matching connection pair is found, a new transitive connection  $t3 \equiv (X, Z)$  (cf. line 14) is constructed. If  $t3$  is not yet known (line 15), it will be added to the set of newly discovered transitive connections (line 16). Finally, newly discovered entries are merged into the *tc* relation and promoted to be the *deltaPath* set for the next iteration (line 19 and line 20).

Except the insert operation in line 16, all operations within the nested for-loop between line 9 and line 18 are read-only operations or targeting nonshared memory locations. Thus, the recursive rule can be made parallel (eg, by merely parallelizing the outermost for-loop in line 9) by providing a synchronized insert operation. However, STL's set implementation is not intended for concurrent use and thus does not support synchronized insertion.

As outlined, for optimum performance of Datalog evaluation, indexed data structures are utilized, allowing the fast insertion and retrieval of tuples. However, due to the complexity of data access patterns in real-world Datalog programs, multiple indices may be required for each relation. The construction of these indices is covered in the work of Subotić et al,<sup>16</sup> where a polynomial time algorithm is described for the selection of optimal indices. Each index then uses a set data structure to store the underlying tuples in the order described by the index, and the B-tree and BRIE data structures presented in this paper are two candidates suitable for this task.

Datalog's rule evaluation has particular use patterns of the underlying set data structure. These use patterns can be exploited for better parallel performance. For example, the usage of a set data structure has two phases<sup>17</sup>: (1) either there are multiple writers but no readers or (2) there are multiple readers but no writers. For the example in Figure 1, the nested loops in line 9 and line 12 read pairs from the relations *deltaTc* and *connection* but do not insert new pairs into *deltaTc/connection*. Newly found pairs are inserted in relation *newTc*, but no read operations on *newTc* exist in this loop nest. The seminaïve evaluation guarantees these two phases for all involved sets in its rule evaluation, ie, there are no interleaved reads and writes. Hence, read operations do not need synchronization, and only the write operation requires synchronization. Furthermore, the tuples of the relations in the body of a rule are traversed in a sorted fashion since the relations are sorted. Hence, newly found tuples are generated in a sorted way as well. Hence, the set data structure can exploit a tuple order to minimize the search overheads in the data structure.

### 3 | ABSTRACT RELATION SPECIFICATION

In this section, we describe the properties required for DER data structures that can be effectively employed for parallelized Datalog evaluation. The DER interface for data-structures facilitates the isolation of query processing logic from the underlying concurrent data structures. In general, the workloads differ significantly depending on the application, and hence, a single data structure cannot meet the needs of all applications. Therefore, we introduce a data structure framework that facilitates the design of new data structures in parallel Datalog engines. Our framework dictates properties for data structures to store Datalog relations so that they can seamlessly be deployed and swapped in a parallel Datalog engine. The resulting DER contract provides the foundation for Datalog engines built on top. It comprises basic API interface specifications as well as restrictions on the tuple types to be maintained and concurrent interaction patterns to be supported.

#### 3.1 | Functional requirements

A DER data structure stores a set of same-size tuples  $\{t_1, \dots, t_n\}$ . The functional requirement of a set container is that no two tuples are the same in the set, ie,  $\forall i \neq j : t_i \neq t_j$ . We assume that a tuple is a fixed-length vector of numbers. Note that tuples may have elements that describe strings, symbols, or even more complex elements such as record.<sup>8</sup> However, we assume that there exists a string/symbol table that translates symbols or strings to numbers. This translation mechanism permits a compact and fixed-length representation of tuple elements even for more complex objects.

Any data structure that uses our framework implements an abstract data-type interface for a set container storing vectors of numbers. A DER data structure implements a set of operation that conforms to the operations that are necessary for executing a parallelized version of seminaïve evaluation.<sup>8,18</sup> A DER must implement the following operations.

- *insert(t)* inserts a fixed sized n-ary integer tuple  $t$  into a set of n-ary tuples concurrently, ignoring duplicates;
- *begin()* and *end()* provides iterators to traverse the set concurrently;
- *lower\_bound(a)* and *upper\_bound(a)* provides iterators to lower and upper bound values of  $a$  stored in the set, according to a set instance specific order;
- *find(t)* obtains an iterator to the tuple  $t$  in the set if present, and
- *empty()* determines whether the set is empty.

This interface mimics the STL set-container interface so that existing data-structure can be adapted quickly to a DER data-structure. The interface operations cover the functionality of storing tuples of fixed-length, set-enforcement, iterators for whole set traversals, range queries,

membership/existence tests (find), and correct parallel execution. Besides the set-container functionality, a DER data structure requires a high-performance implementation of these operations. The DER operations are to be performed with little or no synchronization overhead on shared-memory multicore computers. This can only be achieved by specializing a generic set-container for this purpose. In the following, we discuss the necessary specializations.

### 3.2 | Specializing for parallel tuple insertion

In this section, we describe how a data structure can be specialized for tuple insertion, ie, the  $insert(t)$  operation. Note that a tuple insertion already combines two primitive operations of DER: (1) the existence check whether the tuple already exists in the relation and (2) the actual insertion of the tuple. The existence check is necessary for enforcing that the container remains a set, ie, without the existence check, a tuple may be introduced more than once into the set.

The combination of the two primitive operations imposes already a synchronization challenge for a DER data structure since, when a large number of multiple insertions are performed in parallel, a race condition may be imposed between the parallel insertions. One of the consequences of the data race will be the violation of the set property if the insertions are not synchronized properly. To obtain satisfactory performance and scalability, synchronization techniques with low synchronization overheads are essential.

Due to the existence check for set enforcement, classical linear data structures such as linked lists and arrays are ill-suited for the purpose of DER. The main reason is that the tuple existence check would exhibit a worst-case runtime complexity of  $O(n)$ , where  $n$  is the number of tuples. Variations of balanced search trees such as B-Trees<sup>18</sup> perform the existence in  $O(\log n)$ , and the variations of tries<sup>15</sup> perform the existence check in  $O(W)$ , where  $W$  is the (constant) bit length of a tuple.

Another observation in the seminaïve evaluation is that a DER data structure has *two distinct phases*. In the first phase, the data structure is only read, and hence, no synchronization is required between multiple threads accessing the data structure simultaneously. In the second phase, parallel multiple insertions occur without reading/traversing the tuples of a relation. Between the two phases, the evaluation sets in place barriers so that no interleavings between tuple insertion operations and read-operations can happen. Naturally, these two phases can be exploited by a DER data structure, specializing the data structure purely for the purpose of Datalog evaluation.

Specifically, we have four auxiliary relations capturing the new, current, previous, and delta knowledge of a recursive relation.<sup>6</sup> For evaluating a rule for a recursive relation, the auxiliary relations current, previous, and delta knowledge are only read but not written to, ie, they are used to iterate tuples over relations / and existence checks only. The auxiliary relation for new knowledge is the only one which receives tuples in a loop-nest and is not read from. For example, in our motivating example in Figure 1, the relation  $newTc$  is only written to inside the loop (cf. line 9 to line 18) but never read inside the loop-nest that evaluates new tuples for the rule  $tc(X, Z) :- tc(X, Y), connection(Y, Z)$ .

As a consequence, we can specialize a DER data structure for synchronization purposes, ie, we do not require any synchronization between read and write operations. Another consequence is that the order of the write operations is not important as long as all tuples that are inserted in the set will materialize after the write phase. Other optimizations that are based on these two distinct phases are also possible. They are specific to a concrete DER data structure, as outlined in Sections 4 and 5. For example, we exploit this assumption to propose a locking mechanism that avoids the overhead of a general purpose B-Tree for insertion.

### 3.3 | Specializing for nested-loop joins

In this section, we describe how the order of attributes that tuples are stored in, eg, a B-Tree can be specialized in the data structure to facilitate low overhead range lookups, ie,  $lower\_bound(a)$  and  $upper\_bound(a)$  operators. These range lookups are the building blocks of evaluating rules.

To achieve peak performance, modern Datalog engines generally assume an execution model that keeps all relations in memory.<sup>16</sup> In this sense, every relation is indexed. A Datalog engine requires that for all operations on a relation, the operations can be performed on an index to achieve improved performance. Hence, several replications of a relation are stored in memory that have different attribute orders. This *cluster of indexes* of a relation facilitates fast searches filtering tuples of relations. Such filtering operations are called *primitive searches*.

A primitive search can be naively achieved by conducting a linear scan over all tuples in a relation and checking the search predicate against each tuple. However, the time complexity of linear scans with  $n$  tuples is  $\mathcal{O}(n)$ , which is too costly for large relations considering that each primitive search is invoked repeatedly many times.

To enable faster searches on a relation, we require that DER data structures have an order among tuples. Given an ordered set of tuples, tuple lookups can be performed efficiently depending on the specifics of the data structure. For example, in some variations of a balanced-search tree, tuples can be found in logarithmic time rather than in linear time.

Since a tuple may have several elements, an order of tuples is imposed by element-wise comparison using a sequence over all attributes of the relation; this comparison is known as a *lexicographical order*. Let us denote an attribute sequence by  $\ell = x_1 < x_2 < \dots < x_m$ , where  $<$  denotes a chaining of elements to form a sequence.

Then, given  $\ell$  that is formed by all attributes of a relation, a lexicographical order  $\sqsubseteq_\ell \mathcal{D} \times \mathcal{D}$  is a total order (ie, *reflexive, asymmetric, transitive*) defined over the domain  $\mathcal{D}$  of the relation with respect to  $\ell$ . For two tuples  $a, b \in \mathcal{D}$ , when  $(a, b) \in \sqsubseteq_\ell \mathcal{D} \times \mathcal{D}$ , we write  $a \sqsubseteq_\ell b$ , and we say that  $a$  is smaller than  $b$  with respect to  $\ell$ . Note that  $a \sqsubseteq_\ell a$ , and for any two different tuples  $a, b \in \mathcal{D}$ , we either have  $a \sqsubseteq_\ell b$  or  $b \sqsubseteq_\ell a$  but not both.

A DER data structure has a single order that can support several primitive searches. However, not all primitive searches can be covered by a single index. Thus, the relation is replicated with different orderings. An evaluation heuristic such as literal orderings in rules impacts which data orderings are most memory efficient in the replicated relations. Note that different attribute sequences usually result in different lexicographical orders and thus different indexes. That is, for tuples  $a, b \in \mathcal{D}$  and attribute sequences  $\ell$  and  $\ell'$ , it is possible that  $a \sqsubseteq_{\ell} b$  and  $b \sqsubseteq_{\ell'} a$ .

Both data structures presented in Sections 4 and 5 support arbitrary attribute orderings and thus allow for modern Datalog engines to improve tuple search time complexity while reducing memory overhead using index covering techniques.<sup>16</sup>

### 3.4 | Specializing for sequences of operations

Another observation is that a DER data structure may be traversed consecutively by multiple threads to perform a single operation each. Although some traversals are log-linear time, the traversal will become expensive for a large number of operations that are in the billions. To amortize the cost of traversal, DER data structures may provide a hint mechanism that caches previous traversals. The generic hint mechanism uses state associated with a DER storing summaries of previous traversals. In case of a new traversal, the state is revisited first and checked whether the old traversals can be reused. If this is the case, the information from the previous operations reduces the execution time of the current operation. Especially, when accessed data elements are approximately sorted, the hint mechanism becomes more effective for ordered data structures such as B-Trees. The hint mechanism can be used for a wide range of operations such as tuple insertion, range queries, and existence checks.

#### 3.4.1 | Interface specifications

To describe interfaces of abstract data types like our DER, a type-based notation can be utilized to provide a more formal description. Let  $D$  be the set of all instances of a DER, and the set  $E$  denotes the set of all its potential elements. Then,

$$\text{contains} : (D, E) \rightarrow \mathbb{B}$$

and

$$\text{insert} : (D, E) \rightarrow D$$

are type specifications for the required operators. The *contains* function takes a pair of a data structure instance  $d \in D$  and an element  $e \in E$  and determines whether  $e$  is present in  $d$  by producing a boolean value  $\text{contains}(d, e) \in \{\text{true}, \text{false}\} = \mathbb{B}$ . The *insert* function, on the other hand, accepts a data structure instance  $d \in D$  and an element  $e \in E$  and computes a potentially new instance  $\text{insert}(d, e) = d_r \in D$  such that the predicate

$$\text{contains}(\text{insert}(d, e), e)$$

holds. Thus, the data structure instance  $d_r$  obtained by processing the insertion function contains the inserted element  $e$ .

While this formalism is based on a functional perspective of the data structure interface, it serves as a reasonable abstraction of operations on data structures maintained within random access machines (RAM) underlying imperative or object-oriented programming languages. Modeling operations following such paradigms would require a more elaborate set of definitions regarding memory state and side effects of operations. By utilizing the functional point of view, such details, not required for our purposes, can be omitted.

Nevertheless, it is important to see that in nonfunctional context, implementations of the *insert* operation above are typically destructive. Thus, the data structure passed as an argument is mutated over the course of function evaluation, and a modified version is returned. The original data structure and the resulting data structure thus typically occupy mostly overlapping regions of the main memory.

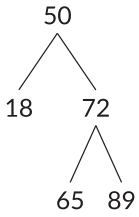
Each operation on a data structure has a given runtime complexity, parameterized by the size of the data structure. For instance, the insertion of an element into a balanced tree requires  $\mathcal{O}(\log(N))$  steps, where  $N$  is the number of elements stored in the tree.

To perform sequences of operations, like the insertion of a range of elements into a given tree, those operations are chained. Thus, to insert the elements  $e_1, \dots, e_n \in E$  into a preexisting data structure  $d_0 \in D$ , the operations

$$\begin{aligned} d_1 &:= \text{insert}(d_0, e_1) \\ d_2 &:= \text{insert}(d_1, e_2) \\ &\dots \\ d_n &:= \text{insert}(d_{n-1}, e_n) \end{aligned}$$

need to be evaluated. If the insertion operation is  $\mathcal{O}(\log N)$ , the overall sequence of operations will cost  $\mathcal{O}(n \log N)$ . Thus, the overall cost of the sequence is the sum of its individual steps.

However, considering the concrete example of a balanced tree, performing sequences of single-element insertion operations could be conducted more efficiently by sharing information among insertions.



**FIGURE 2** Binary search tree example

For instance, consider the insertion of the value 32 into the binary search tree outlined in Figure 2. The new value is to be inserted as the right child of the node 18. At this point, it is known that the subtree to be rooted by the new value 32 is to contain all values between 18 and 50. Thus, if the next value to be inserted is in this range, navigation through the tree can be omitted, and the new element can be directly inserted as a child node of 32. Under these circumstances, the costs for inserting the second element is reduced from  $\mathcal{O}(\log N)$  to  $\mathcal{O}(1)$ . To forward this kind of information from one insertion operation to the next, we introduce operator hints.

Let  $H$  be the set of all those hints. Then, an insertion operation utilizing hints is specified by

$$\text{insert}_H : (D, E, H) \rightarrow (D, H).$$

An insertion operation thus takes a data structure, the element to be inserted, and some operation hints as an input, and produces an updated pair of a data structure and a hint for subsequent operations. Analogous, the *contains* operation can be modified to include operator hints such that

$$\text{contains}_H : (D, E, H) \rightarrow (\mathbb{B}, H).$$

The steps to perform sequences of element insertions get modified to

$$\begin{aligned} (d_1, h_1) &:= \text{insert}(d_0, e_1, h_0) \\ (d_2, h_2) &:= \text{insert}(d_1, e_2, h_1) \\ &\dots \\ (d_n, h_n) &:= \text{insert}(d_{n-1}, e_n, h_{n-1}) \end{aligned}$$

using some initial operator hint  $h_0 \in H$ . Thereby, the overall runtime complexity can be reduced. In the binary search tree example outlined above, inserting a sequence of  $n$  consecutive elements into a gap of the initial tree reduced from  $\mathcal{O}(n \log N)$  to  $\mathcal{O}(\log N + n)$ .

Naturally, one could argue to integrate hints into data structures directly, thereby eliminating the need for the extra parameter. However, keeping them apart allows client code to manage the context of insertion operations. For instance, assume two lists of consecutive elements  $a_1, \dots, a_n \in E$  and  $b_1, \dots, b_n \in E$  have to be inserted into a binary search tree in an interleaved fashion. By maintaining operator hints within the data structure, the proximity of the elements ending up in the tree is lost, and the beneficial effect is not materializing. However, by maintaining two separate operator hints external to the data structure, those sequences can be inserted using their respective hint values, thus benefiting from this mechanism. While seemingly being a contrived example, this is precisely the situation to be faced in a concurrent environment where those two sequences are to be inserted by two distinct threads into a shared data structure.

In a concurrent environment, multiple threads might attempt to insert sequences of elements into a shared data structure. Allowing externally managed operator hints does not only enable those threads to benefit from the proximity of their processed elements but also to maintain the actual operator hints within thread-private storage. Thus, updates to operator hints do not need to be synchronized as it would be the case if hints would be managed internally.

### 3.4.2 | Data structure perspective

Operator hints effectively facilitate the realization of caches for partial results of former operator evaluations. For instance, in the binary search tree insertion case, the pointer to node 32 rooting the subtree covering the range  $(18, \dots, 50)$  can be maintained as a hint, thereby avoiding tree-navigation costs in future insert operations. These partial results of operations may be retrieved through caches faster than they might be recomputed using the management information included in the data structure.

Since hints thus exhibit the characteristics of caches, their design needs to scope with a range of questions: what information should be cached, what key is utilized to address the cached information, how is cached information invalidated, how much information is cached without losing its performance benefit, and what replacement policy for cache entries shall be used.

In the binary search tree example above, we outlined the utilization of a single cache entry formed by a node pointer addressing the most recently introduced tree node. We utilized the boundaries of its covered subtree to test the relevance of this information for future inserts. Fortunately, since simple binary search trees are never restructured, this information is never invalidated by future inserts; hence, no invalidation



has to be realized. Finally, since we only store a pointer to the last inserted node, we utilize a single-slot cache with least-recently-used (LRU) replacement policy.

Naturally, all those design decisions can be explored. In particular, by adding additional slots, more intricate patterns of insertions could benefit from operator hints. While a single hint entry improves the performance of one sequence of consecutive elements, caching two hints using an LRU policy will lead to  $\mathcal{O}(\log(N) + n)$  insertion complexity also in the case of two interleaved sequences.

## 4 | OPTIMISTIC B-TREE

B-trees<sup>19</sup> are among the most widely utilized data structures for maintaining large quantities of keys. Their sorted and balanced nature provides  $\mathcal{O}(\log n)$  worst-case execution time for insert and lookup operations,  $\mathcal{O}(n)$  traversals, as well as  $\mathcal{O}(\log n + r)$  complexity for range queries, where  $r$  is the size of the range. The mostly consecutive storage of keys in inner nodes and leaves further facilitates efficient utilization of memory caches due to a high degree of spatial locality of consecutive accesses. Like other data structures, B-trees require  $\mathcal{O}(n)$  space. However, by blocking the storage of keys in nodes, the overhead of management information like pointers is effectively amortized, even for small key types.

While providing outstanding performance in many use cases demanding sequential accesses, state-of-the-art concurrent implementations of B-trees cannot cope well with insert-intensive workloads as imposed by parallel Datalog evaluation. In particular, within parallel sections, all threads are continuously attempting to insert data into a single shared B-tree instance, constituting a near worst-case scenario for typical locking schemes. Maintaining the consistency of the internal data organization requires the utilization of sophisticated locking mechanisms and schemes, typically limiting parallel scalability. Frequent lock state and data updates need to be propagated among CPU cores sharing access to B-trees, causing high communication demands on interconnects, easily leading to resource contention.

While imposing challenges on the scalability of insert operations, the Datalog use case also exhibits properties leading to relaxation of synchronization requirements. In particular, the Datalog evaluation scheme guarantees that no data structure is ever read and written concurrently. Membership tests, range queries, and traversals are never processed while performing insertions on the same tree. Thus, these operations do not need to be shielded from each other.

The remainder of this section outlines a concurrent B-tree design providing highly scalable insertion performance on multicore and multsocket architectures. It is based on a new optimistic fine-grained locking scheme using *optimistic read-write locks*.<sup>†</sup> By avoiding memory modifications in the most frequent read-only case, even when acquiring and releasing locks, the utilization of intrasocket and intersocket communication channels is vastly reduced. Furthermore, by caching leaf-nodes reached by insertion and lookup operations within thread-private variables, a hint mechanism shortcutting consecutive tree navigation steps has been integrated to exploit temporal locality of leaf-node accesses.

### 4.1 | Synchronization

The two phases of the parallel seminaïve evaluation of Datalog<sup>17</sup> guarantee the strict separation of concurrent read and write accesses to shared data structures. Data structures are either (concurrently) read by multiple threads or written to. Thus, since lookups, emptiness checks, and range queries do not modify any shared state of a B-tree, only multiple concurrent executions of the B-tree's insertion operation need to be synchronized to isolate their effects from each other.

Previous work has been exploring a wide variety of locking schemes for B-trees.<sup>12</sup> Designs range from global locks, over per-node fine-grained locking using mutex or read/write locks, to implementations utilizing the capabilities of hardware transactional memory.<sup>20</sup> Many of those are limited in their scalability due to high bandwidth requirements caused by frequent cache line invalidations when updating lock state information. Especially, the lock protecting the root node, which is involved at the beginning of every insert operation, causes a high amount of communication, resulting in a significant performance penalty, and thus in severely limited scalability.

#### 4.1.1 | Optimistic read-write lock

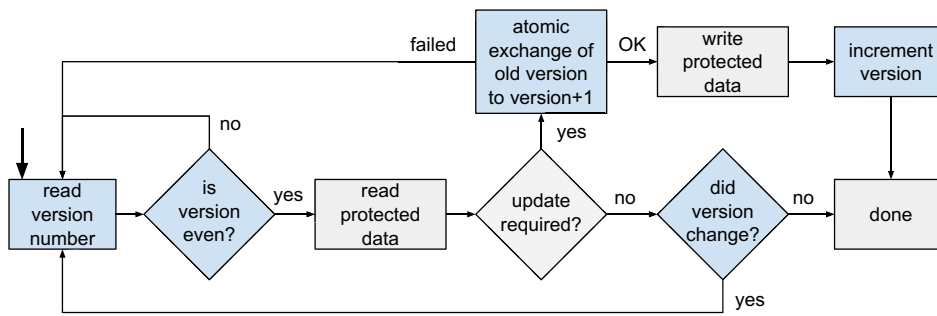
The *seqlock*<sup>21</sup> design constitutes a highly scalable lock solution for such frequent read / seldom write use case. Its design is based on the optimistic assumption that conflicts are very unlikely to occur. Under these circumstances, it can be cheaper to focus on reactive conflict detection than conventional proactive conflict prevention, an idea also previously exploited in database applications.<sup>22</sup>

The central idea is to associate a *sequence number* to each shared blocks of data like, eg, a B-tree node. The sequence number covers two pieces of information: a version number incremented whenever the data is updated and one bit to indicate whether an update is currently in progress. The progress bit is typically the last bit of the sequence number, turning even numbers into indications of stable states, without active writers, and odd numbers into indicators of ongoing updates.

Before accessing shared data, the sequence number is retrieved. If it is odd, a write operation is in progress, and the accessing thread (busy-) waits for its completion. Once even, the sequence number is recorded, and the read operation on the shared data performed. To complete the read, the sequence number after reading the data is read again and compared to the recorded value. If identical, no concurrent update has occurred. If different, a conflict has been detected, and the read operation may have to be repeated.

<sup>†</sup> Hence, the name optimistic B-tree.





**FIGURE 3** Outline of optimistic locking scheme

Write operations start by attempting to atomically increase an even sequence number to an odd number. On success, exclusive write access is guaranteed though concurrent read accesses may be present. The writing thread may then continue by modifying the shared data before finally increasing the now odd sequence number to the next even number. This unblocks new read or write operations and allows ongoing read operations to detect conflicts.

However, during B-tree insert operations, threads are not clearly separated into reading or writing threads. Each thread passing through a node is a reading thread but may later be mutated to a writing thread in cases a key needs to be inserted in a leaf node or an inner node gets split. Thus, the basic seqlock scheme needed to be customized to what we refer to as an *optimistic read-write* lock. Each thread starts as a reading thread but may optionally upgrade its read permission to write permission, without invalidating read data. Figure 3 outlines the state diagram of every thread accessing a protected region using our optimistic read-write locks.

Each access follows the idea of a *transaction*, with a clearly indicated start and end. Each transaction starts by reading data, followed by an optional write phase, and consistency check at the end. Upon failure, the transaction is restarted.

To that end, the optimistic lock offers several operations. For read accesses, *start\_read*, *valid*, and *end\_read* are provided. The *valid* operation may be used during ongoing transactions to validate data read so far. The operations *try\_upgrade\_to\_write*, *end\_write*, and *abort\_write* are offered to control the start and termination of the optional write phase. Write phases may be aborted in case no modification is performed although write access has been acquired. In this case, the sequence number is reverted back to its initial state (decremented) to avoid invalidating ongoing concurrent reads. Furthermore, the nonblocking operation *try\_start\_write* and the blocking operation *start\_write* are offered to skip the initial read phase and enter directly into a write phase.

#### 4.1.2 | Optimistic B-tree insertion

In our optimistic B-Tree, each node is equipped with an optimistic read-write lock to protect its internal state. Furthermore, the root node pointer is as well protected by an additional *root\_lock*. The orchestration of these locks to facilitate a properly synchronized key insertion procedure is outlined in Algorithm 1.

Lines 2-9 cover the thread safe creation of a root node, only triggered when attempting to insert into an empty tree. In this case, an initial root-node is created by the first thread to pick up the root lock. After ensuring that at least one node is present, lines 11-49 cover the actual insertion procedure for inserting keys in non-empty trees.

Lines 13-17 are tasked to obtain a valid reference to the current root node, after which lines 20-33 navigate down the tree, chaining locks. When reaching a leaf node, the new key is inserted (lines 35-47). To that end, read permissions are upgraded to write permissions (line 36), and the new key is inserted (line 46) if there is sufficient space. If, however, the targeted leaf node is full, it needs to be split (line 40), and the insertion is restarted at the top. The splitting procedure is summarized in Algorithm 2.

Before performing modifications on the tree, the split procedure acquires write permissions on all nodes to be modified, bottom-up. The involved nodes comprise all nodes on the path from the targeted leaf node back up to the root node up and including the first nonfull node. Lines 2-23 perform the corresponding navigation and locking steps. Once permissions have been acquired, nodes are split (line 26), and locks are released in reverse order (lines 28-35). As typical for B-trees, by splitting nodes, the lower half of the keys are retained by the split node, while the upper half is moved to a new node, which is registered in the parent-nodes next-pointer list, referenced by the median key value.

The following invariants are obeyed at all time: almost all data stored in a node, including keys in leaf nodes as well as keys and pointers in inner nodes, are protected by the enclosing nodes' read/write lock. The only exception is given by a nodes parent pointer, which is protected by the parent node's lock or, in case of the root node, by the root pointer lock. Whenever reading any of those values, a read permission is acquired from its protecting lock, and after the value is obtained, its validity is verified before consuming the value by, eg, dereferencing a node pointer (eg, line 28 in Algorithm 1). Analogous, modifications are encapsulated by the acquisition and release of write permissions. In all cases, detected conflicts are handled by restarting the overall insertion algorithm.

The outlined locking scheme detects read/write conflicts but sustains from spending time (and memory accesses) on preventing those. It is thus optimistic. However, conflicts may occur, and invalid values may be read from shared memory locations. It is thus essential to validate version numbers before using any obtained data. In particular, dereferencing corrupted values extracted for child-node pointers would cause segfaults in a best case scenario, and invalid, next to impossible to detect subtle memory corruptions in the worst case.

**Algorithm 1** Optimistic B-tree insertion procedure

---

```

1: procedure insert (tree, val)
2:   // safely initialize root node pointer
3:   while tree->root == null do
4:     if !(try_start_write(tree->root_lock)) continue
5:     if tree->root == null then
6:       tree->root ← < create new node >
7:     end if
8:     end_write(tree->root_lock)
9:   end while
10:
11: restart:
12:   // safely obtain root node and its lock
13:   repeat
14:     root_lease ← start_read(tree->root_lock)
15:     cur ← tree->root
16:     cur_lease ← start_read(cur->lock)
17:   until end_read(root_lease)
18:
19:   // descent into the tree
20:   while true do
21:     // if value to be inserted is present => done
22:     if contains(cur, val) and valid(cur_lease) return
23:
24:     // process inner node
25:     if cur->inner then
26:       next ← find_next(cur, val)
27:       if !valid(cur_lease) goto restart
28:       next_lease ← start_read(next->lock)
29:       if !valid(cur_lease) goto restart
30:       cur ← next
31:       cur_lease ← next_lease
32:     continue
33:   end if
34:
35:   // request write access to located leaf node
36:   if !try_upgrade_to_write(cur_lease) goto restart
37:
38:   // make some space, if necessary
39:   if full(cur) then
40:     SPLIT(tree, cur)
41:     end_write(cur_lock)
42:     goto restart
43:   end if
44:
45:   // insert value into this leaf node
46:   < insert value in current node >
47:   end_write(cur_lock)
48:   return
49: end while
50: end procedure

```

---

Since exclusive write-permissions are only acquired bottom up, and held read-permissions are not preventing other thread from acquiring write permissions, deadlocks are effectively avoided since no dependency cycles can be formed this way. However, rare lock contention situations may occur for inner nodes forming the root of multiple paths to leaf nodes only comprising fully occupied nodes.

The advantage of the presented optimistic locking scheme over similar fine-grained schemes using conventional locks is derived from optimizing the most frequent use case of a mere read of the state of an inner node. The navigation through an inner node only involves memory reads for handling lock operations, key search steps, and child-node pointer retrieval. None of these reads to shared memory cause any cache line invalidation and thus no bus communication. Using conventional locks, at least the lock state itself would have to be mutated twice, when acquiring the read permissions and when releasing it.

#### 4.1.3 | Optimistic B-tree implementation

Seqlocks impose a particular challenge to be implemented correctly using contemporary programming languages. Their basic idea is to permit parallel accesses and detect and compensate for conflicts in rare cases. However, by doing so, seqlocks depend on properly defined semantics for race conditions, a case frequently excluded from the definitions governing the memory model underlying programming languages. C++, in particular, states that any program with data races exhibits *undefined* behavior. Implementing seqlocks thus requires a careful encoding.

The problem of implementing seqlocks using modern C++ has been invested by Boehm in great detail.<sup>23</sup> The key idea is to exploit C++'s exception of data being accessed through atomic wrapper from any sort of race condition, by definition. Thus, as long as all shared data that is protected by seqlocks is accessed through atomic data type wrappers, a race condition-free valid seqlock implementation can be realized. In his

**Algorithm 2** Optimistic B-tree node splitting procedure

---

```

1: procedure split (tree,node)
2:   // write-lock path bottom-up
3:   cur ← node
4:   parent ← cur->parent
5:   path ← ε
6:   while true do
7:     if parent != null then
8:       while true do
9:         start_write(parent->lock)
10:        if parent == cur->parent break
11:        abort_write(parent->lock)
12:        parent ← cur->parent
13:      end while
14:    else
15:      start_write(tree->root_lock)
16:    end if
17:    path ← append(path,parent)
18:
19:    // stop at root or non-full inner node
20:    if parent == null or !full(parent) break
21:    cur ← parent
22:    parent ← cur->parent
23:  end while
24:
25:  // conduct actual split
26:  < split node and propagate to parents >
27:
28:  // unlock path top-down
29:  for all parent in reverse(path) do
30:    if parent != null then
31:      end_write(parent->lock)
32:    else
33:      end_write(tree->root_lock)
34:    end if
35:  end for
36: end procedure

```

---

paper, Boehm outlines several different refinements of this idea, minimizing the performance penalty introduced by those atomic wrappers. By carefully utilizing memory order arguments and synchronization fences, the native memory model of a targeted architecture can be exploited to the point where accesses through atomic wrappers yield (almost) identical sequences of assembly instructions as nonwrapped access, without sacrificing correctness.

For our implementation, we adopted a solution presented by Boehm. When entering a critical region, the sequence number is obtained using `memory_order_acquire`. Shared data is then read using `memory_order_relaxed`. Finally, to validating retrieved data, a `memory_fence_acquire` is included to force all those relaxed reads to be completed, before reading the lock's sequence number using `memory_order_relaxed` to complete the validation.

## 4.2 | Operation hints

The data produced during the evaluation of Datalog queries using nested loop joins<sup>6</sup> may, in certain occasions, trigger membership queries, and inserts following the same order as data are stored in tables. Thus, consecutive operations may require to navigate to consecutive keys stored within the same leaf node of the B-tree. This spatial locality can be exploited to reduce the number of times the B-tree has to be navigated top-to-bottom by consecutive calls. This is realized through the integration of an operation hint mechanism into our B-tree design.

Essentially, each thread maintains pointers to the last leaf nodes reached in a B-tree instance when performing an insertion operation or a membership test. Both hints are maintained separately within a thread-private variable, allowing effects to be bound to the operation type, tree instance, and thread context. When performing inserts (or membership tests), those hints referencing leaf nodes are tested first. If they reference a node in which the inserted key is to be located, the navigation from the root node to the leaf can be omitted. In case of a membership test, the leaf referenced by the matching hint can be inspected directly. In case of an insertion operation, keys can be directly inserted in the leaf node. In either case, if the provided operator hints are invalid for the operation at hand, they are ignored, and the default membership test or insertion procedure is executed.

For example, when inserting the key (3, 4) into a B-tree using standard lexicographical ordering, it may end up in a node storing keys ranging from (2, 9) to (4, 1). By maintaining the reference to this node, a subsequent insertion of, eg, the key (3, 6) can be immediately performed on the node located by the previous insertion since the key is as well contained in the range of elements covered by this leaf node. However, if this insertion is followed by inserting an out-of-range key like (5, 2), the maintained hint is ignored, and a regular insertion starting at the B-tree's root node is initiated. Hints for lookup operations operate analogously.

By storing operator hints in thread-private memory locations, typically maintained on the processing thread's stack, and passing them along as extra arguments to insert and lookup operations, thread-safe hint manipulations are automatically derived. Since insertion hints may cause the navigation-down phase to be skipped, no read-locks along the path from the root-node to the leaf node are obtained. These locks are, in the regular version of the insertion procedure as outlined in Algorithm 1, required to ensure that the correct leaf node has been reached. However, since from the boundaries of the leaf node in case of the utilization of a hint this property can be ensured, these missing read-locks can be tolerated. Nevertheless, write locks for splitting nodes need to be acquired bottom-up as usual. This bottom-up exclusive-lock-acquiring characteristic of our locking scheme thus facilitates the efficient utilization of insertion hints even in concurrent access contexts.

Finally, our C++ implementation covers various additional minor optimizations to improve its performance in the Datalog context: (1) most recursive operations are replaced by iterative versions to avoid excessive function-call overheads, (2) the involvement of node-rebalancing operations to reduce the number of necessary node splits and increase memory efficiency, (3) a linear B-tree merge algorithm, and (4) the utilization of customized three-way comparison functions for B-tree entries. All these optimizations cover implementation details tuned to improve real-world execution performance. However, they do not contribute new algorithmic insights beyond what has been presented above. The implementation of our B-tree is open-source and forms one of the foundations of the Soufflé Datalog engine.<sup>24</sup> The implementation can be found in the file `BTree.h` and is licensed under UPL V1.0.

## 5 | BRIE

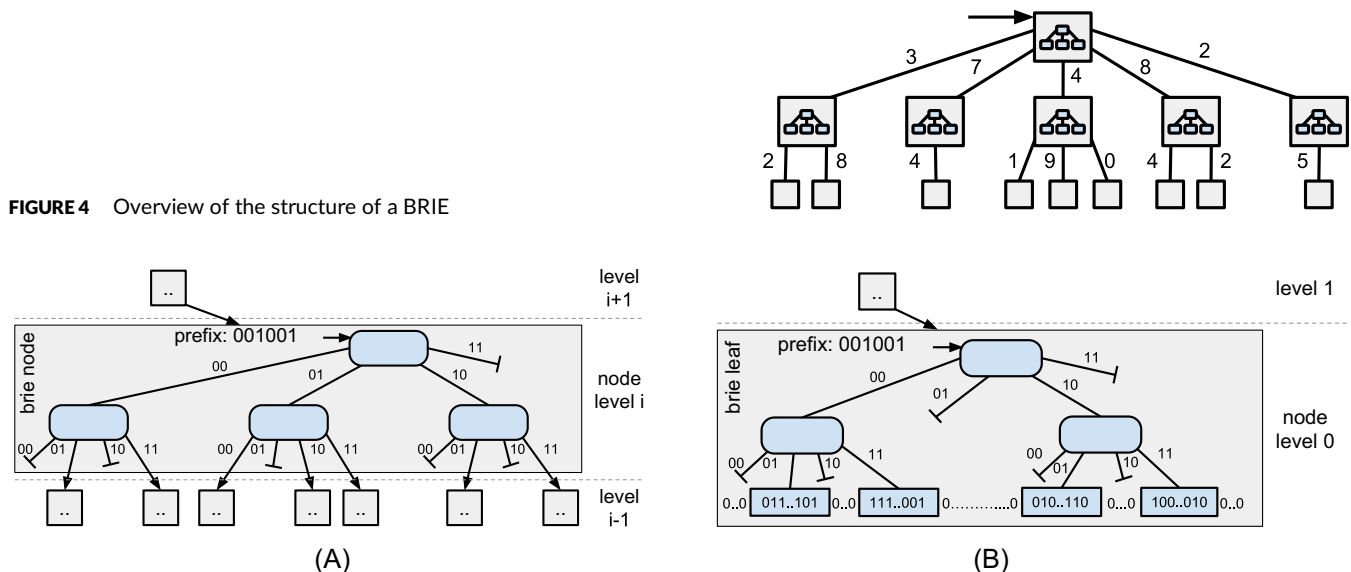
The optimistic B-tree is specialized for the concurrent access scheme exhibited by the Datalog evaluation scheme. Another approach is to specialize a data structure by considering the properties of the expected data to be maintained. For example, in the case of points-to analysis the data inserted in the most performance relevant relations are often *dense*, meaning that within the boundaries of a covered ( $n$ -dimensional) range, a high number of elements will be present. This property can be utilized for designing a more efficient  $n$ -dimensional point data structure.

In this section, we introduce the *Brie*. The BRIE is structurally based on a trie and uses a node-blocking scheme found internally in B-trees. The trie-based structure eliminates the need for tree reshuffling and facilitates lock-free insertion operations, while the node-blocking scheme leads to an efficient cache friendly consecutive storage of elements. Combined, these BRIE traits enable improved performance on dense data sets.

In Figure 4, the overall structure of a BRIE is shown. The foundation of the BRIE design is a fixed-height trie. Unlike B-trees or classic binary search trees, the BRIE encodes entries through its edges instead of explicit values in nodes. Each edge is labeled by an integer, forming a component of a stored tuple. Technically, this is realized by maintaining a map between labels and child node pointers. This mapping is based on quadrees extended by the node blocking feature of B-trees. Moreover, bit-masks are used on the leaf node level to compress the information of the presence or absence of values into a single bit.

Figure 5A depicts the inner structure of an inner node of a BRIE. The mapping between tuple component values and child node pointers is realized by a blocked quadtree. The binary encoding of the map's keys is used for navigating the quadtree. On the leaf level of the inner node's quadtree, pointers to the next level of the overall BRIE are stored.

To improve cache efficiency, multiple levels of the inner quadtree are collapsed into single nodes. Figure 5A illustrates this collapse for the two-level case. However, we note that the actual implementation collapses a configurable number of levels into a single node. For the experiments in Section 6, an empirically determined optimal value of 6 levels is selected.



**FIGURE 5** A, Internal structure of an inner node of a BRIE; B, Internal structure of a leaf node of a BRIE

Another optimization is performed by pruning branches that do not contain any child node pointers. Moreover, the root is formed by the first node with more than a single child. In this way, a chain of nodes representing a common prefix of all keys at the top of the quadtree is avoided. Instead, the common prefix and its length is stored in an extra variable. This prefix is shortened as necessary during insert operations.

The structure of leaf level nodes in the BRIE is shown in Figure 5B. While conceptually identical to inner nodes, the child node pointers are replaced by machine word-sized bitmaps indicating the presence or absence of entries. Thus, effectively, the last  $\log_2(W)$  levels of the index structure are collapsed into a single machine word  $W$  bits wide.

The insert operation for the bribe is shown in Algorithm 3. The insert is synchronized using lock-free techniques that are based on atomic operators. These techniques facilitate the necessary means to set individual bits in the bitmaps maintained by leaf nodes and to atomically compare-and-swap pointers in inner node levels.

---

**Algorithm 3** Brie insertion procedure (simplified)

---

```

1: // node handling utility
2: procedure ensureExistence<T>(ptr)
3:   // handle none-existing root
4:   if ptr == nullptr then
5:     new ← new T()
6:     if !CAS(ptr, nullptr, new) then
7:       // somebody else was faster
8:       delete new
9:     end if
10:   end if
11:   return ptr
12: end procedure
13:
14: // inner-node tree navigation utility
15: procedure getLeaf (root, val)
16:   // handle potentially none-existing root
17:   ensureExistence<IndexNode> (root)
18:
19:   // adapt prefix if necessary
20:   ensurePrefixCovered (root, val)
21:
22:   // navigate down the internal quadtree
23:   info = atomic_load (root->prefixInfo)
24:   cur ← info.root
25:   level ← info.level
26:   while level >= 0 do
27:     next ← cur->next[(val >> level) & 0x3F]
28:     cur ← ensureExistence<IndexNode> (next)
29:     level ← level - 6
30:   end while
31:   return cur->value
32: end procedure
33:
34: // entry point for insertion of value
35: procedure insert (node, val, level)
36:   // handle leaf bribe node
37:   if level == 0 then
38:     mask ← GETLEAF(node->index, val[0] >> 6)
39:     atomic_or (mask, 1 << (val[0] & 0x3F))
40:     return
41:   end if
42:
43:   // navigate through inner bribe node
44:   next ← getLeaf (node->index, val[level])
45:
46:   // if next does not exist, create it
47:   ensureExistence<BrieNode> (next)
48:
49:   // insert recursively
50:   insert (next, val, level-1)
51: end procedure

```

---

The *insert* operation adds a new element to a given BRIE by utilizing several auxiliary functions. The function *getLeaf* performs the node-internal quadtree navigation and manipulation within nodes. Furthermore, the generic function *ensureExistence* is used to create new nodes on demand while navigating the node-internal quadtree (lines 17 and 28) or the overall trie based metastructure (line 47). Finally, the *ensurePrefixCovered* function adapts the common prefixes stored as part of the node-internal quadtrees and grows those trees bottom-up whenever needed (line 20). Algorithm 4 outlines the details of this third function.

Whenever a new element is inserted into an inner quadtree, its key value may exceed the coverage of the currently maintained common prefix value (line 36). To include the new key, the number of leading significant bits of the common prefix is gradually reduced until the new key

**Algorithm 4** Brie inner node shared prefix utilities

---

```

1: // a utility to check whether a given prefix is sufficient
2: procedure covers (info, val)
3:   i ← atomic_load (info)
4:   if i.root == nullptr return false
5:   mask ← (-1) << i.level
6:   return val & mask == i.prefix
7: end procedure
8:
9: // a utility to shorten the prefix by 6 bits
10: procedure raiseLevel (info, val)
11:   // get current state (atomic)
12:   ol ← atomic_load (info)
13:
14:   // create new prefix info struct (on stack)
15:   nl ← PrefixInfo()
16:   nl.level ← ol.level + 6
17:   nl.root ← new Node()
18:   nl.root->next[(ol.prefix >> ol.level) & 0x3F] ← ol.root
19:
20:   // initialize or update prefix
21:   if ol.root == nullptr then
22:     nl.prefix ← val & ((-1) << 6)
23:   else
24:     nl.prefix ← ol.prefix & ((-1) << nl.level)
25:   end if
26:
27:   // compare and swap prefix index struct (atomic)
28:   if !CAS(info, ol, nl) then
29:     // somebody else was faster
30:     delete nl.root
31:   end if
32: end procedure
33:
34: // adapts the shared prefix
35: procedure ensurePrefixCovered (node, val)
36:   while !covers (node->prefixInfo, val) do
37:     raiseLevel (node->prefixInfo, val)
38:   end while
39: end procedure

```

---

is covered as well. Simultaneously, in each step, a new root node is added. The function *raiseLevel* covers the synchronized introduction of a new root node and the corresponding update of management information. The latter comprises the root node pointer, the level counter, and the common prefix. Algorithm 4 uses a hard-coded increment of 6 bits per level, thus the utilization of  $2^6 = 64$  child node pointers in each quadtree node. Overall, this results in inner quadtrees to facilitate high-output degree nodes grow bottom-up similar to a B-tree. However, no reshuffling of data nor node splitting is involved. Pointers once set within the BRIE are never updated or moved.

A synchronized implementation of the insert function has to protect three critical regions: (1) the insertion of a new BRIE node, (2) the setting of bits in leaf nodes, and (3) the update of the common prefix and thus the raising of the height of nested quadtrees.

Line 6 of Algorithm 1 realizes the synchronized addition of BRIE nodes. In this line, a pointer to a newly created subtree is replacing a null-pointer marking an empty subtree. The exchange is conducted using a compare-and-swap (CAS) operation. This atomic operation ensures that only one thread inserts a new subtree ever successfully at any given position in the tree. Thus, in case multiple threads attempt to insert newly created subtrees at the same position simultaneously, only one of those will succeed. The remaining threads will detect the collision, discard the nonaccepted temporary subtrees, and continue using the winning instance.

The second operation to be synchronized to handle concurrent modifications is in line 39 of Algorithm 3. It synchronizes threads concurrently attempting to set bits in leaf-node level bit masks. Here, we mark the presence of tuples in leaf nodes of the BRIE. The utilization of an atomic *boolean\_or* operator is sufficient to ensure correct updates. Furthermore, update conflicts do not have to be detected since their effects are implicitly aggregated.

The third critical region, located in the function *raiseLevel* in Algorithm 4, handles the growing of the BRIE-node local quadtrees. Conceptually, its synchronization also follows the CAS approach. However, unlike the previous case handling the first critical region, the value of the memory location to be updated is not a simple scalar but a small struct maintaining the common prefix, its length, and the root node pointer. In line 12, this current state of this struct is retrieved atomically and used for computing a new prefix, prefix length, and root node. Once completed, it is compare-and-swapped in as a replacement in line 28. This update of the prefix information contains concurrent updates detection and the discarding of temporary values in case of conflicts. The support for CAS operations on wider data types comprising multiple fields as utilized in this example depends on widely available platform specific support. C11 and C++14, as well as many other languages, offer atomic utilities to access those features or functionally equivalent substitutes.

The updates (1) and (3) may induce expensive conflicts; however, the probability of such conflicts occurring is low. Typical work load patterns will result in threads inserting elements in distinct parts of a BRIE, thus avoiding any conflicts. Nonconflicting inserts have very limited overhead over sequential nonsynchronized versions. In essence, the overhead is given by the utilization of a small number of atomic CAS operations instead of ordinary assignments.

When comparing conventional B-trees to the BRIE, our BRIE offers a number of key benefits.

Firstly, insertion operations are computationally cheaper. Since the navigation steps in the BRIE do not depend on search operations over ranges of keys, in each step, the following node can be directly addressed. Moreover, insert operations do not require insertions into ordered arrays and the movement of lists of keys as required in B-tree implementations.

Secondly, in the case of memory usage, the sharing of common prefix on both the trie and quadtree structures, as well as the utilization of a single bit to mark the presence of an entry on the leaf node level contribute to a reduction in the amount of memory required to store a given set of entries. However, as we demonstrate in Section 6, the efficiency of this compression capability depends heavily on the correlation of entries. Therefore, if common prefixes are sparse, the memory utilization may be significantly higher compared to B-trees.

Lastly, with the presented synchronization scheme, synchronization is conceptually lock-free (assuming sufficient atomic operation support) and linearizable. The latter can be shown based on the sequentially consistent order of successful CAS operations.

## 6 | EVALUATION

Besides our B-tree implementation (denoted as *btree* or *optimistic btree*) and our BRIE implementation (denoted as *brie*), we include the following additional data structures in our evaluation. For the evaluation of memory consumption and sequential performance, we include the following.

- C++'s *std::set*, denoted as *STL rbtree*, as an example of a balanced tree-based in-memory data structure (red-black tree) satisfying all requirements stated for Datalog relations;
- C++'s hash-based *std::unordered\_set*, for clarity denoted as *STL hashset*, providing theoretically superior insertion and lookup performance of  $\mathcal{O}(1)$ , yet no efficient support for range queries;
- A state-of-the-art B-tree implementation provided by Google,<sup>25</sup> denoted as *google btree*, to evaluate the quality of our optimistic B-tree implementation;
- A sequential version of our B-tree, to evaluate the impact of our locking scheme on the performance of operators with and without operation hints, respectively, denoted as *seq btree* and *seq btree (n/h)*.

For the evaluation of parallel operations, we evaluate the following.

- The *concurrent\_unordered\_set* implementation of Intel's TBB library version 2017\_U7<sup>26</sup> denoted as *TBB hashset*, representing an industry standard, state-of-the-art concurrent hash-based set implementation;
- An implementation of our B-tree without operation hints denoted as *btree (n/h)*
- a parallel reduction based set implementation where insertions take place on thread private set instances, before being merged in a parallel reduction step; the implementation uses Google's B-tree with OpenMP's user-defined reduction operation support; it is denoted as *reduction btree*.

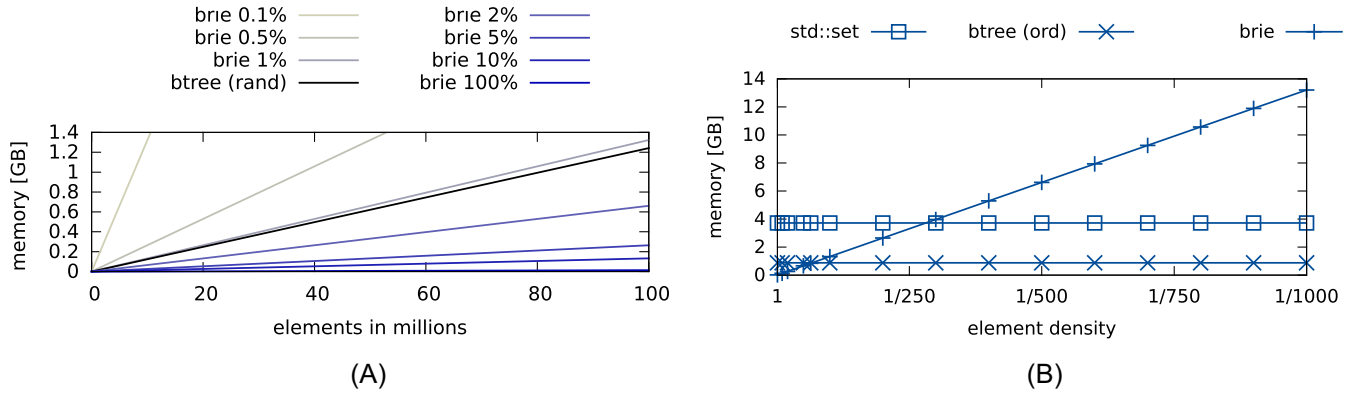
Table 1 provides a summary of all data structures included in our evaluation.

Most experiments presented in this section have been conducted on a four-socket Intel(R) Xeon(R) CPU E5-4650 system (8 cores each, 32 total) equipped with 256 GB memory using GCC 5.4.0 with -O3 optimization. For these multithreaded experiments, GCC's OpenMP implementation is used as the underlying parallel runtime system, threads are pinned to cores, and sockets are filled before threads are assigned to additional sockets. Single-threaded experiments have been conducted on an Intel Xeon Gold 6130 system equipped with 192 GB memory, using GCC 5.5.0 with -O3 optimization.

Designation	Thread safe	Description
<i>STL rbtree</i>	no	C++ standard library's set, implementing a red-black tree
<i>STL hashset</i>	no	C++ standard library's unordered set; a hash based set
<i>google btree</i>	no	Google's B-tree
<i>TBB hashset</i>	yes	Intel Threading Building Blocks' concurrent unordered set
<i>seq btree</i>	no	A sequential version of our B-tree
<i>seq btree (n/h)</i>	no	Our sequential B-tree without hints
<i>reduction btree</i>	yes	Google's B-tree combined with bulk inserts through parallel reduction
<i>btree</i>	yes	Our optimistic B-tree
<i>btree (n/h)</i>	yes	Our optimistic B-tree without hints
<i>brie</i>	yes	Our Brie

**TABLE 1** Summary of investigated data structures





**FIGURE 6** Memory usage of the B-tree and BRIE. A, B-tree vs BRIE with different point densities (10 000 samples each); B, Memory requirements for 100 M points depending on data point density

## 6.1 | Memory requirements

The first aspect to be investigated in order to determine the suitability of our data structures for representing relations within a Datalog query processor is their memory consumption. In theory, both our data structures exhibit  $\mathcal{O}(N)$  memory requirements. However, the memory usage of the BRIE depends on the *density* of those entries, while the memory usage of the B-tree is approximately consistent regardless of the characteristics of the data.

In the BRIE, the data points (10, 11) and (10, 101) will be represented through different 1 bit in different leaf nodes. On the contrary, the elements (10, 11) and (10, 12) would be represented through 1 bit in the same leaf node, consuming less memory. Since the values of stored entries are used to address the bit to be set for their representation, the actual values of the stored data points become relevant for the overall memory consumption.

Let  $S \subseteq \mathbb{N}^n$  be a set of  $n$ -dimensional points and

$$B(S) = \{ \bar{x} \in \mathbb{N}^n \mid \forall 1 \leq i \leq n, \exists \bar{I} \in S. \exists \bar{u} \in S. \bar{I}_i \leq x_i \leq \bar{u}_i \}$$

the set of points contained in the axis aligned bounding box of  $S$ . Then, we define the *density*  $d(S)$  of  $S$  by

$$d(S) = \frac{|S|}{|B(S)|} \in (0, \dots, 1].$$

Therefore, to evaluate the memory requirements of the BRIE data structure, the density of data points needs to be considered.

For brevity, we focus on 2D data in this evaluation section. However, as additional evaluation showed, similar results are obtained when evaluating the memory consumption or performance of one- or higher-dimensional point sets.

The memory usage of the BRIE filled with data of various densities is illustrated in Figure 6A. Results labeled “brie X%” are obtained by maintaining a data density of X%. These results demonstrate a linear scaling of memory usage with respect to the number of elements stored in the data structures. It also demonstrates that the BRIE becomes more efficient with higher data density, with the cut-off point being at approximately 1% density where the memory usage of BRIE and B-tree are equal. Therefore, with sufficiently high-density data, the BRIE provides higher storage efficiency compared to the B-tree; however, with lower densities, the B-tree vastly outperforms the BRIE.

To quantify the impact of the data density on the memory requirements of a BRIE, Figure 6B summarizes the total memory usage of three different data structures to store 100 M 2D data points of varying point density. Only the memory requirements of BRIES are affected by the density. In total, to store a set of 2D points ( $2 \times 4$  bytes each), C++’s `std::set` requires 40 bytes, Google’s B-tree between 8.8 and 10.5 bytes, Soufflé’s B-tree 9.4 to 13.3 bytes, and our BRIE  $\frac{0.142}{d(S)}$  bytes per point. Therefore, even at 2% density, our BRIE uses 6.4 bytes per point, which is smaller than a direct encoding using 8 bytes.

## 6.2 | Sequential performance

In the first step, we evaluate the execution time required for the three most performance critical operations: inserts, membership tests, and range queries. The result of this evaluation is a base-line reference of the performance of our data structures as well as insights on the overhead induced by the employed locking mechanisms.

For evaluating the performance of the insert operation, we insert varying numbers of 2D points<sup>‡</sup> into initially empty sets, measure the overall execution time, and compute the achieved throughput in insert/s. Thereby, we distinguish between an ordered and unordered use case. In

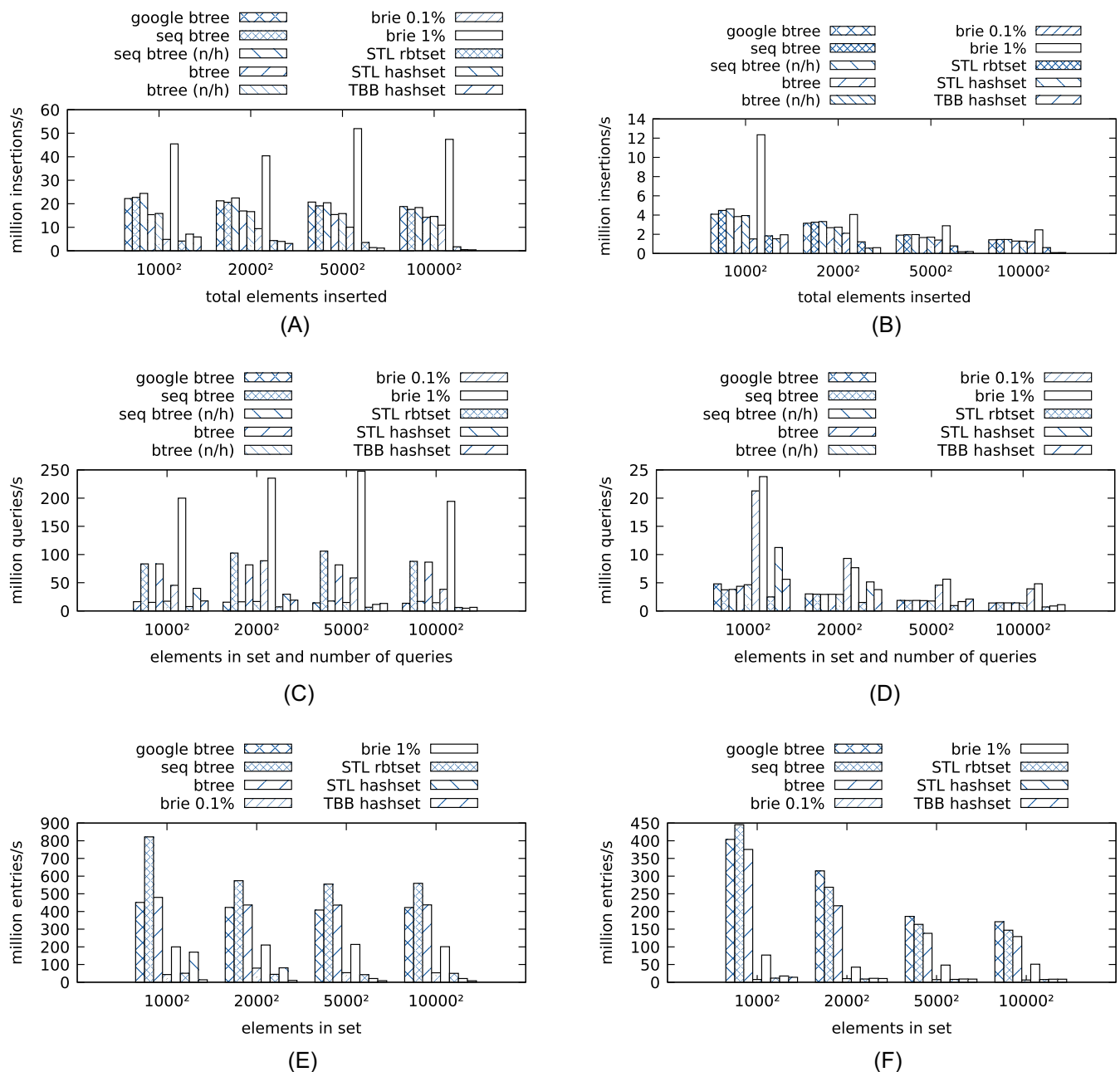
<sup>‡</sup>2D data is the most relevant case in many Datalog queries; besides, results remain similar for other dimensions

the ordered, the elements are inserted in their lexicographical order; in the unordered, a random order is employed. Furthermore, due to the sensitivity of the BRIE toward the density of data, varying point densities were tested for the BRIE.

For the membership query benchmark, we insert the same sets of elements into our candidate data structures, followed by querying for each element once in order and in a random sequence. The computation time of all queries is recorded, and the amortized query performance in queries/s obtained.

Finally, for the evaluation of the range query operation, we are focusing on the cost of scanning (aka iterating) through a range of elements since the cost of locating the start of a range is already covered by the membership test. Thus, for this benchmark, we are measuring the number of entries visited per second while iterating once through the entire set of elements within a relation. Note that hints are not applicable to the iteration operation.

Figure 7 summarizes the collected performance data. The first row illustrates the sequential insertion performance. As can be observed, the BRIE outperforms all alternative data structures at 1% density; however, with 0.1% density, the BRIE is less efficient than B-tree-based data structures. Further experiments at 100% density showed that BRIE's outperformed our B-tree by up to 5× in the ordered case and 11× in the random order case. Out of the remaining candidates, the B-tree data structures outperform the STL data structures and TBB's hash-based set.



**FIGURE 7** Sequential performance of performance critical set operations. A, Sequential insertion (ordered); B, Sequential insertion (random order); C, Membership test (ordered); D, Membership test (random order); E, Full-range scan (after ordered insert); F, Full-range scan (after random insert)

Despite the superior asymptotic runtime complexity of hash-based data structures, the cache-friendly organization of data in B-trees causes considerably fewer cache misses. This leads to superior performance compared to the random memory access pattern inherent in hash-based data structures. Among the B-tree data structures, the ordered insertions result in approximately 5× higher performance than random order insertions partly due to the improved cache utilization and the reduced complexity of inserting elements within leaves in order. Furthermore, in the random order case, larger dataset sizes lead to weaker performance. While in both cases, inner nodes need to be passed to reach insertion points; in the random order case, these inner nodes are more likely to trigger cache misses. This effect is even stronger in the BRIE, where cache misses are more likely to occur when the structures exceed the available cache sizes. Among our B-tree implementation, the operation hints cannot amortize their maintenance costs, and overhead caused by the locking mechanism can be observed (up to ~ 25% in ordered and ~ 15% in random).

Figure 7C and Figure 7D illustrate the observed query performance for ordered and unordered sequences of accesses. We observe that with sufficient data density, the BRIE outperforms all other data structures. With the random order case, the density has less effect, with 0.1% and 1% density performing similarly. This is due to cache effects, where the path in the tree up to an element may be cached and reused in the ordered case but cannot provide benefits in the random order case. We also observe that hints provide up to a 6× performance boost for membership tests for B-trees in the ordered case since navigating the tree can almost always be circumvented by using the hint. Secondly, the STL unordered set provides high performance for small datasets (up to 2× faster than TBB's hashset); however, this advantage vanishes with larger data set sizes.

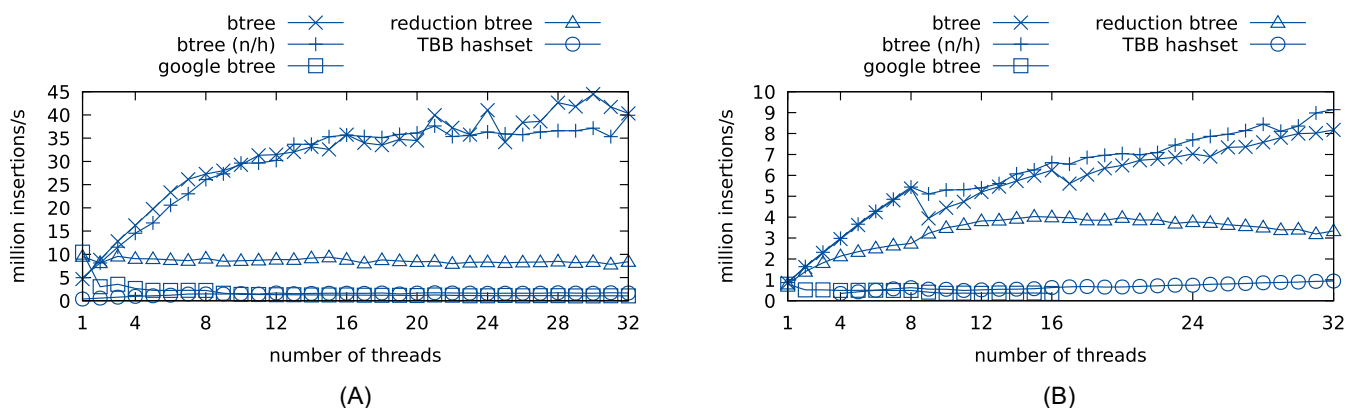
Finally, Figure 7E and Figure 7F summarize the rate at which elements stored in a data structure can be iterated through. Here, the bit-encoded values of the BRIE require reconstruction during the scan, and thus, performance is lower than our B-tree by a factor of 2. However, the compact storage of data in B-trees facilitates efficient iterations. The filling order, however, has an impact on the filling grade of leaf nodes, affecting the overall efficiency. Higher filling rates, as caused by in-order inserts, lead to a more compact tree, involving fewer nodes. Since every node switch is a potential cache miss, this leads to improved iteration speed. Finally, as can be observed, the integration of synchronization techniques and thus the necessary wrapping of key elements into atomic types are causing a performance deficit for our optimistic B-tree compared to its sequential equivalent.

Overall, the data underlines the superiority of BRIES with high-density data, with performance gains of up to 11× for insertions and membership tests. However, this advantage disappears when iterating through elements and with lower data densities. Regardless of data density, B-tree-based data structures are more performant compared to hash or red/black tree-based structures due to their cache efficiency. Moreover, these results demonstrate that our optimistic B-tree implementation exhibits sequential performance characteristics comparable to Google's state-of-the-art B-tree implementation. However, while Google's solution is thread unsafe, our concurrent B-tree is able to scale well in a parallel environment.

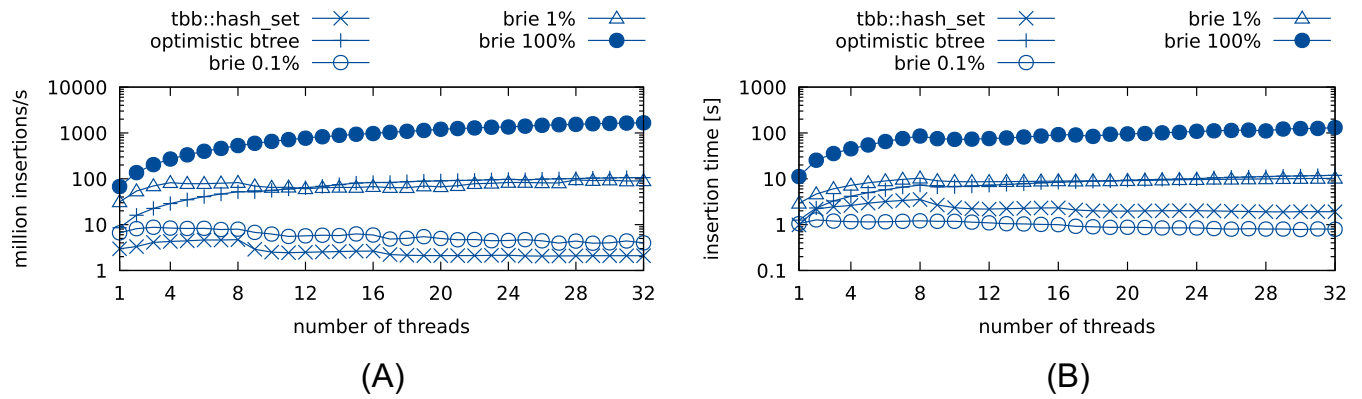
### 6.3 | Parallel performance

To evaluate the parallel insertion performance, we evaluate the parallel scalability of our candidate data structures. We insert 100 M 2D points into an initially empty set using a varying number of threads (strong scaling). Figure 8 summarizes the insertion throughput of our contestants with when gradually scaling up computation with multiple threads. Besides the three thread-safe data structures (btree, btree (n/h), and TBB hashset), we included two configurations realizing synchronization through external means: one variant with a global lock synchronizing insertions (*google btree*) and another where inserts are performed by threads on thread-private copies and merged through subsequent reduction step (*reduction btree*). While both techniques could be applied to any set implementation, for our evaluation, we chose the fastest sequential external option available to us, the Google B-tree.

Among the five contestants, only the global-lock-based approach failed, predictably, to gain performance improvements through the utilization of more than a single thread. The remaining data structures manage to do so in at least one scenario. In absolute terms, however, TBB's



**FIGURE 8** Parallel performance of insert operations. A, Parallel insertion (ordered, multisocket, 4x8 cores); B, Parallel insertion (random, multisocket, 4x8 cores)



**FIGURE 9** Comparison of insertion times of 100M 2D points, using ideal density for each data structure. A, Parallel insertion (ordered, multisocket, 4x8 cores; B, Parallel insertion (random, multisocket, 4x8 cores)

hash-based set is not able to compensate for its sequential performance deficit through parallel scaling. Our optimistic B-tree outperforms TBB's hashset by at least a factor of 8.5 with 1 thread and up to a factor of 59 with more threads.

The reduction-based approach manages to obtain speedups over sequential performance in the random order cases, where the computation effort for the thread-local insertions is dominating the concluding merge. By reducing this domination, either by making the thread-local insertions more efficient (ordered) or by reducing the number of elements inserted by each individual thread (more threads), the potential of this approach is reduced.

Regarding our B-trees, it can be observed that the operation hints cause little differences in the parallel insertion performance in the given use cases, as it has been observed in the sequential experiment above. Both manage to obtain a speedup of approximately 8.8 and 10 in the two use cases. Furthermore, Figure 8A and Figure 8B illustrate the scalability of the optimistic locking approach beyond socket boundaries, while TBB's hashset sustains significant performance losses beyond the first socket boundary (not visible in Figure 8B). Figure 8A thereby constitutes a use case performing most operations within NUMA boundaries. This is due to the partitioning of the elements to be inserted among the threads in the benchmarks, their pinning to cores, as well as the default first-touch NUMA policy. It thus demonstrates the achievable performance when being NUMA aware, while Figure 8B does not have this advantage. Consequentially, clear performance drops whenever growing beyond a single NUMA domain are observed. Our data structure is not inherently NUMA-aware.

We also evaluate the performance of the BRIE compared to the best (our B-tree) and worst (TBB's hashset) thread-safe data structures. For a fair comparison, we evaluate three different data densities for the BRIE: 0.1%, 1%, and 100% density.

Figure 9 summarizes the obtained results. We observe that at peak efficiency, the BRIE achieves more than an order of magnitude faster performance compared to our B-tree, in both ordered and unordered cases. However, with sparse data, this advantage disappears, and at 0.1% density performance is below even that of TBB's hashset in the unordered case. In terms of parallel scalability, BRIES achieve a parallel efficiency of up to 80% for 32 cores, in the best case with 100% density. However, when the data becomes too sparse, the high memory usage and associated memory management overhead and low locality result in a loss of efficiency.

We also observe that TBB's hashset implementation experiences performance penalties for each additional socket used in the computation. Both the BRIE and the B-tree also suffer slightly from crossing socket boundaries; however, they continue to benefit from additional parallel resources.

Overall, in the ordered as well as in the unordered insertion benchmarks, B-trees and BRIES achieve their best performance with the full 32 cores. In the ordered insertion benchmark, B-trees outperform TBB's state-of-the-art concurrent set implementations by a factor of 22, while BRIES outperform B-trees by an additional factor of 15, making BRIES up to 350x faster than TBB's concurrent set. Therefore, the provided data demonstrates the superior parallel capabilities of our data structures compared to the reference data structures.

## 6.4 | Processing Datalog queries

For our last experiment, we integrate both our B-tree and BRIES into the Soufflé Datalog engine. We compare their performance to our other candidate data structures when used for evaluating real-world large-scale Datalog analyses.

Our evaluation focuses on the resulting system's ability to decrease overall processing time with increasing availability of parallel resources while keeping the processed problem size fixed (strong scaling). This corresponds to the most common use case of attempting to analyze a fixed data set, in the case of program analysis a fixed code base, as quickly as possible. On the other hand, the complex relationship between input data and inflicted workload renders designing weak scaling experiments impractical for all but the most trivial analysis.

### 6.4.1 | B-tree

For our evaluation, we utilized two real-world benchmarks: a context-sensitive var-points-to analysis using the DOOP framework<sup>1</sup> and a security vulnerability analysis for an Amazon EC2 network. The DOOP analysis was run on the suite of DaCapo benchmarks, comprising 11 different Java programs. Both analyses comprise 100s of relations and rules. Table 2 summarizes the properties of those two benchmarks, along with runtime statistics for the inserts, membership tests, and lower/upper bound calls.

Figure 10 illustrates the collected performance data for our two benchmarks. Note that Figure 10A presents the total time for analysis of all 11 DaCapo benchmarks for brevity. For the DOOP static program analysis, we observe that the performance of our B-tree with a single thread is approximately 1.5× faster than the nearest reference data structure, the Google B-tree. In parallel workloads, our B-tree implementation provides better scalability than the reference data structures with global locks. Compared to the concurrent TBB hashset-equipped Datalog engine, the optimistic B-tree-based counterpart maintains similar scalability; however, performance is approximately 4× better than the TBB hashset for all numbers of threads. Moreover, the usage of operation hints improved performance by up to 10%.

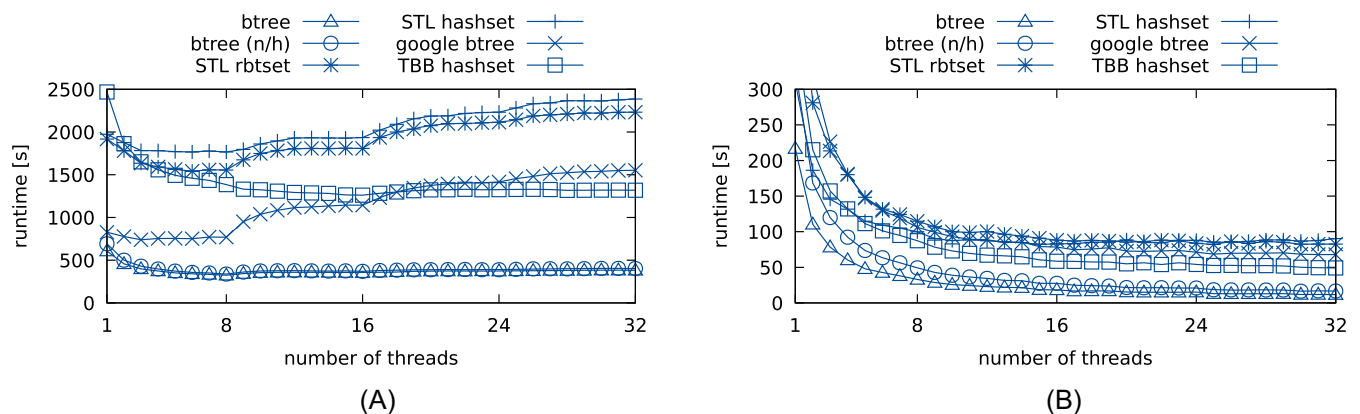
With the security vulnerability analysis, we observe that our B-tree performs approximately 2× better than the TBB hashset. Interestingly, with this workload, even the global locked data structures demonstrate some scalability since it is read-heavy rather than write-heavy. We also note that this security analysis contains more relations with fewer tuples (1.2e7 out of the 1.6e7 produced tuples were concentrated in a single relation), and as a result, the hash-based data structures perform better compared with the DOOP program analysis. Moreover, for this analysis, our B-tree with operation hints performs almost 1.5× as well as our B-tree without hints, indicating that ordered queries play an important part in this workload.

With regards to operation hints, collected statistics show that for the single-threaded DOOP analysis, up to 54% of operations resulted in an operation hint hit and up to 52% hits for the 16 thread case. With the security vulnerability analysis, these rates reach 77% for the single-threaded case and 76% for the 16 thread case. This further suggests that the security vulnerability analysis heavily involves ordered data and therefore demonstrates a larger performance improvement with the addition of operation hints.

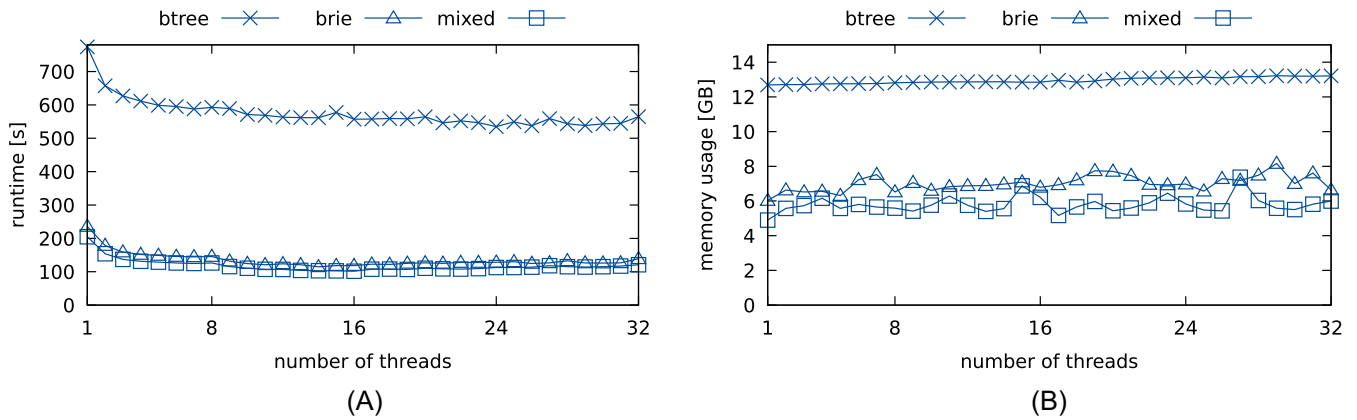
In both cases, the utilization of parallel resources provides performance improvements over the best sequential version. In the DOOP case, speedup of 1.9× can be obtained, and in the vulnerability case, a speedup of 8.4× compared to our B-tree run with 1 thread. These improvements are on top of the 1.3–3× sequential improvement compared to other reference data structures.

**TABLE 2** Real-world Datalog benchmark properties

Datalog property	DOOP on DaCapo (avg. per benchmark)	Amazon EC2 security vulnerability
Relations	493	287
Rules	810	236
Evaluation statistics	DOOP on DaCapo (avg. per benchmark)	Amazon EC2 security vulnerability
Inserts	8.3e7	2.1e7
Membership tests	1.5e8	4.2e9
Lower_bound calls	2.1e8	2.5e9
Upper_bound calls	2.1e8	2.5e9
Input tuples	8.3e6	3515
Produced tuples	2.5e7	1.6e7



**FIGURE 10** Comparison of different data structures for two real-world applications relying on large scale Datalog queries. A, Context-sensitive var-points-to analysis (insertion-heavy); B, Security vulnerability analysis (read-heavy)



**FIGURE 11** Comparison of Soufflé running a points-to analysis of OpenJDK with BRIE and B-tree data structures, with *mixed* denoting automatic selection of data structures. A, Runtime; B, Memory

### 6.4.2 | Brie

As demonstrated by the previous experiments, given the right circumstances, BRIEs offer vastly superior performance for (parallel) insertion and data query operations. Therefore, in practice, the BRIE data structure performs well when the workload consists of large volumes of high-density data. One example of such a workload is a points-to analysis of the OpenJDK dataset.

Figure 11 summarizes the performance of our BRIE data structure compared to B-tree when running Soufflé on this points-to analysis. We also include the performance when using an automatic *mixed* selection of data structures, whereby BRIE is used for relations with two or fewer dimensions and B-trees in all other cases. The rationale for such a mechanism is that lower dimension data is more likely to exhibit high density, and therefore, the BRIE is more suitable than the B-tree.

We observe a runtime improvement of up to 4× by using BRIE instead of B-tree and a memory usage reduction of up to 2×. The result suggests that much of the data in this workload is of high-density. Indeed, for the largest relation, *VarPointsTo*, we observed that B-tree used 11.9 bytes per tuple, while BRIE used 2.2 bytes per tuple. Being a binary relation, directly storing the data would use 8 bytes per tuple, and therefore, the BRIE exploits density to store tuples 3.6× more efficiently than even a direct encoding. We also observe a slight performance improvement by using a mixed data structure selection, suggesting that the workload also contains some higher dimension relations, with low data density. For these relations, BRIE performed worse than B-tree, and thus, optimal performance is obtained by using a combination of the two.

### 6.4.3 | Comparison with concurrent tree data structures

In this section, we compare our B-tree data structure to not alternative data concurrent data structures. These alternatives include *PALM tree*,<sup>27</sup> *Masstree*,<sup>28</sup> and *B-slack tree*.<sup>29</sup> We outline the limitations of these data structure's suitable for Datalog evaluation.

*PALM tree*<sup>27</sup> is a concurrent lock-free B+ tree implementation. It uses an internal synchronization strategy, where elements to be inserted are added to an internal queue and processed concurrently by the data structure itself. Moreover, *PALM tree* uses AVX instructions and therefore only supports single-integer keys and not tuples as required for Datalog processing.

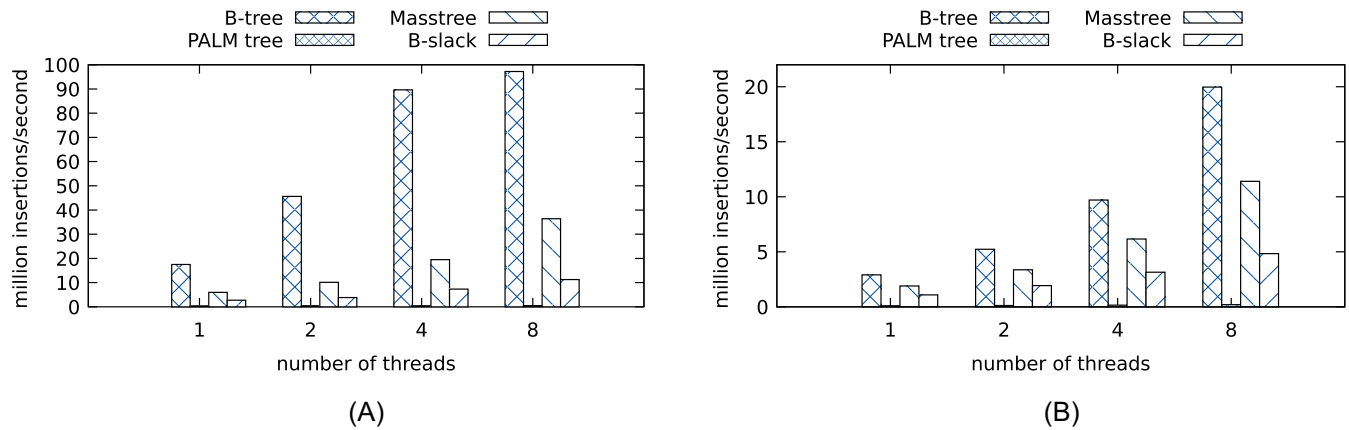
*Masstree*<sup>28</sup> is built as a client/server architecture, designed for use-cases requiring persistence across reboots. Hence, it is not optimized for use in an in-memory Datalog engine. Moreover, *Masstree* only supports strings as keys and thus would require a significant redesign of soufflé.

*B-slack tree*<sup>29</sup> constitutes a variation of B-tree that constrains the overall fill grade of all child nodes. As a result, B-slack trees exhibit better worst-case space complexity than conventional B-trees. B-slack trees weakens structural constraints and decouples the insertion from the rebalancing steps making the locking of B-trees more local. However, B-slack trees do not specify the locking scheme. Absolute insertion performance or parallel scalability have not been investigated in the work of Brown.<sup>29</sup>

We compared the performance of alternative data structures to our B-tree. Due to the limitations of the data structures for our Datalog workload, we provide an additional set of benchmarks inserting and reading a set of 10 000 000 fixed-size 32-bit integers, in both sequential and random order. The results are in Figures 12 and Table 3. Note that *Masstree* was evaluated using the included benchmark utility as the client/server architecture could not be integrated with our microbenchmark suite. The results demonstrate that our B-tree exhibits up to 3× better sequential insertion performance and up to 1.5× better random insertion performance compared to the next best data structure (*Masstree*). Moreover, our data structure exhibited up to 8× higher throughput than B-slack tree, and 200× higher throughput than *PALM tree*.

The repo <http://github.com/souffle-lang/ppopp19> contains the experiments without some industrial benchmarks whose source-code could not be disclosed.





**FIGURE 12** Comparison of parallel 32-bit integer insertion throughput of our B-tree, PALM tree, Masstree, and B-slack (see Table 3). A, Runtime; B, Memory

**TABLE 3** Throughput inserting 32-bit integers

Threads	Throughput (ordered/random) [ $10^6$ elements/second]			
	B-tree	PALM tree	Masstree	B-slack
1	17.5/2.91	0.38/0.10	5.99/1.90	2.73/1.09
2	45.61/5.23	0.44/0.13	10.11/3.37	3.83/1.94
4	89.64/9.71	0.47/0.16	19.52/6.17	7.32/3.15
8	97.19/16.97	0.49/0.20	36.38/11.41	11.29/4.84

## 7 | RELATED WORK

### 7.1 | Data structures for Datalog

There have been a number of approaches for Datalog evaluation based on various alternative data structures. For example, *bddbdb*<sup>7</sup> employs binary decision diagrams to store relations. While BDDs are particularly applicable for program analyses written in Datalog, their use also introduces the new issue of finding optimal variable orderings. Hash-based sets have also been used in systems such as  $\mu Z$ <sup>10</sup> and *Socialite*.<sup>3</sup> Liu et al propose a novel approach based on nesting arrays and linked lists, demonstrating a worst-case constant running time. However, experiments show that in practice, this technique is not scalable for large datasets. In our experience, B-trees (as used in PA-Datalog/LogicBlox<sup>30</sup> version 3 and Soufflé<sup>24</sup>) have been the most suitable data structure for large ruleset/dataset benchmarks.<sup>8</sup> In particular, the work of Subotić et al<sup>16</sup> demonstrates the advantages of B-trees with novel indexing schemes specifically designed for these use cases. Moreover, by exploiting the property of seminaïve evaluation that data structures are never read from and written to at the same time, we are able to further optimize the design of the data structure, using ideas from Shun and Blelloch.<sup>17</sup>

While B-trees demonstrate scalable performance for the general case, certain properties of the use case (eg, density) may be exploited by specialized data structures. Therefore, engines such as  $\mu Z$  also include a framework for integrating different backends depending on the use case. For example, the work of Lopes et al<sup>2</sup> implements a BDD-based and a bit vector-based backend to more effectively compute networking problems. For the program analysis use case, data with high-density is prevalent, and the Brie demonstrates an advantage in these situations.

### 7.2 | Parallel Datalog engines

The parallelization of Datalog engines has been well studied in the past.<sup>3,31–38</sup> These approaches are mainly focused on using distributed systems and thus developing novel communication methods to efficiently evaluate Datalog in parallel. While most distributed Datalog approaches focus on top-down evaluation other works,<sup>3,32</sup> propose methods to parallelize seminaïve evaluation. However, in a distributed setting, the main bottleneck for large-scale Datalog evaluation becomes the network communication overhead. Hence, approaches targeted to shared-memory multicore machines have been developed, employing fine-grain parallelism to bottom-up evaluation schemes. For instance, Datalog-MC<sup>39</sup> uses hash tables internally, with range-partitioning providing the capability to execute each partition on separate cores of a shared-memory multicore system, using a variant of hash-join. Each partition is then protected by a lock during evaluation, thus incurring appreciable overheads. To evaluate Datalog programs, Datalog-MC represents rules as and/or trees, which are compiled to Java. LogicBlox version 4<sup>11</sup> uses persistent functional data structures that avoid the need for synchronization. Instead, by virtue of their immutability, insertions efficiently replicate state via the persistent data structure. Another approach for parallelizing Datalog is to evaluate using GPUs. Martinez-Angeles et al implement such a GPU-based engine.<sup>9</sup> Their basic data structure is an array of tuples, allowing for duplicates. After each relational operation, explicit duplicate elimination is performed, which leads to large overheads in some situations. Moreover, the potentially high number of duplicates occurring in temporary results quickly exceeded the memory budget on GPUs. This approach has been extended by Red Fox,<sup>40</sup> which relies on a more sophisticated key-value



store and allows more complicated queries than the earlier work. While GPU-based Datalog evaluation demonstrates significant speedups over traditional CPU-based techniques, they are limited by the amount of available memory for large-scale benchmarks. Moreover, complex use cases such as program analysis or network analysis have not been studied in GPU engines. For further performance comparisons between Datalog engines, we point the reader to other works.<sup>18,41</sup>

### 7.3 | Concurrent tree data structures

B-trees are among the most successful data structures for indexed data. The development of locking techniques for concurrent usage has been a well studied area.<sup>12</sup> However, most research focuses on the secondary storage use case.<sup>42</sup> For in-memory usage, a modified B-tree known as B-link tree managed to eliminate the need for read locks.<sup>43</sup> By adding a “link” pointer to each node, pointing to the next node on the same level of the tree, Lehman and Yao demonstrate that correct concurrent operation occurs even before the usual tree pointers are changed for a new split node. Unfortunately, we have not been able to obtain an implementation for comparison. An alternative approach has obtained good results by utilizing hardware transactional memory features available on some recent Intel architectures for synchronizing B-tree insert operations.<sup>20</sup> Their evaluation shows comparable parallel scalability to our optimistic approach. However, special hardware support is required for those, and multsocket systems have not been evaluated. Concurrency has been investigated in several tree-like data structures for general use. For example, the data structure in the work of Bronson et al<sup>44</sup> is similar in spirit to our work with an optimistic concurrency, which allows invisible readers. The approach in the work of Drachsler et al<sup>45</sup> maintains interval information to determine the placement of data. The data structure in the work of Howley and Jones<sup>46</sup> is based on single-word reads, writes, and compare-and-swap, where its algorithm operations only contend if concurrent updates affect the same nodes. Other concurrent tree-like structures have been presented in other works.<sup>13,14,47-53</sup> The involved design decisions are orthogonal to the locking scheme presented by our work. Realizing a C++ version of the B-slack tree utilizing our seq-lock-based synchronization scheme has the potential of yielding a highly scalable concurrent implementation. A recent approach<sup>29</sup> constrains the overall fill grade of all child nodes providing better worst-case space complexity than conventional B-trees and weakens the structural properties of the tree to improve the locality of modifications. However, the design does not specify a node-locking scheme for a concurrent implementation.

## 8 | CONCLUSION

In this paper, we have presented a new approach to parallel Datalog evaluation that relies on a data structure framework that can incorporate multiple specialized data structures. The framework, aside from allowing easy integration of novel data structures, represents a research platform to permit the investigation of new data structures for Datalog evaluation.

We have argued that since the appropriate choice of the data structure depends on the Datalog application, there is no need to find a single data structure for all applications and workloads. This resembles a very practical solution to fill the gap between the logic evaluation and the actual implementation of relations in the form of data structures.

To demonstrate the effectiveness of our approach, we have presented two data structures that follow our framework: (1) an in-memory B-tree data structure providing an efficient base-line solution for a wide variety of relations and (2) the Brie data structure specialized for high-volume and highly dense relations. Our experiments have demonstrated that significant speed-ups (up to 6×) can be achieved by varying specialized data structure for given use-cases.

We will leave more detailed investigations of NUMA effects, heterogeneous architectures, and scaling beyond a single node to future work.

### ACKNOWLEDGMENTS

We thank Cristina Cifuentes and Oracle Labs in Brisbane where early versions of this work was done, Byron Cook and the ARG team at Amazon Web Services. This research was supported partially by the Australian Government through the ARC Discovery Project funding scheme under grant DP180104030, the Austrian Research Promotion Agency (FFG) under contract 868018, and a research donation from AWS.

### ORCID

Herbert Jordan  <https://orcid.org/0000-0002-0884-270X>

Bernhard Scholz  <https://orcid.org/0000-0002-7672-7359>

### REFERENCES

1. Bravenboer M, Smaragdakis Y. Strictly declarative specification of sophisticated points-to analyses. *ACM SIGPLAN Not.* 2009;44(10):243-262.
2. Lopes NP, Bjørner N, Godefroid P, Jayaraman K, Varghese G. Checking beliefs in dynamic networks. Paper presented at: 12th USENIX Symposium on Networked Systems Design and Implementation; 2015; Oakland, CA.
3. Seo J, Park J, Shin J, Lam MS. Distributed socialite: a datalog-based language for large-scale graph analysis. *Proc VLDB Endow.* 2013;6(14):1906-1917.
4. Brent L, Jurisevic A, Kong M, et al. Vandal: a scalable security analysis framework for smart contracts. *arXiv preprint arXiv:1809.03981*; 2018.
5. Jordan H, Subotić P, Zhao D, Scholz B. A specialized B-tree for concurrent datalog evaluation. In: *Proceedings of the 24th Symposium on Principles and Practice of Parallel Programming*; 2019; Washington, DC.

6. Abiteboul S, Hull R, Vianu V, eds. *Foundations of Databases: The Logical Level*. 1st ed. Boston, MA: Addison-Wesley Longman Publishing; 1995.
7. Whaley J, Avots D, Carbin M, Lam MS. Using datalog with binary decision diagrams for program analysis. Paper presented at: Asian Symposium on Programming Languages and Systems; 2005; Tsukuba, Japan.
8. Jordan H, Scholz B, Subotić P. Soufflé: on synthesis of program analyzers. Paper presented at: International Conference on Computer Aided Verification; 2016; Toronto, Canada.
9. Martinez-Angeles CA, Dutra I, Costa VS, Buenabad-Chávez J. A datalog engine for GPUs. In: *Declarative Programming and Knowledge Management*. Cham, Switzerland: Springer; 2014:152-168.
10. Hoder K, Bjørner N, de Moura L.  $\mu Z$  – an efficient engine for fixed points with constraints. Paper presented at: International Conference on Computer Aided Verification; 2011; Snowbird, UT.
11. LogicBlox. Declarative cloud platform for applications that combine transactions & analytics. <http://www.logicblox.com>
12. Graefe G. A survey of B-tree locking techniques. *ACM Trans Database Syst*. 2010;35(3). Article No. 16.
13. Oshman R, Shavit N. The SkipTrie: low-depth concurrent search without rebalancing. Paper presented at: Proceedings of the 2013 ACM Symposium on Principles of Distributed Computing; 2013; Montréal, Canada.
14. Prokopec A, Bronson NG, Bagwell P, Odersky M. Concurrent tries with efficient non-blocking snapshots. *ACM SIGPLAN Not*. 2012;47(8):151-160.
15. Jordan H, Subotić P, Zhao D, Scholz B. Brie: a specialized trie for concurrent datalog. In: Proceedings of the 10th International Workshop on Programming Models and Applications for Multicores and Manycores; 2019; Washington, DC.
16. Subotić P, Jordan H, Chang L, Fekete A, Scholz B. Automatic index selection for large-scale datalog computation. *Proc VLDB Endow*. 2018;12(2):141-153.
17. Shun J, Blelloch GE. Phase-concurrent hash tables for determinism. In: Proceedings of the 26th ACM Symposium on Parallelism in Algorithms and Architectures; 2014; Prague, Czech Republic.
18. Scholz B, Jordan H, Subotić P, Westmann T. On fast large-scale program analysis in datalog. In: Proceedings of the 25th International Conference on Compiler Construction; 2016; Barcelona, Spain.
19. Bayer R, McCreight E. Organization and maintenance of large ordered indexes. *Acta Informatica*. 1972;1:173-189.
20. Karnagel T, Dementiev R, Rajwar R, et al. Improving in-memory database index performance with Intel® transactional synchronization extensions. Paper presented at: 2014 IEEE 20th International Symposium on High Performance Computer Architecture; 2014; Orlando, FL.
21. Lameter C. Effective synchronization on Linux/NUMA systems. In: Gelato Conference; 2005; Palo Alto, CA.
22. Menascé DA, Nakanishi T. Optimistic versus pessimistic concurrency control mechanisms in database management systems. *Information Systems*. 1982;7(1):13-27.
23. Boehm H-J. Can seqlocks get along with programming language memory models? In: Proceedings of the 2012 ACM SIGPLAN Workshop on Memory Systems Performance and Correctness; 2012; Beijing, China.
24. Developers TS. Soufflé – a datalog engine. <http://www.github.com/souffle-lang/souffle>. Accessed January 5, 2019.
25. Niebler E. B-tree containers from Google. <https://isocpp.org/blog/2013/02/b-tree-containers-from-google>. Accessed February 14, 2017.
26. Reinders J. *Intel Threading Building Blocks: Outfitting C++ for Multi-Core Processor Parallelism*. Sebastopol, CA: O'Reilly Media, Inc; 2007.
27. Sewall J, Chhugani J, Kim C, Satish N, Dubey P. Palm: parallel architecture-friendly latch-free modifications to B+ trees on many-core processors. *PVLDB*. 2011;4:795-806.
28. Mao Y, Kohler E, Morris RT. Cache craftiness for fast multicore key-value storage. In: Proceedings of the 7th ACM European Conference on Computer Systems; 2012; Bern, Switzerland.
29. Brown T. B-slack trees: space efficient B-trees. In: Ravi R, Gørtz IL, eds. *Algorithm Theory - SWAT 2014*. Cham, Switzerland: Springer International Publishing; 2014:122-133.
30. Aref M, ten Cate B, Green TJ, et al. Design and implementation of the LogicBlox system. In: Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data; 2015; Melbourne, Australia.
31. Shkapsky A, Yang M, Interlandi M, Chiu H, Condie T, Zaniolo C. Big data analytics with datalog queries on spark. In: Proceedings of the 2016 International Conference on Management of Data; 2016; San Francisco, CA.
32. Shaw M, Koutris P, Howe B, Suciu D. Optimizing large-scale semi-naïve datalog evaluation in Hadoop. In: *Datalog in Academia and Industry*. Berlin, Germany: Springer-Verlag; 2012:165-176.
33. Hulin G. Parallel processing of recursive queries in distributed architectures. In: *Proceedings of the Fifteenth International Conference on Very Large Data Bases*. San Francisco, CA: Morgan Kaufmann Publishers Inc; 1989:87-96.
34. Cohen S, Wolfson O. Why a single parallelization strategy is not enough in knowledge bases. In: Proceedings of the 8th ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems; 1989; Philadelphia, PA.
35. Ganguly S, Silberschatz A, Tsur S. A framework for the parallel processing of datalog queries. In: Proceedings of the 1990 ACM SIGMOD International Conference on Management of Data; 1990; Atlantic City, NJ.
36. Seib J, Lausen G. Parallelizing Datalog programs by generalized pivoting. In: Proceedings of the 10th ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems; 1991; Denver, CO.
37. Wolfson O, Silberschatz A. Distributed processing of logic programs. *ACM SIGMOD Rec*. 1988;17(3):329-336.
38. Wolfson O, Ozeri A. A new paradigm for parallel and distributed rule-processing. *ACM SIGMOD Rec*. 1990;19(2):133-142.
39. Yang M, Shkapsky A, Zaniolo C. Scaling up the performance of more powerful datalog systems on multicore machines. *VLDB J*. 2017;26(2):229-248.
40. Wu H, Damos G, Sheard T, et al. Red fox: an execution environment for relational query processing on GPUs. In: Proceedings of Annual IEEE/ACM International Symposium on Code Generation and Optimization; 2014; Orlando, FL.
41. Antoniadis T, Triantafyllou K, Smaragdakis Y. Porting doop to Soufflé: a tale of inter-engine portability for datalog-based analyses. In: Proceedings of the 6th ACM SIGPLAN International Workshop on State Of the Art in Program Analysis; 2017; Barcelona, Spain.
42. Kung HT, Robinson JT. On optimistic methods for concurrency control. *ACM Trans Database Syst*. 1981;6(2):213-226.
43. Lehman PL, Yao SB. Efficient locking for concurrent operations on B-trees. *ACM Trans Database Syst*. 1981;6(4):650-670.
44. Bronson NG, Casper J, Chafi H, Olukotun K. A practical concurrent binary search tree. *ACM SIGPLAN Not*. 2010;45(5):257-268.

45. Drachsler D, Vechev M, Yahav E. Practical concurrent binary search trees via logical ordering. *ACM SIGPLAN Not.* 2014;49(8):343-356.
46. Howley SV, Jones J. A non-blocking internal binary search tree. In: *Proceedings of the 24th Annual ACM Symposium on Parallelism in Algorithms and Architectures*; 2012; Pittsburgh, PA.
47. Brown T, Helga J. Non-blocking  $k$ -ary search trees. In: *Principles of Distributed Systems*. Berlin, Germany: Springer-Verlag; 2011:207-221.
48. Braginsky A, Petrank E. A lock-free B+tree. In: *Proceedings of the Twenty-Fourth Annual ACM Symposium on Parallelism in Algorithms and Architectures*; 2012; Pittsburgh, PA.
49. Crain T, Gramoli V, Raynal M. A speculation-friendly binary search tree. In: *Proceedings of the 17th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*; 2012; New Orleans, LA.
50. Levandoski JJ, Lomet DB, Sengupta S. The Bw-tree: a B-tree for new hardware platforms. Paper presented at: 2013 IEEE 29th International Conference on Data Engineering (ICDE); 2013; Brisbane, Australia.
51. Arbel M, Attiya H. Concurrent updates with RCU: search tree as an example. In: *Proceedings of the 2014 ACM Symposium on Principles of Distributed Computing*; 2014; Paris, France.
52. Natarajan A, Mittal N. Fast concurrent lock-free binary search trees. In: *Proceedings of the 19th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*; 2014; Orlando, FL.
53. Basin D, Bortnikov E, Braginsky A, et al. Kiwi: a key-value map for scalable real-time analytics. *ACM SIGPLAN Not.* 2017;52(8):357-369.

**How to cite this article:** Jordan H, Subotić P, Zhao D, Scholz B. Specializing parallel data structures for Datalog. *Concurrency Computat Pract Exper.* 2020;e5643. <https://doi.org/10.1002/cpe.5643>