

TREETOASTER: Towards an IVM-Optimized Compiler

Anonymous Author(s)

Abstract

A compiler’s optimizer operates over abstract syntax trees (ASTs), continuously applying rewrite rules to replace subtrees of the AST with more efficient ones. Especially on large source repositories, even simply finding opportunities for a rewrite can be expensive, as optimizer traverses the AST naively. In this paper, we leverage the need to repeatedly find rewrites, and explore options for making the search faster through indexing and incremental view maintenance (IVM). Concretely, we consider bolt-on approaches that make use of embedded IVM systems like DBToaster, as well as two new approaches: Label-indexing and TREEToASTER, an AST-specialized form of IVM. We integrate these approaches into an existing just-in-time data structure compiler and show experimentally that TREEToASTER can significantly improve performance with minimal memory overheads.

CCS Concepts: • Information systems → Database views; Query optimization.

Keywords: Abstract Syntax Trees, Compilers, Indexing, Incremental View Maintenance

ACM Reference Format:

Anonymous Author(s). 2021. TREEToASTER: Towards an IVM-Optimized Compiler. In *Proceedings of ACM Conference (Conference’17)*. ACM, New York, NY, USA, 13 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 Introduction

Typical database query optimizers, like Apache Spark’s Catalyst [3] and Greenplum’s Orca [34], work with queries encoded as abstract syntax trees (ASTs). A tree-based encoding makes it possible to specify optimizations as simple, composable, easy-to-reason-about pattern/replacement rules. Unfortunately, a significant portion of the optimizer’s time is spent simply searching the AST for nodes that match one of these patterns. In Apache Spark (Figure 1), between 20-80% of the optimizer’s time¹ is spent finding AST nodes eligible for rewriting². This can be tens or even hundreds of seconds on large queries.

In this paper, we propose TREEToASTER, an compiler-specific form of incremental view maintenance that virtually eliminates the cost of finding nodes eligible for a rewrite. In lieu of repeated linear scans through the AST for eligible

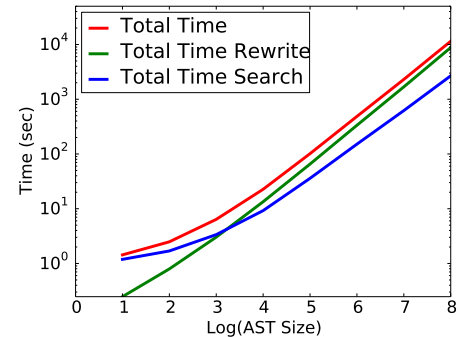


Figure 1. Apache Spark Catalyst’s optimizer spends 20-80% of its time searching for rewrites depending on AST size.

nodes, TREEToASTER materializes a view for each rewrite rule containing all nodes eligible for the rule, and incrementally maintains it as the tree evolves through the optimizer.

Naively, we might implement this incremental maintenance scheme by simply reducing the compiler’s pattern matching logic to a standard relational query language, and “bolting on” a standard database view maintenance system [21, 32]. This simple approach typically reduces search costs to a (small) constant, while adding only a negligible overhead to tree updates. However, classical view maintenance systems come with a significant storage overhead. As we show in this paper, TREEToASTER improves on the “bolt-on” approach by leveraging the fact that both ASTs and pattern queries are given as trees. As we show, when the data and query are both trees, TREEToASTER achieves similar maintenance costs without the memory overhead of caching intermediate results (Figure 2). TREEToASTER further reduces memory overheads by taking advantage of the fact that the compiler already maintains a copy of the AST in memory with pointers linking nodes together. TREEToASTER combines these compiler-specific approaches with standard techniques for view maintenance like inlining and compiling to C++ [21] to produce an incremental-view maintenance engine that meets or beats state-of-the-art view maintenance systems on AST pattern-matching workloads, while using significantly less memory.

To best illustrate the advantages of TREEToASTER, we apply it in the context of a recently proposed Just-in-Time Data Structure compiler [4, 18] that treats tree-based index data structures as ASTs. Like most AST-based optimizers, this compiler uses pattern/replacement rules to asynchronously identify opportunities for incremental reorganization, as in database cracking [16] or log-structured merge trees [28]. We implement TREEToASTER within JUSTINTIME DATA and show that it virtually eliminates AST search costs — the

¹The graph shows a 2^n -way union of 2^n -way joins, where n is on the x-axis.

²This is smaller in Orca, but still a significant 5-20% of the total time.

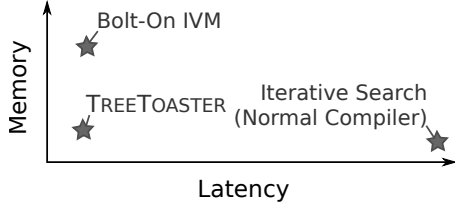


Figure 2. TREETOASTER achieves AST pattern-matching performance competitive with “bolting-on” an embedded IVM system, but with negligible memory overhead.

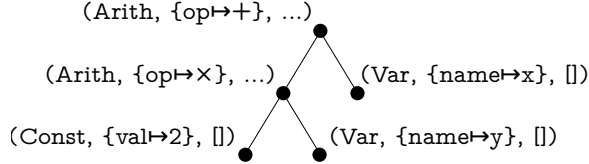


Figure 3. An AST for the expression $2 * y + x$

dominant cost of optimization in this system with minimal memory overheads.

Concretely, the contributions of this paper are: (i) We formally model AST pattern-matching queries and present a technique for incrementally maintaining precomputed views over such queries; (ii) We show how declaratively specified rewrite rules can be further inlined into view maintenance to further reduce maintenance costs; (iii) As a proof of concept, we “bolt-on” DBToaster, an embeddable IVM system, onto a just-in-time data-structure compiler [4, 18]. This modification dramatically improves performance, but adds significant memory overheads; (iv) We present TREETOASTER, a novel form of IVM optimized for compiler construction. TREETOASTER that avoids the high memory overheads of bolt-on IVM; (iv) We present experiments that show that TREETOASTER significantly outperforms “bolted-on” state-of-the-art IVM systems and is beneficial to the just-in-time data-structure compiler.

2 Notation and Background

In its simplest form, a typical compiler’s activities break down into three steps: parsing, optimizing, and output.

Parsing. First, a parser converts input source code into a structured Abstract Syntax Tree (AST) encoding the code.

Example 2.1. Figure 3 shows the AST for the expression $2 * y + x$. AST nodes have labels (e.g., *Arith*, *Var*, or *Const*) and attribute maps (e.g., $\{\text{op} \mapsto +\}$ or $\{\text{val} \mapsto 2\}$).

We formalize an AST as a tree with labeled nodes and annotated with zero or more attributes.

Definition 1 (Node). An Abstract Syntax Tree node $N = (\ell, A, \bar{N})$ is a 3-tuple, consisting of (i) a label ℓ drawn from an alphabet \mathcal{L} ; (ii) annotations $A : \Sigma_M \rightarrow \mathbb{D}$, a partial map from

an alphabet of attribute names Σ_M to a domain \mathbb{D} of attribute values; and (iii) an ordered list of children \bar{N} .

We define a leaf node (denoted $\text{isleaf}(N)$) as a node that has no child nodes. We assume that nodes follow a schema $\mathcal{S} : \mathcal{L} \rightarrow 2^{\Sigma_M} \times \mathbb{N}$; For each label ($\ell \in \mathcal{S}$), we fix a set of attributes that are present in all nodes with the label ($\bar{x} \in 2^{\Sigma_M}$), as well as an upper bound on the number of children ($c \in \mathbb{N}$).

Optimization. Next, the optimizer rewrites the AST, iteratively applying pattern-matching expressions and deriving replacement subtrees. We note that even compilers written in imperative languages frequently adopt a declarative style for expressing pattern-matching conditions. For example, ORCA [34] (written in C++) builds small ASTs to describe pattern matching structures, while Catalyst [3] (written in Scala) relies on Scala’s native pattern-matching syntax.

Example 2.2. A common rewrite rule for arithmetic eliminates no-ops like addition to zero. For example, the subtree

$$(\text{Arith}, \{\text{op} \mapsto +\}, [(\text{Const}, \{\text{val} \mapsto 0\}, []), (\text{Var}, \{\text{name} \mapsto b\}, [])])$$

can be replaced by $(\text{Var}, \{\text{name} \mapsto b\}, [])$ If the optimizer encounters a subtree with an *Arith* node at the root, *Const* and *Var* nodes as children, and a 0 as the value attribute of the *Const* node; it replaces the entire subtree by the *Var* node.

The optimizer continues searching for subtrees matching one of its patterns until no further matches exist (a fixed point), or an iteration threshold or timeout is reached.

Output. Finally, the compiler uses the optimized AST as appropriate by generating bytecode, a physical plan, etc....

2.1 Pattern Matching Queries

We formalize pattern matching in the following grammar:

Definition 2 (Pattern). A pattern query $q \in \mathcal{Q}$ is one of

$$\mathcal{Q} : \text{AnyNode} \mid \text{Match}(\mathcal{L}, \Sigma_I, \bar{Q}, \Theta)$$

The symbol $\text{Match}(\ell_q, i, \bar{Q}, \theta)$ indicates a structural match that succeeds iff (i) The matched node has label ℓ_q , (ii) the children of the matched node recursively satisfy $q_i \in \bar{Q}$, and (iii) the constraint θ over the attributes of the node and its children is satisfied. The node variable $i \in \Sigma_I$ is used to identify the node in subsequent use, for example to reference the node’s attributes. The symbol *AnyNode* matches any node. Figure 5 formalizes the semantics of \mathcal{Q} .

The grammar for constraints is given in Figure 4, and its semantics are typical. A variable atom $i.x$ is a 2-tuple of a Node name ($i \in \Sigma_I$) and an Attribute name ($x \in \Sigma_M$), respectively, and evaluates to $\Gamma(i)(x)$, given some scope $\Gamma : \Sigma_I \rightarrow \Sigma_M \rightarrow \mathbb{D}$. This grammar is expressive enough to capture the full range of comparisons ($>$, \geq , \leq , $<$, $=$, $<$), and so we use these freely throughout the rest of the paper.

$$\Theta : \text{atom} = \text{atom} \mid \text{atom} < \text{atom} \mid \Theta \wedge \Theta \mid \Theta \vee \Theta \mid \neg \Theta \mid \mathbb{T} \mid \mathbb{F}$$

$$\text{atom} : \text{const} \mid \Sigma_I . \Sigma_M \mid \text{atom} [+ , - , \times , \div] \text{atom}$$
Figure 4. Constraint Grammar

$$\llbracket q(\ell, A, [N_1 \dots N_n]) \rrbracket = \begin{cases} \mathbb{T}, \emptyset & \text{if } q = \text{AnyNode} \\ \mathbb{T}, \Gamma & \text{if } q = \text{Match}(\ell_q, i, [q_1 \dots q_n], \theta) \\ & \ell_q = \ell, \quad \theta(\Gamma), \\ & \llbracket q_1(N_1) \rrbracket = \mathbb{T}, \Gamma_1 \\ & \dots \llbracket q_n(N_n) \rrbracket = \mathbb{T}, \Gamma_n, \\ & \Gamma = \{ i \rightarrow A \} \cup \bigcup_{k \in [n]} \Gamma_k \\ \mathbb{F}, \emptyset & \text{otherwise} \end{cases}$$

Figure 5. Semantics for pattern queries ($q \in \mathcal{Q}$)

Example 2.3. Returning to Example 2.2, only Arith nodes over Const and Var nodes as children are eligible for the simplification rule. The corresponding pattern query is:

$$\text{Match}(\text{Arith}, A, [\text{Match}(\text{Const}, B, [], \{B.\text{val} = 0\}), \\ \text{Match}(\text{Var}, C, [], \mathbb{T})], \{A.\text{op} = +\})$$

Note the constraint on the Const match pattern; This sub-pattern only matches a node who's val(ue) attribute is 0.

We next formalize pattern matching over ASTs. First, we define the descendants of a node (denoted $\text{Desc}(N)$) to be the set consisting of N and its descendants:

$$\text{Desc}(N) \triangleq \{ N \} \bigcup_{k \in [n]} \text{Desc}(N_k) \text{ s.t. } N = (\ell, A, [N_1, \dots, N_n])$$

Definition 3 (Match). A match result, denoted $q(N)$, is the subset of N or its descendents on which q evaluates to true.

$$q(N) \triangleq \{ N' \mid N' \in \text{Desc}(N) \wedge \exists \Gamma : q(N') = \mathbb{T}, \Gamma \}$$

Pattern Matching is Expensive. Optimization is a tight loop in which the optimizer searches for a pattern match, applies the corresponding rewrite rule to the matched node, and repeats until convergence. Pattern matching typically requires iteratively traversing the entire AST. Every applied rewrite creates or removes opportunities for further rewrites, necessitating repeated searches for the same pattern. Even with intelligent scheduling of rewrites, the need for repeated searches can not usually be eliminated outright, and as shown in Figure 1 can take up to 80% of the optimizer's time.

Example 2.4. Continuing the example, the optimizer would traverse the entire AST looking for Arith nodes with the appropriate child nodes. A depth-first traversal ensures that any replacement happens before the optimizer checks the parent for eligibility. However, another rewrite may introduce new opportunities for simplification (e.g., by creating new Const nodes), and the tree traversal must be repeated.

$$\bar{R}_q \triangleq \begin{cases} \emptyset & \text{if } q = \text{AnyNode} \\ \{ (R_\ell \text{ AS } i) \} \bigcup_{x \in [n]} \bar{R}_{q_x} & \text{if } q = \text{Match}(\ell, i, [q_1, \dots, q_n], \theta) \end{cases}$$

$$\theta_q \triangleq \begin{cases} \mathbb{T} & \text{if } q = \text{AnyNode} \\ \theta \bigwedge_{x \in [n]} \theta_{q_x} \wedge \text{join}(i.\text{child}_x, q_x) & \text{if } q = \text{Match}(\ell, i, [q_1, \dots, q_n], \theta) \end{cases}$$

$$\text{join}(a, q) \triangleq \begin{cases} \mathbb{T} & \text{if } q = \text{AnyNode} \\ a = i.\text{id} & \text{if } q = \text{Match}(\ell, i, [q_1, \dots, q_n], \theta) \end{cases}$$

$$q \equiv \text{SELECT } * \text{ FROM } \bar{R}_q \text{ WHERE } \theta_q$$

Figure 6. Converting a pattern q to an equivalent SQL query.

3 Bolting-On IVM for Pattern Matching

As a warm-up, we start with a simple, naive implementation of incremental view maintenance for compilers by mapping our pattern matching grammar to relational queries and “bolting on” an existing system for incremental view maintenance (IVM). Although this specific approach falls short, it illustrates how IVM maps onto the pattern search problem. To map the AST to a relational encoding, for each label/schema pair $\ell \rightarrow \langle \{ x_1, \dots, x_k \}, c \rangle \in \mathcal{S}$, we define a relation $R_\ell(\text{id}, x_1, \dots, x_k, \text{child}_1, \dots, \text{child}_c)$ with an id field, and one field per attribute or child. Each node $N = (\ell, A, [N_1, \dots, N_c])$ is assigned a unique identifier id_N and defines a row of relation R_ℓ .

$$\langle \text{id}_N, A(x_1), \dots, A(x_k), \text{id}_{N_1}, \dots, \text{id}_{N_c} \rangle$$

A pattern q can be reduced to an equivalent query over the relational encoding, as shown in Figure 6. A pattern with k Match nodes becomes a k -ary join over the relations \bar{R}_q corresponding to the label on each Match node. Each relation is aliased to its node variable. Join constraints are given by parent/child relationships, and pattern constraints transfer directly to the WHERE clause.

Example 3.1. Continuing Example 2.2, the AST nodes are encoded as relations: Arith(id, op, child₁, child₂), Const(id, val), and Var(id, name). The corresponding pattern match query, following the process in Figure 6 is:

```
SELECT * FROM Arith a, Const b, Var c
WHERE a.child1 = b.id AND a.child2 = c.id
AND a.op = '+' AND b.val = 0
```

3.1 Background: Incremental View Maintenance

Materialized views are used in production databases to accelerate query processing. If a view is accessed repeatedly, database systems materialize the view query Q by precomputing its results $Q(D)$ on the database D . When the database changes, the view must be updated to match: Given a set of changes, ΔD (e.g., insertions or deletions), a naive approach

would be to simply recompute the view on the updated database $Q(D + \Delta D)$. However, if ΔD is small, most of this computation will be redundant. A more efficient approach is to derive a so-called “delta query” (ΔQ) that computes a set of updates to the (already available) $Q(D)$. That is, denoting the view update operation by \Leftarrow :

$$Q(D + \Delta D) \equiv Q(D) \Leftarrow \Delta Q(D, \Delta D)$$

Example 3.2. Recall $Q(\text{Arith}, \text{Const}, \text{Var})$ from the prior example. After inserting a row c into Const , we want:

$$\begin{aligned} & Q(\text{Arithmetic}, \text{Const} \uplus c, \text{Var}) \\ &= \text{Arith} \bowtie (\text{Const} \uplus c) \bowtie \text{Var} \\ &= (\text{Arith} \bowtie \text{Const} \bowtie \text{Var}) \uplus (\text{Arith} \bowtie c \bowtie \text{Var}) \\ &= Q(\text{Arith}, \text{Const}, \text{Var}) \uplus (\text{Arith} \bowtie c \bowtie \text{Var}) \end{aligned}$$

Instead of computing the full 3-way join, we can replace c with a singleton and compute the cheaper query $(\text{Arith} \bowtie c \bowtie \text{Var})$, and union the result with our original materialized view to obtain an updated view.

The cost of $\Delta Q(D, \Delta D)$ and \Leftarrow is generally lower than re-running the query, making this a win when database updates are small and infrequent. However, ΔQ can still be expensive. For larger or more frequent changes, we further reduce the cost of computing ΔQ by caching intermediate results. Ross et. al. [32] proposed a form of cascading IVM that caches all intermediate results in the physical plan of the view query.

Example 3.3. Continuing the example, we use the following execution order:

$$(\text{Arithmetic} \bowtie \text{Var}) \bowtie \text{Const}$$

In addition to materializing $Q(\cdot)$, Ross’ scheme also materializes the results of $Q_1 = (\text{Arithmetic} \bowtie \text{Var})$. When c is inserted into Const , the update only requires a simple 2-way join $c \bowtie Q_1$. However, updates are now (slightly) more expensive as multiple views may need to be updated.

Ross’ approach of caching intermediate state is analogous to typical approaches to fixpoint computation (e.g., Differential Dataflow [24]), but penalizes updates to tables early in the query plan. With DBToaster [21], Koch et. al. proposed instead materializing all possible query plans. Counter-intuitively, this added materialization significantly reduces the cost of view maintenance. Although far more tables need to be updated with every database change, the updates are generally small and efficiently computable.

3.2 Bolting DBToaster onto a Compiler

DBToaster [21] in particular is designed for embedded use. It compiles a set of queries down to a C++ or Scala data structure that maintains the query results. The data structure exposes insert, delete, and update operations for each source relation; and materializes the results of each query into an iterable collection. One strategy for improving compiler performance is to make the minimum set of changes

required (i.e., “bolt-on”) to allow it to use an incremental view maintenance data structure generated by DBToaster:

1. The reduction above generates SQL queries for each pattern-match query used by the optimizer.
2. DBToaster builds a view maintenance data structure.
3. The compiler is instrumented to register changes in the AST with the view maintenance data structure.
4. Iterative searches in the optimizer for candidate AST nodes are replaced with a constant-time lookup on the view maintenance data structure.

As we show in Section 7, this approach significantly outperforms naive iterative AST scans. Although DBToaster requires maintaining supplemental data structures, the overhead of maintaining these structures is negligible compared to the benefit of constant-time pattern match results.

Nevertheless, there are three major shortcomings to this approach. First, DBToaster effectively maintains a shadow copy of the entire AST — at least the subset that affects pattern-matching results. Second, DBToaster aggressively caches intermediate results. For example, our running example requires materializing 2 additional view queries, and this number grows combinatorially with the join width. Finally, DBToaster-generated view structures register updates at the granularity of individual node insertions/deletions, making it impossible for them to take advantage of the fact that most rewrites follow very structured patterns. For a relatively small number of pattern-matches, the memory use of the compiler with a DBToaster view structure bolted in increases by a factor of 2.5×. Given that memory consumption is already a pain point for large ASTs, this is not viable.

Before addressing these pain points, we first assess why they arise. First, DBToaster-generated view maintenance data structures are self-contained. When an insert is registered, the structure needs to preserve state for later use. Although unnecessary fields are projected away, this still amounts to a shadow copy of the AST. Second, DBToaster has a heavy focus on aggregate queries. Caching intermediate state allows aggressive use of aggregation and selection push-down into intermediate results, both reducing the amount of state maintained and the work needed to maintain views.

Both benefits are of limited use in pattern-matching on ASTs. Pattern matches are SPJ queries, mitigating the value of aggregate push-down. The value of selection push-down is mitigated by the AST’s implicit foreign key constraints: each child has a single parent and each child attribute references at most one child. Unlike a typical join where a single record may join with many results, here a single node only participates in a single join result³. This also limits the value of materializing for the sake of cache locality.

In summary, for ASTs, the cached state is either redundant or minimally beneficial. Thus a view maintenance scheme

³To clarify, a node may participate in multiple join results in different positions in the pattern match, but only in one result at the same position.

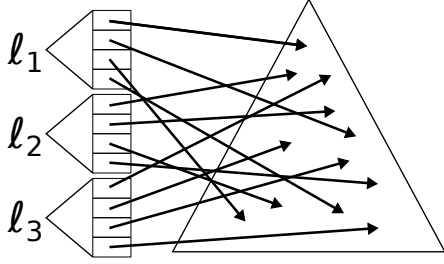


Figure 7. Indexing the AST by Label

designed specifically for compilers should be able to achieve the same benefits, but without the memory overhead.

4 Pattern Matching on a Space Budget

We have a set of patterns q_1, \dots, q_m and an evolving abstract syntax tree N . Our goal is, given some q_k , to be able to obtain a single, arbitrary element of the set $q_k(N)$ as quickly as possible. Furthermore, this should be possible without significant overhead as N evolves into N' , N'' , and so forth. Recall that there are three properties that have to hold for a node N to match q : (i) The node and pattern labels must match, (ii) Any recursively nested patterns must match, and (iii) The constraint must hold over the node and its descendants.

4.1 Indexing Labels

A standard first approach to accelerating queries is indexing, for example by building a secondary index over the node labels, as illustrated in Figure 7. For each node label, the index maintain a set of pointers to all nodes in the AST with that label. Updates to the AST are propagated into the index. Pattern match queries can use this index to scan a subset of the AST that includes only nodes with the appropriate label, as shown in Algorithm 1.

Algorithm 1: IndexLookup(N, q, Index_N)

Input: $N \in \mathcal{N}, q \in \mathcal{Q}, \text{Index}_N : \ell \rightarrow \{ \text{Desc}(N) \}$

Output: $N_{\text{match}} \in \text{Desc}(N)$

```

1 if  $q = \text{AnyNode}$  then
2   return  $N_{\text{match}} \leftarrow N$ 
3 else if  $q = \text{Match}(\ell, i, [q_1, \dots, q_n], \theta)$  then
4   for  $N_{\text{idx}} \in \text{Index}_N[\ell]$  do
5     if  $q(N) = \mathbb{T}, \Gamma$  then
6       return  $N_{\text{match}} \leftarrow N_{\text{idx}}$ 

```

Indexing the AST by node label is simple, and has a relatively small memory overhead: approximately 28 bytes per AST node using the C++ standard library `unordered_set`, with significant space for improvement. Similarly, the maintenance overhead is low — one hash table insert and/or remove per AST node changed.

Example 4.1. To find matches for the rule of Example 2.2, we retrieve a list of all `Arith` nodes from the index and iteratively check each for a pattern match. Note that this approach only supports filtering on labels; Recursive matches and constraints both need to be re-checked with each iteration.

4.2 Incremental View Maintenance

While indexing works well for single-node patterns, recursive patterns require a heavier-weight approach. Concretely, when a node in the AST is updated, we need to figure out which new pattern matches the update creates, and which pattern matches it removes. As we saw in Section 3, this could be accomplished by “joining” the updated node with all of the other nodes that could participate in the pattern. However, to compute these joins efficiently DBToaster and similar systems need to maintain a significant amount of supporting state: (i) The view itself, (ii) Intermediate state needed to evaluate subqueries efficiently (iii) A shadow copy of the AST. The insight behind TREEToASTER is that the latter two sources of state are unnecessary when the AST is already available: (i) Subqueries (inter-node joins) reduce to pointer chasing when the AST is available, and (ii) A shadow copy of the AST is unnecessary if the IVM system can navigate the AST directly.

We begin to outline TREEToASTER in Section 5 by defining IVM for immutable (functional) ASTs. This simplified form of the IVM problem has a useful property: When a node is updated, all of its ancestors are updated as well. Thus, we are guaranteed that the root node of a pattern match will be part of the change set (i.e., ΔD) and can simply look for pattern matches rooted at these nodes. We then refine the approach to mutable ASTs, where only a subset of the tree is updated. Then, in Section 6 we describe how declarative specifications of rewrite rules can be used to streamline the derivation of update sets and to eliminate unnecessary checks. Throughout Sections 5 and 6, we focus on the case of a single pattern query, but note that this approach generalizes trivially to multiple patterns.

5 IVM for ASTs

We first review a generalization of multisets proposed by Blizard [7] that allows for elements with negative multiplicities. A generalized multiset $\mathbf{M} : \text{dom}(\mathbf{M}) \rightarrow \mathbb{Z}$ maps a set of elements from a domain $\text{dom}(\mathbf{M})$ to an integer-valued multiplicity. We assume finite-support for all generalized multisets: only a finite number of elements are mapped to non-zero multiplicities. Union on a generalized multiset, denoted \oplus , is defined by summing multiplicities.

$$(\mathbf{M}_1 \oplus \mathbf{M}_2)(x) \triangleq \mathbf{M}_1(x) + \mathbf{M}_2(x)$$

Difference, denoted \ominus , is defined analogously:

$$(\mathbf{M}_1 \ominus \mathbf{M}_2)(x) \triangleq \mathbf{M}_1(x) - \mathbf{M}_2(x)$$

We write $x \in \mathbf{M}$ as a shorthand for $\mathbf{M}(x) \neq 0$. When combining sets and generalized multisets, we will abuse notation and lift sets to the corresponding generalized multiset, where each of the set's elements is mapped to 1.

A view View_q is a generalized multiset. We define the correctness of a view relative to the root of an AST. Without loss of generality, we assume that any node appears at most once in any AST.

Definition 4 (View Correctness). *A view View_q is correct for N if it maps exactly the subset of N and its descendants that match q to 1:*

$$\text{View}_q = \{ \{ N' \rightarrow 1 \mid N' \in q(N) \} \}$$

If we start with a view View_q that is correct for the root of an AST N and rewrite the AST's root to N' , we would like to update the view accordingly. We assume for the moment that we have an easy way to obtain a delta between the two ASTs: the difference: $\text{Desc}(N') \ominus \text{Desc}(N)$. This delta is generally small for a rewrite, including only the nodes of the rewritten subtree and their ancestors. We revisit this assumption in the following section. Algorithm 2 shows a simple algorithm for maintaining the View_q , given a small change Δ , expressed as a generalized multiset.

Algorithm 2: $\text{IVM}(q, \text{View}_q, \Delta)$

Input: $q \in \mathcal{Q}, \text{View}_q \in \{ \{ N \} \}, \Delta \in \{ \{ N \} \}$

Output: View'_q

```

1  $\text{View}'_q \leftarrow \text{View}_q$ 
2 for  $N_i \in \Delta$  /* nodes with multiplicity  $\neq 0$  */
3 do
4   if  $q(N_i) = \mathbb{T}, \Gamma$  then
5      $\text{View}'_q \leftarrow \text{View}'_q \oplus \{ \{ N_i \rightarrow \Delta(N_i) \} \}$ 

```

Example 5.1. Consider the AST of Figure 3, which contains five nodes, and our ongoing example rule. Let us assume that the left subtree is replaced by $\text{Const}(0)$ (e.g., if $\text{Var}(y)$ is resolved to 0). The multiset of the corresponding delta is:

$$\begin{aligned} & \{ \{ \text{Const}, \{ \text{val} \mapsto 0 \}, [] \} \mapsto 1, \{ \text{Arith}, \{ \text{op} \mapsto + \}, [\dots] \} \mapsto 1, \\ & \{ \text{Const}, \{ \text{val} \mapsto 2 \}, [] \} \mapsto -1, \{ \text{Var}, \{ \text{name} \mapsto y \}, [] \} \mapsto -1, \\ & \{ \text{Arith}, \{ \text{op} \mapsto \times \}, [\dots] \} \mapsto -1 \} \end{aligned}$$

Only one of the nodes with nonzero multiplicities matches, making the update: $\{ \{ \text{Arith}, \{ \text{op} \mapsto + \}, [\dots] \} \mapsto 1 \}$

For Algorithm 2 to be correct, we need to show that it computes exactly the update to View_q .

Lemma 5.2 (Correctness of IVM). *Given two ASTs N and N' and assuming that View_q is correct for N , then the generalized multiset returned by $\text{IVM}(q, \text{View}_q, \text{Desc}(N') \ominus \text{Desc}(N))$ is correct for N' .*

Proof Sketch. We examine the multiplicity of each node N'' over different cases of $N'' \in \text{Desc}(N), N'' \in \text{Desc}(N')$ and $q(N'') = \mathbb{T}/\mathbb{F}, \Gamma$ and show equivalence between the incrementally maintained and naively recomputed view. If $q(N'') = \mathbb{F}, \Gamma$ then the multiplicity remains unchanged for node N'' . Otherwise, The condition on line 2 would apply and the multiplicity gets recalculated according to line 3. The full proof appears in an associated tech report.

5.1 Mutable Abstract Syntax Trees

Although correct, IVM assumes that the AST is immutable: When a node changes, each of its ancestors must be updated to reference the new node as well. **Even when TREE_TOASTER is built into a compiler with immutable ASTs, many of these pattern matches will be redundant. By lifting this restriction (if in spirit only), we can decrease the overhead of view maintenance by reducing the number of nodes that need to be checked with each AST update.** To begin, we create a notational distinction between the root node N and the node being replaced R . For clarity of presentation, we again assume that any node R occurs at most once in N . $N[R \setminus R']$ is the node resulting from a replacement of R with R' in N :

$$N[R \setminus R'] = \begin{cases} R' & \text{if } N = R \\ (\ell, A, [N_1[R \setminus R'], \dots, N_n[R \setminus R']]) & \\ & \text{if } N = (\ell, A, [N_1, \dots, N_n]) \end{cases}$$

We also lift this notation to collections:

$$\text{View}[R \setminus R'] = \{ \{ N[R \setminus R'] \rightarrow c \mid (N \rightarrow c) \in \text{View} \} \}$$

We emphasize that although this notation modifies each node individually, this complexity appears only in the analysis. The underlying effect being modeled is a single pointer swap.

Example 5.3. The replacement of Example 5.1 is written:

$$N[(\text{Arith}, \{ \text{op} \mapsto \times \}, [\dots]) \setminus (\text{Const}, \{ \text{val} \mapsto 0 \}, [])]$$

In the mutable model, the root node itself does not change.

Definition 5 (Pattern Depth). *The depth $D(q)$ of a pattern q is the number of edges along the longest downward path from root of the pattern to an arbitrary pattern node q_i .*

$$D(q) = \begin{cases} 0 & \text{if } q = \text{AnyNode} \\ 1 + \max_{i \in [n]} (D(q_i)) & \text{if } q = \text{Match}(\ell, i, [q_1, \dots, q_n], \theta) \end{cases}$$

The challenge posed by mutable ASTs is that the modified node may make one of its ancestors eligible for a pattern-match. However, as we will show, only a bounded number of ancestors are required. Denote by $\text{Ancestor}_i(N)$ the i th ancestor of N ⁴. The maximal search set, which we now define, includes all nodes that need to be checked for matches.

⁴We note that ASTs do not generally include ancestor pointers. The ancestor may be derived by maintaining a map of node parents, or by extending the AST definition with parent pointers.

Definition 6 (Maximal Search Set). *Let R and R' be an arbitrary node in the AST and its replacement. The maximal search set for R and R' and pattern q , $\lceil R, R' \rceil_q$ is the difference between the generalized multiset of the respective nodes, their descendants, and their ancestors up to a height of $D(q)$.*

$$\lceil R, R' \rceil_q \triangleq \text{Desc}(R) \oplus \{ \lceil \text{Ancestor}_i(R) \rightarrow 1 \mid i \in [n] \rceil \} \\ \ominus \text{Desc}(R') \ominus \{ \lceil \text{Ancestor}_i(R') \rightarrow 1 \mid i \in [n] \rceil \}$$

Lemma 5.4. *Let N be the root of an AST, q be a pattern, and R and R' be an arbitrary node in the AST and its replacement. If View_q is correct for N . and $\text{View}'_q = \text{IVM}(q, \text{View}_q, \lceil R, R' \rceil_q)$, then $\text{View}'_q[R \setminus R']$ is correct for $N[R \setminus R']$*

Proof Sketch. Shown by building on the proof of Lemma 5.2 with a recursive proof that changing a subtree can not affect the matchability of a node more than $D(q)$ ancestors above.

Example 5.5. The pattern depth of our running example is 1. Continuing the prior example, only the node, its 1-ancestor (i.e., parent), and the 1-descendants (i.e., children) of the replacement node would need to be examined for view updates.

6 Inlining into Rewrite Rules

Algorithm 2 takes the set of changed nodes as an input. In principle, this information could be obtained by manually instrumenting the compiler to record node insertions, updates, and deletions. However, many rewrite rules are structured: The rule replaces exactly the matched set of nodes with a new subtree. Unmodified descendants are re-used as-is, and with mutable ASTs a subset of the ancestors of the modified node are re-used as well. TREE_TOASTER provides a declarative language for specifying the output of rewrite rules. This language serves two purposes. In addition to making it easier to instrument node changes for TREE_TOASTER, declaratively specifying updates opens up several opportunities for inlining-style optimizations to the view maintenance system. The declarative node generator grammar follows:

$$\mathcal{G} : \text{Gen}(\mathcal{L}, \overline{\text{atom}}, \overline{\mathcal{G}}) \mid \text{Reuse}(\Sigma_I)$$

A Gen term indicates the creation of a new node with the specified label, attributes, and children. Attribute values are populated according to a provided attribute scope $\Gamma : \Sigma_I \rightarrow \Sigma_M \rightarrow \mathbb{D}$. A Reuse term indicates the re-use of a subtree from the previous AST, provided by a node scope $\mu : \Sigma_I \rightarrow \mathcal{N}$. Node generators are evaluated by the $\llbracket \cdot \rrbracket_{\Gamma, \mu} : \mathcal{G} \rightarrow \mathcal{N}$ operator, defined as follows:

$$\llbracket g \rrbracket_{\Gamma, \mu} = \begin{cases} \mu(i) & \text{if } g = \text{Reuse}(i) \\ (\ell, \{a_1(\Gamma), \dots, a_k(\Gamma)\}, \llbracket g_1 \rrbracket_{\Gamma, \mu}, \dots, \llbracket g_n \rrbracket_{\Gamma, \mu}) & \text{if } g = \text{Gen}(\ell, [a_1, \dots, a_k], [g_1, \dots, g_n]) \end{cases}$$

A declaratively specified rewrite rule is given by a 2-tuple $\langle q, g \rangle \in \mathcal{Q} \times \mathcal{G}$, a match pattern describing the nodes to be removed from the tree, and a corresponding generator describing the nodes to be re-inserted into the tree. As a simplification for clarity of presentation, we require that

Reuse nodes reference nodes matched by AnyNodepatterns. Define the set of matched node pairs as the set

$$\text{pair}(q, R) = \{ \langle q, R \rangle \} \cup \dots \\ \dots \begin{cases} \{ \langle \text{AnyNode}, R \rangle \} & \text{if } q = \text{AnyNode} \\ \bigcup_{k \in [n]} \text{pair}(q_k, N_k) & \text{if } q = \text{Match}(\ell, x, [q_1, \dots, q_n], \theta) \\ & R = (\ell, A, [N_1, \dots, N_n]) \end{cases}$$

A set of generated node pairs $\text{pair}(g, \Gamma, \mu)$ is defined analogously relative to the node $\llbracket g \rrbracket_{\Gamma, \mu}$

Definition 7 (Safe Generators). *Let N be an AST root, q be a pattern query, and $R \in q(N)$ be a node of the AST matching the pattern. We call a generator $g \in \mathcal{G}$ safe for $\langle q, R \rangle$ iff g reuses exactly the wildcard matches of q . Formally:*

$$\langle \text{AnyNode}, N \rangle \in \text{pair}(q, R) \Leftrightarrow \langle \text{Reuse}(N), N \rangle \in \text{pair}(g, \Gamma, \mu)$$

Let $g \in \mathcal{G}$ be a generator that is safe for $\langle m, R \rangle$, where $m \in \mathcal{Q}$ is a pattern. The mutable update delta from N to $N[R \setminus \llbracket g \rrbracket_{\Gamma, \mu}]$ is:

$$\Delta = \{ \lceil N' \rightarrow 1 \mid \langle g', N' \rangle \in \text{pair}(g, \Gamma, \mu) \rceil \} \ominus \\ \{ \lceil N' \rightarrow 1 \mid \langle q', N' \rangle \in \text{pair}(m, R) \rceil \}$$

Note that the size of this delta is linear in the size of g and m .

6.1 Inlining Optimizations

Up to now, we have assumed that no information about the nodes in the update delta is available at compile time. For declarative rewrite rules, we are no longer subject to this restriction. The labels and structure of the nodes being removed and those being added are known at compile time. This allows TREE_TOASTER to generate more efficient code by eliminating impossible pattern matches.

Algorithm 3: $\text{Inline}_{gen}(q, g)$

Input: $q \in \mathcal{Q}, g \in \mathcal{G}$
Output: $f : \mathcal{N} \mapsto \{ \mathcal{N} \}$

- 1 **if** $q = \text{AnyNode} \vee g = \text{Reuse}(\mu)$ **then**
- 2 $f' \leftarrow (N \mapsto \{ N \})$
- 3 **else if** $q = \text{Match}(\ell, i, [q_1, \dots, q_n], \theta)$ **then**
- 4 $g = \text{Gen}(\ell', i', [g_1, \dots, g_n]);$
- 5 **if** $\text{Align}_0(q, g)$ **then**
- 6 $f'' \leftarrow (N \mapsto \{ N \})$
- 7 **else**
- 8 $f'' \leftarrow (N \mapsto \emptyset)$
- 9 **for** $i \in [n]$ **do**
- 10 $f_i \leftarrow \text{Inline}_{gen}(q, g_i)$
- 11 $f' \leftarrow (N \mapsto f''(N) \cup \bigcup_{i \in [n]} f_i(N))$
- 12 $\mathcal{A} = \{ i \mid i \in [D(q)] \wedge \text{Align}_i(q, g) \};$
- 13 $f \leftarrow (N \mapsto f'(N) \cup \{ \text{Ancestor}_i(N) \mid i \in \mathcal{A} \})$

The process of elimination is outlined for generated nodes in Algorithm 3. A virtually identical process is used for

matching removed nodes. The algorithm outputs a function that, given the generated replacement node (i.e., $\llbracket g \rrbracket_{\Gamma, \mu}$) that is not available until runtime, returns the set of nodes that could match the provided pattern. Matching only happens by label, as attribute values are also not available until runtime. If the pattern matches anything or if the node is reused (i.e., its label is not known until runtime), the node is a candidate for pattern match (Lines 1-2). Otherwise, the algorithm proceeds in two stages. It checks if a newly generated node can be the root of a pattern by recursively descending through the generator (Lines 3-11). Finally, it checks if any of the node's ancestors (up to the depth of the pattern) could be the root of a pattern match by recursively descending through the pattern to see if the root of the generated node could match (Lines 12-13). On lines 5 and 12, Algorithm 3 makes use of a recursive helper function: `Align`. In the base case `Align0` checks if the input pattern and generator align – whether they have equivalent labels at equivalent positions.

$$\text{Align}_0(q, g) = \begin{cases} \mathbb{T} & \text{if } q = \text{AnyNode} \vee g = \text{Reuse}(\mu) \\ \mathbb{F} & \text{if } q = \text{Match}(\ell, A, [\dots], \theta) \\ & g = \text{Gen}(\ell', i, [\dots]) \wedge \ell \neq \ell' \\ \forall k : \text{Align}_0(q_k, g_k) & \text{if } q = \text{Match}(\ell, A, [\dots], \theta) \\ & g = \text{Gen}(\ell, i, [\dots]) \end{cases}$$

The recursive case `Alignd` checks for the existence an alignment among the d th level descendants of the input pattern.

$$\text{Align}_d(q, g) = \exists k : \text{Align}_{d-1}(q_k, g)$$

Example 6.1. Continuing the running example, only the `Var` node appears in both the pattern and replacement. Thus, when a replacement is applied we need only check the parent of a replaced node for new view updates.

7 Evaluation

To evaluate `TREE_TOASTER`, we incorporated four IVM mechanisms into the `JUSTINTIMEDATA` [4, 5] compiler, a JIT compiler for data structures built around a complex AST that can *adapt* at runtime. The `JUSTINTIMEDATA` compiler naturally works with large ASTs and requires low latencies, making it a compelling use case. As such `JUSTINTIMEDATA`'s provide an infrastructure to test `TREE_TOASTER`. Our tests compare: (i) The `JUSTINTIMEDATA` compiler's existing **Naive** iteration-based optimizer, (ii) **Indexing** labels, as proposed in Section 4.1, (iii) **Classical** incremental view maintenance implemented by bolting on a view maintenance data structure created by `DBTOASTER` with the `-depth=1` flag, (iv) `DBToaster`'s full recursive view maintenance bolted onto the compiler, and (v) `TREE_TOASTER (TT)`'s view maintenance built into the compiler.

Our experiments confirm the following claims: (i) `TREE_TOASTER` significantly outperforms `JUSTINTIMEDATA`'s naive iteration-based optimizer, (ii) `TREE_TOASTER` matches or outperforms bolt-on IVM systems, while consuming significantly less memory, (iii) On complex workloads, `TREE_TOASTER`'s view maintenance latency is half of bolt-on approaches,

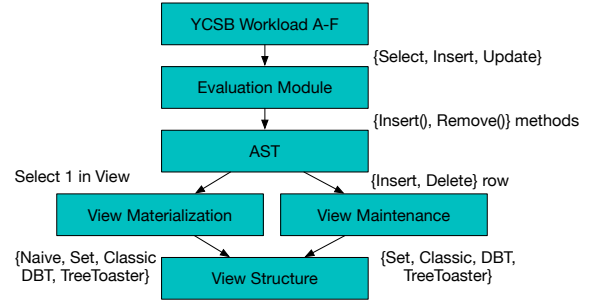


Figure 8. Benchmark Infrastructure

7.1 Workload

To evaluate `TREE_TOASTER`, we rely on a benchmark workload created by `JUSTINTIMEDATA` [5, 19], an index designed like a just-in-time compiler. `JUSTINTIMEDATA`'s underlying data structure is modeled after an AST, allowing a JIT runtime to incrementally and asynchronously rewrite it in the background using pattern-replacement rules [4] to support more efficient reads. Data is organized through 5 node types that closely mimic the building blocks of typical index structures:

- (Array, data: Seq[<key: Int, value: Int>], \emptyset)
- (Singleton, data: <key: Int, value: Int>, \emptyset)
- (DeleteSingleton, key: Int, N_1)
- (Concat, \emptyset , N_1 , N_2)
- (BinTree, sep: Int, N_1 , N_2)

`JUSTINTIMEDATA` was configured to use five pattern-replacement rules that cause it to mimic Database Cracking [16] incrementally building a tree, while pushing updates (`Singleton` and `DeleteSingleton` respectively) down into the tree.

CrackArray: This rule matches `Array` nodes and partitions them on a randomly selected pivot $\text{sep} \in \text{data}$.

$$\begin{aligned} \text{Match}(\text{Array}, [\text{data}], \emptyset, \mathbb{T}) &\rightarrow \text{Gen}(\text{BinTree}, [\text{sep}], [\\ &\quad \text{Gen}(\text{Array}, [\{x \mid x.\text{key} < \text{sep}\}], []), \\ &\quad \text{Gen}(\text{Array}, [\{x \mid x.\text{key} \geq \text{sep}\}], [])]) \end{aligned}$$

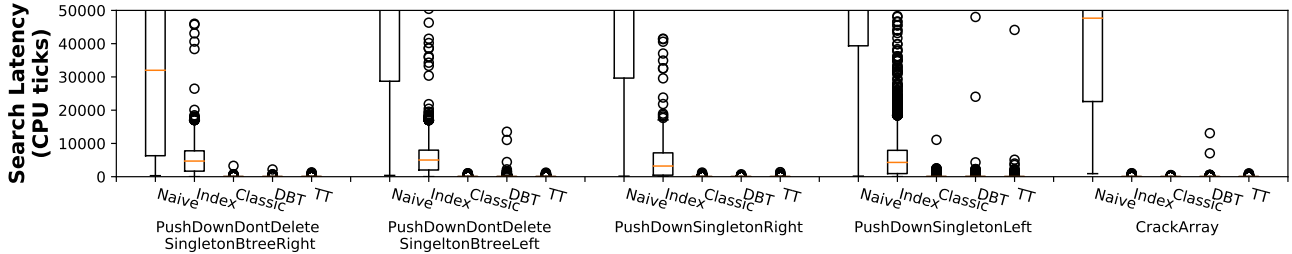
PushDownSingletonBtreeLeft/Right: These rules push `Singleton` nodes down into `BinTree` depending on the separator.

$$\begin{aligned} \text{Match}(\text{Concat}, C, [\text{Match}(\text{BinTree}, B, q_1, q_2, \emptyset), \\ \text{Match}(\text{Singleton}, S, \emptyset, \emptyset)], S.\text{key} < \text{sep}) &\rightarrow \\ \text{Gen}(\text{BinTree}, [\text{sep}], [\text{Gen}(\text{Concat}, [], [\\ &\quad \text{Reuse}(q_1), \text{Reuse}(S),], \text{Reuse}(q_2),)]) \end{aligned}$$

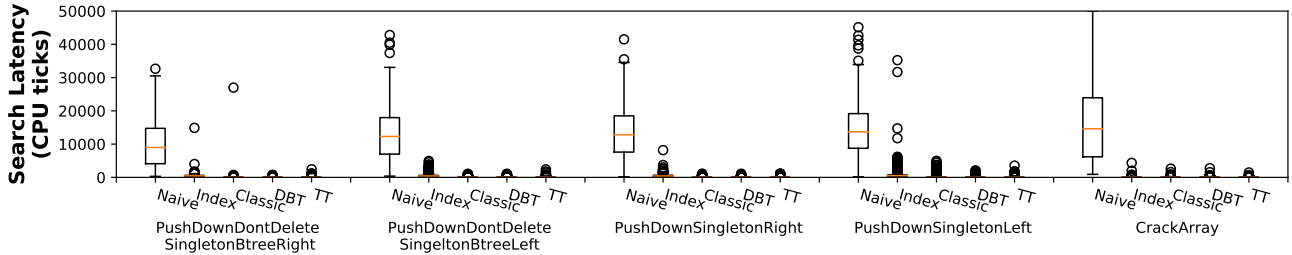
`PushDownSingletonBtreeRight` is defined analogously.

PushDownDeleteSingletonBtreeLeft/Right: These rules push `DeleteSingleton` nodes depending on the separator and are defined analogously to `PushDownSingletonBtreeLeft`.

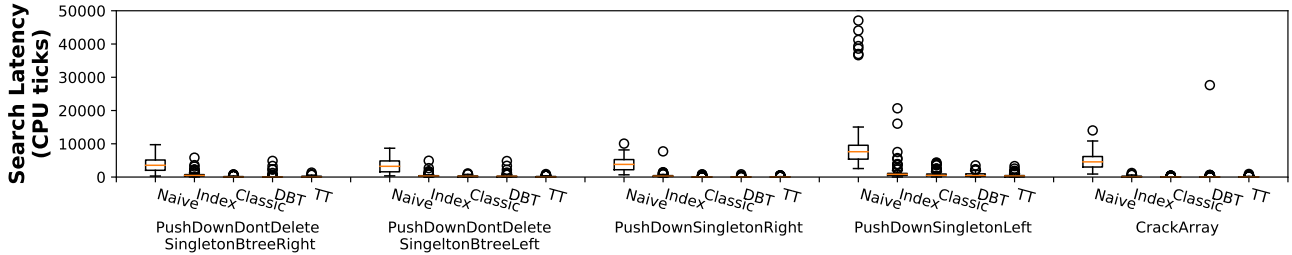
Although these rewrite rules appear relatively simple, their pattern structures are representative of the vast majority of optimizer in both Apache Spark [3] and ORCA [34].



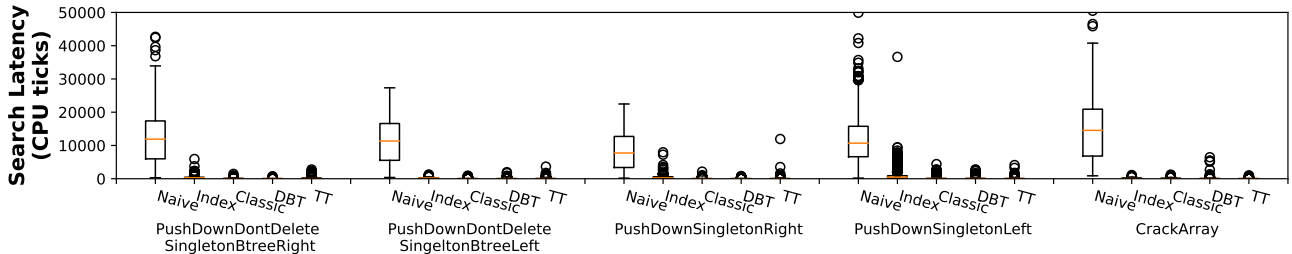
(a) Workload A



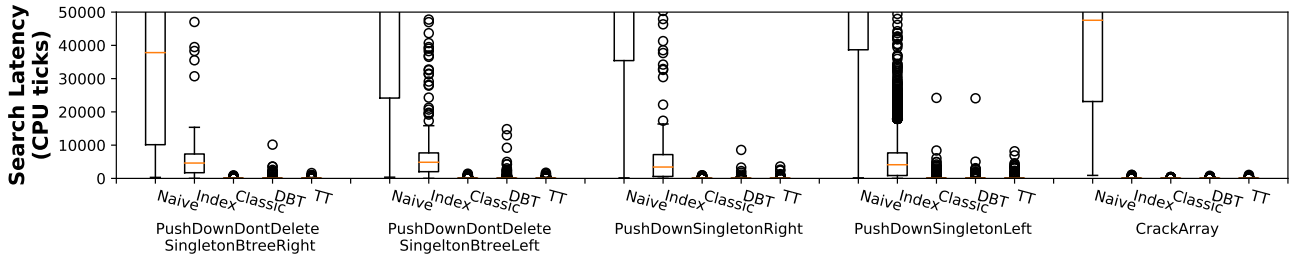
(b) Workload B



(c) Workload C



(d) Workload D



(e) Workload F

Figure 9. Relative Average Search Technique Performance by Rewrite Rule

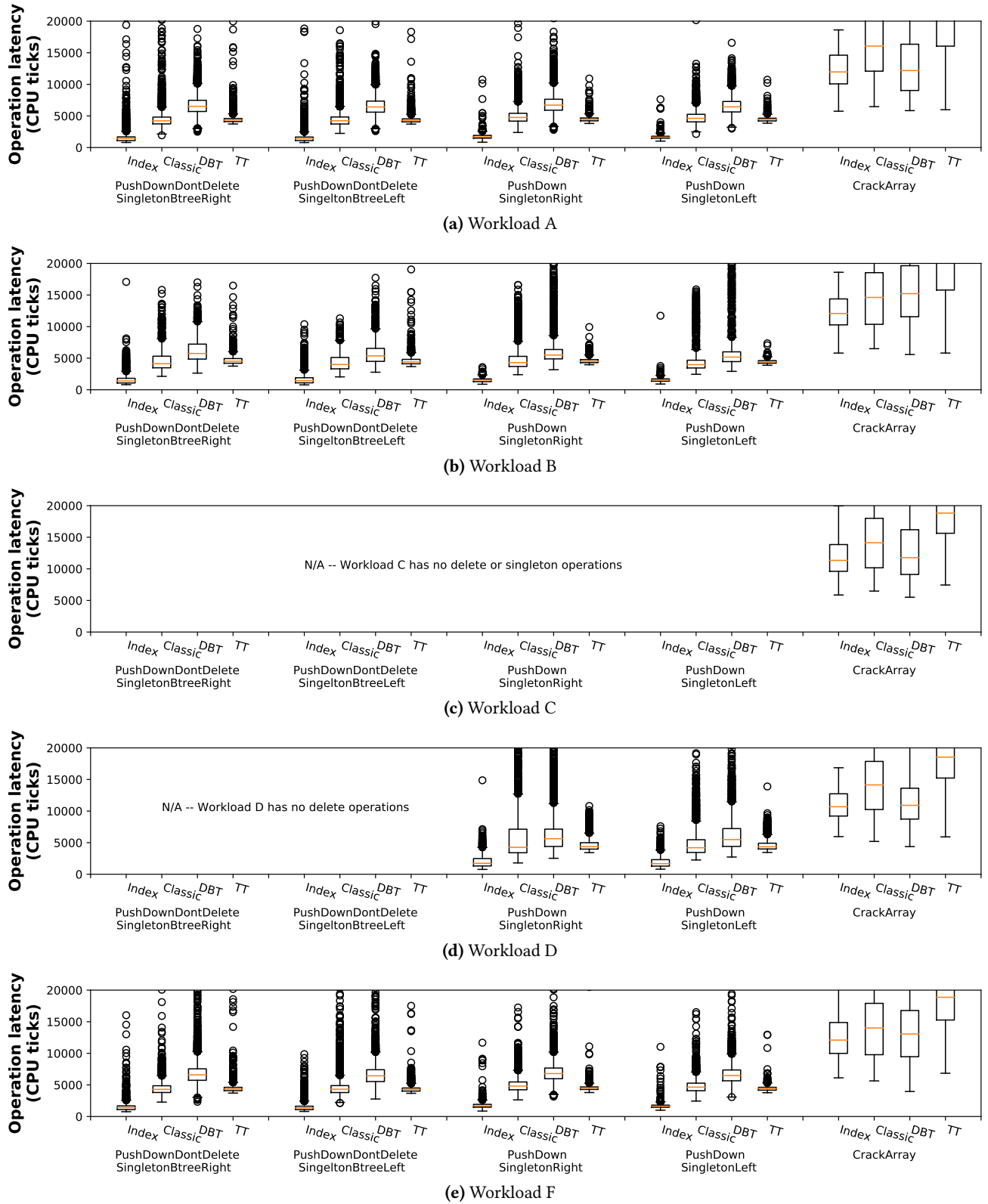


Figure 10. Relative Total Maintenance Cost by Rewrite Rule

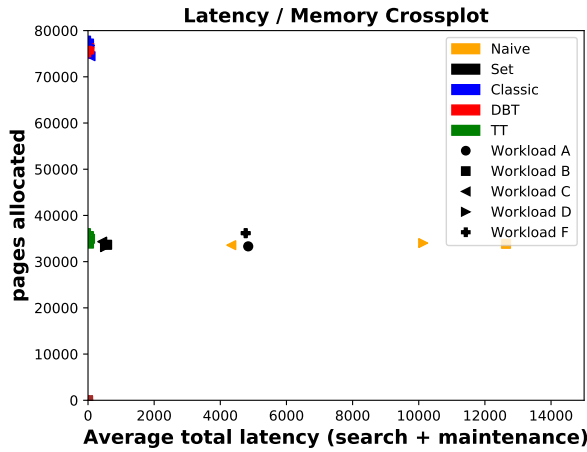


Figure 11. Total Latency (search cost + maintenance operations) and Memory Use, by method and node type

7.2 Data Gathering and Measurement

We instrumented JUSTINTIME DATA to collect updates to AST nodes as instantiations (`insert()`) and garbage collections (`remove()`) operations. To vary the distribution of optimization opportunities we used each of the six baseline YCSB [11] benchmark workloads as input to JUSTINTIME DATA. Each workload exercises a different set of node operations, resulting in ASTs composed of different node structures, patterns, and the applicability of different rewrite rules. To facilitate our experimental comparison we built a testing module in C++, allowing us to replace JUSTINTIME DATA’s naive tree traversal with each the view maintenance schemes described above for an apples-to-apples comparison. Figure 8 illustrates the benchmark generation process.

Views for TREETOASTER and label indexing were generated by declarative specification as described in Section 6 and views for DBTOASTER were generated by hand, translating rules to equivalent SQL as described in Section 2.

For data gathering we instrument TREETOASTER to collect in-structure information pertaining to view materialization and maintenance: the time required to identify potential JUSTINTIME DATA transform operations (rows in the materialized views), and the time to perform view maintenance operations upon each structure reorganization step. The test harness also records database operation latency and process memory usage, as reported by the Linux `/proc` interface. To summarize, we measure performance along three axes: (i) Time spent finding a pattern match, (ii) Time spent maintaining support structures (if any), and (iii) Memory allocated.

7.3 Evaluation

As noted in Section 7.1, JUSTINTIME DATA is configured to use 5 representative rewrite rules. Detailed results are grouped by the triggering rule. Each combination was run 10 times, with the search and operation results aggregated. We obtained

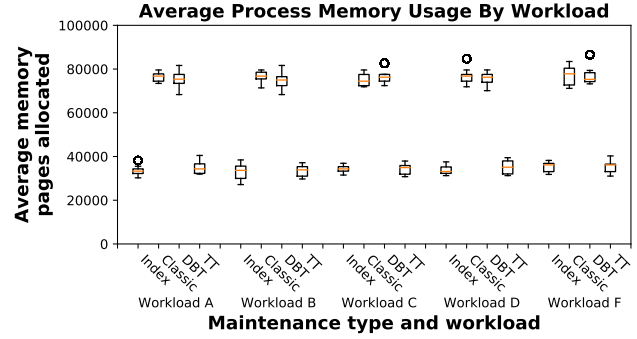


Figure 12. Average Process Memory Usage Summary.

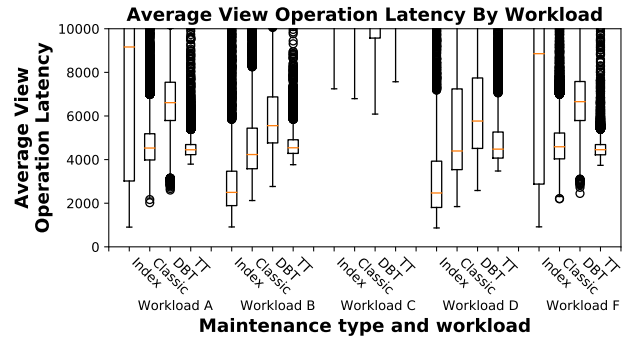


Figure 13. Average IVM operation performance Summary.

our results on a server running Ubuntu 16.04.06 LTS with 192GB RAM and 24 core Xeon 2.50GHz processors.

7.4 Results

We first evaluate how IVM performs relative to other methods of identifying target nodes in tree structure views. To do so, we compare the latency in identifying potential nodes (i.e. materializing views) using the 5 methods described above.

Figure 9 shows results obtained for several YCSB workloads using 300M keys. There were 5 sets of views that were materialized, representing target nodes in the underlying tree structure, which were candidates for each of 5 structure reorganization transforms (e.g. CrackArray). The 5 boxplot clusters compare the relative average performance of identifying 1 such node, using each of our 5 identification methods (e.g. Naive, Index). In each case, the naive search approach exhibits the worst performance. The label index approach also yields worse results than either of the IVM approaches. For identifying target nodes (i.e. emitting 1 row in the view), we thus conclude that an IVM approach performs better.

We compare TREETOASTER to IVM alternatives, including label indexing as well as a classic IVM system and a hash-based IVM implemented using DBTOASTER. Figure 11 shows the results for different loads using 300M keys. For each system, the graph shows both memory and overall performance, in terms of the combined access latency and search costs per

optimizer iteration. `TREETOASTER` easily outperforms naive iteration and has an advantage over the label index, with only a slight memory penalty relative to the former. It slightly outperforms `DBTOASTER` in general, but the total system memory for `TREETOASTER` is less than half of `DBTOASTER`.

Figure 10 shows the average total latency spent searching for a target node for a `JUSTINTIME` reorganization step, plus all maintenance steps in the reorganization, for each of the 4 IVM target node identification methods. While the label index approach starts off well (loads B and D, containing 5% writes), it scales poorly. Under increased pressure (write heavy loads A and F), average total time was significantly worse than that of `TREETOASTER`. Figure 10 also shows that, in terms of total cost, `TREETOASTER` outperforms the reference platforms: slightly better than classic IVM, and significantly better than of `DBTOASTER` IVM.

Finally, we want to measure the importance of differing approaches on system memory usage. Figure 12 shows the average memory usage in pages. For all the methods used, memory footprint did not change appreciably across time within each run. There was a small inter-run variance. Comparing across materialization and maintenance methods, both classic IVM and `DBTOASTER` IVM exhibited significantly greater memory consumption – an expected result due to its strategy of maintaining large pre-computed tables. Despite using significantly less memory to optimize performance, `TREETOASTER` performs as well as if not significantly better than these 2 alternatives. Figure 13 shows an aggregate summary across all workloads. Overall, `TREETOASTER` offers a minimum bound of both memory and latency to IVM across our evaluated alternatives.

8 Related Work

`TREETOASTER` builds on decades of work in Incremental View Maintenance (IVM) – See [10] for a survey. The area has been extensively studied, with techniques developed for incremental maintenance support for of a wide range of data models [6, 9, 32, 38] and query language features [17, 20, 23, 30].

Particularly notable are techniques that obtain performance gains through different forms of dynamic programming [21, 24, 32, 39]. Ross et. al. [32] propose materializing the set of intermediate relations in the query plan, while Koch et. al [21] propose materializing the set of intermediate relations for all possible query plans. A key feature of both approaches is computing the minimal update – or slice – of the query result, an idea also at the root of systems like Differential Dataflow [24]. Both approaches show significant performance gains on general queries. However as we previously discussed, the sources of these gains: selection-pushdown, aggregate-pushdown, and cache locality are all less relevant in the context of abstract syntax trees. Similarly-spirited approaches can be found in other contexts, including graphical inference [39], and fixed point computations [24].

Also relevant is the idea of embedding query processing logic into a compiled application [2, 15, 21, 25, 27, 29, 31, 33]. Systems like BerkeleyDB, SQLite, and DuckDB embed full query processing logic, while systems like `DBToaster` [2, 21, 27] and `LinQ` [25] compile queries along with the application, making it possible to generate native code optimized for the application. Most notably, this makes it possible to aggressively inline SQL and imperative logic, often avoiding the need for boxing, expensive VM transitions for user-defined functions, and more [25, 33, 37]. Major database engines have also recently been extended to compile queries to native code [8, 12, 26], albeit at query compile time.

To our knowledge, IVM over Abstract Syntax Trees specifically has not been studied directly. A related effort, the `Cascades` framework [14] considers streamlined approaches to scheduling rule application, a strategy that is used by the `Orca` [34] compiler. There are also several related efforts in the more general category of IVM for general tree and graph query languages and data models like XPath [13], Cypher [35, 36], and the Nested Relational Calculus [22]. These schemes address recursion, with which join widths are no longer bounded; as well as aggregates without an abelian group representation (e.g., min/max), for which supporting deletion efficiently is more difficult.

However, two approaches aimed at the object exchange model [1, 40], are very closely related to our own approach. One approach proposed by Abiteboul et. al. [1] first determines the potential positions at which an update could affect a view and then uses the update to recompute the remaining query fragments. However, its more expressive query language limits optimization opportunities, creating situations where it may be faster to simply recompute the view query from scratch. The other approach proposed by Zhuge and Molina [40] follows a similar model to our immutable IVM scheme, enforcing pattern matches on ancestors by recursively triggering shadow updates to all ancestors.

9 Conclusion

In this paper we introduce a formalized mechanism for pattern-matching queries over ASTs and IVM over such queries. Our realization of the theory in the just-in-time data-structure compiler highlights the viability of our system as compared to tradition IVM approaches.

Our future work includes extending our approach to work with graphs. This will allow for recursive pattern matches and efficient IVM over graph structures. This is particularly interesting for traditional compilers as there exist many optimizations that rely on fixed-points while traversing control-flow graphs. Similarly, we plan on integrating our approach into compiler, like GCC, which provides a number of AST based optimization passes.

References

- [1] Serge Abiteboul, Jason McHugh, Michael Rys, Vasilis Vassalos, and Janet L. Wiener. 1998. Incremental Maintenance for Materialized Views over Semistructured Data. In *VLDB*. Morgan Kaufmann, 38–49.
- [2] Yanif Ahmad, Oliver Kennedy, Christoph Koch, and Milos Nikolic. 2012. Dbtoaster: Higher-order delta processing for dynamic, frequently fresh views. *arXiv preprint arXiv:1207.0137* (2012).
- [3] Michael Armbrust, Reynold S. Xin, Cheng Lian, Yin Huai, Davies Liu, Joseph K. Bradley, Xiangrui Meng, Tomer Kaftan, Michael J. Franklin, Ali Ghodsi, and Matei Zaharia. 2015. Spark SQL: Relational Data Processing in Spark. In *SIGMOD Conference*. ACM, 1383–1394.
- [4] Darshana Balakrishnan, Lukasz Ziarek, and Oliver Kennedy. 2019. Fluid data structures. In *DBPL*. ACM, 3–17.
- [5] Darshana Balakrishnan, Lukasz Ziarek, and Oliver Kennedy. 2019. Just-in-Time Index Compilation. *arXiv preprint arXiv:1901.07627* (2019).
- [6] José A. Blakeley, Per-Åke Larson, and Frank Wm. Tompa. 1986. Efficiently Updating Materialized Views. In *SIGMOD Conference*. ACM Press, 61–71.
- [7] Wayne D. Blizard. 1990. Negative Membership. *Notre Dame J. Formal Log.* 31, 3 (1990), 346–368.
- [8] Dennis Butterstein and Torsten Grust. 2016. Precision Performance Surgery for PostgreSQL: LLVM-based Expression Compilation, Just in Time. *Proc. VLDB Endow.* 9, 13 (2016), 1517–1520.
- [9] Surajit Chaudhuri, Ravi Krishnamurthy, Spyros Potamianos, and Kyuseok Shim. 1995. Optimizing Queries with Materialized Views. In *ICDE*. IEEE Computer Society, 190–200.
- [10] Latha S. Colby, Timothy Griffin, Leonid Libkin, Inderpal Singh Mumick, and Howard Trickey. 1996. Algorithms for Deferred View Maintenance. In *SIGMOD Conference*. ACM Press, 469–480.
- [11] Brian F Cooper, Adam Silberstein, Erwin Tam, Raghu Ramakrishnan, and Russell Sears. 2010. Benchmarking cloud serving systems with YCSB. In *Proceedings of the 1st ACM symposium on Cloud computing*. 143–154.
- [12] Databricks. 2015. Project Tungsten. <https://databricks.com/glossary/tungsten>. (2015).
- [13] Katica Dimitrova, Maged El-Sayed, and Elke A. Rundensteiner. 2003. Order-Sensitive View Maintenance of Materialized XQuery Views. In *ER (Lecture Notes in Computer Science, Vol. 2813)*. Springer, 144–157.
- [14] Goetz Graefe. 1995. The Cascades Framework for Query Optimization. *IEEE Data Eng. Bull.* 18, 3 (1995), 19–29.
- [15] D. Richard Hipp. 2000. SQLite: Small. Fast. Reliable. Choose any three. <https://sqlite.org/>.
- [16] Stratos Idreos, Martin L. Kersten, and Stefan Manegold. 2007. Database Cracking. In *CIDR*. www.cidrdb.org, 68–78.
- [17] Akira Kawaguchi, Daniel F. Liuwen, Inderpal Singh Mumick, and Kenneth A. Ross. 1997. Implementing Incremental View Maintenance in Nested Data Models. In *DBPL (Lecture Notes in Computer Science, Vol. 1369)*. Springer, 202–221.
- [18] Oliver Kennedy and Lukasz Ziarek. 2015. Just-In-Time Data Structures. In *CIDR*. www.cidrdb.org.
- [19] Oliver Kennedy and Lukasz Ziarek. 2015. Just-In-Time Data Structures. In *CIDR*. Citeseer.
- [20] Christoph Koch. 2010. Incremental query evaluation in a ring of databases. In *PODS*. ACM, 87–98.
- [21] Christoph Koch, Yanif Ahmad, Oliver Kennedy, Milos Nikolic, Andres Nötzli, Daniel Lupei, and Amir Shaikhha. 2014. DBToaster: higher-order delta processing for dynamic, frequently fresh views. *VLDB J.* 23, 2 (2014), 253–278.
- [22] Christoph Koch, Daniel Lupei, and Val Tannen. 2016. Incremental View Maintenance For Collection Programming. In *PODS*. ACM, 75–90.
- [23] Per-Åke Larson and Jingren Zhou. 2007. Efficient Maintenance of Materialized Outer-Join Views. In *ICDE*. IEEE Computer Society, 56–65.
- [24] Frank McSherry, Derek Gordon Murray, Rebecca Isaacs, and Michael Isard. 2013. Differential Dataflow. In *CIDR*. www.cidrdb.org.
- [25] Erik Meijer, Brian Beckman, and Gavin M. Bierman. 2006. LINQ: reconciling object, relations and XML in the .NET framework. In *SIGMOD Conference*. ACM, 706.
- [26] Thomas Neumann. 2011. Efficiently Compiling Efficient Query Plans for Modern Hardware. *Proc. VLDB Endow.* 4, 9 (2011), 539–550.
- [27] Milos Nikolic, Mohammad Dashti, and Christoph Koch. 2016. How to Win a Hot Dog Eating Contest: Distributed Incremental View Maintenance with Batch Updates. In *SIGMOD Conference*. ACM, 511–526.
- [28] Patrick E. O’Neil, Edward Cheng, Dieter Gawlick, and Elizabeth J. O’Neil. 1996. The Log-Structured Merge-Tree (LSM-Tree). *Acta Informatica* 33, 4 (1996), 351–385.
- [29] Oracle. 1994. Oracle BerkeleyDB. <https://www.oracle.com/database/berkeley-db/>.
- [30] Themistoklis Palpanas, Richard Sidle, Roberta Cochrane, and Hamid Pirahesh. 2002. Incremental Maintenance for Non-Distributive Aggregate Functions. In *VLDB*. Morgan Kaufmann, 802–813.
- [31] Mark Raasveldt and Hannes Mühleisen. 2019. DuckDB: an Embeddable Analytical Database. In *SIGMOD Conference*. ACM, 1981–1984.
- [32] Kenneth A. Ross, Divesh Srivastava, and S. Sudarshan. 1996. Materialized View Maintenance and Integrity Constraint Checking: Trading Space for Time. In *SIGMOD Conference*. ACM Press, 447–458.
- [33] Amir Shaikhha. 2013. An Embedded Query Language in Scala. <http://infoscience.epfl.ch/record/213124>
- [34] Mohamed A. Soliman, Lyublena Antova, Venkatesh Raghavan, Amr El-Helw, Zhongxian Gu, Entong Shen, George C. Caragea, Carlos Garcia-Alvarado, Foyzur Rahman, Michalis Petropoulos, Florian Waas, Sivaramakrishnan Narayanan, Konstantinos Krikellas, and Rhonda Baldwin. 2014. Orca: a modular query optimizer architecture for big data. In *SIGMOD Conference*. ACM, 337–348.
- [35] Gábor Szárnyas. 2018. Incremental View Maintenance for Property Graph Queries. In *SIGMOD Conference*. ACM, 1843–1845.
- [36] Gábor Szárnyas, József Marton, János Maginecz, and Dániel Varró. 2018. Reducing Property Graph Queries to Relational Algebra for Incremental View Maintenance. *CoRR* abs/1806.07344 (2018).
- [37] Thomas Würthinger. 2014. Graal and truffle: modularity and separation of concerns as cornerstones for building a multipurpose runtime. In *MODULARITY*. ACM, 3–4.
- [38] Jun Yang and Jennifer Widom. 2003. Incremental computation and maintenance of temporal aggregates. *VLDB J.* 12, 3 (2003), 262–283.
- [39] Ying Yang and Oliver Kennedy. 2017. Convergent Interactive Inference with Leaky Joins. In *EDBT*. OpenProceedings.org, 366–377.
- [40] Yue Zhuge and Hector Garcia-Molina. 1998. Graph Structured Views and Their Incremental Maintenance. In *ICDE*. IEEE Computer Society, 116–125.