
FluidDB: Adaptive storage layout using reversible relational operators

Christos Perivolaropoulos



Doctor of Philosophy

THE UNIVERSITY OF EDINBURGH

January 2022

If the scientists, brought to heel by self-interested rulers, limit themselves to piling up knowledge for knowledge's sake, then science can be crippled and your new machines will lead to nothing but new impositions. You may in due course discover all that there is to discover, and your progress will nonetheless be nothing but a progress away from mankind. The gap between you and it may one day become so wide that your cry of triumph at some new achievement will be echoed by a universal cry of horror. Life of Galileo – Bertolt Brecht

Abstract

It is a popular practice to use materialized intermediate results to improve performance of RDBMSes. Work in this area has focused either on optimizers matching existing results or selecting useful intermediate results from a plan, but few attempts have been made to create plans with intermediate results in mind, and none that make any deduplicate the stored data to alleviate the storage cost of maintaining possibly large queries.

We built *FluidB* to explore a novel approach to integrating the selection of materialized results with the planner in order to optimize the logical representation of data in memory. FluidB materializes hot intermediate results and deduplicates data to alleviate the cost of maintaining them. This is achieved by introducing *reversible operations*, versions of normal relational operators that optionally pass enough data to the output to make the input relations reconstructable. A planner aware of such operations can build query plans that dynamically adapt the data layout to the plan under constrained memory budget. This thesis revolves around four main chapters each of which describes in detail a different part of FluidB and a final one that goes into evaluation of the system.

The first chapter focuses on query processing and the relational algebra semantics that FluidB operates under. FluidB parses queries into graphs of sub-queries connected by reversible operators. Each such graph of the workload is merged into a global query graph that is used to infer properties of each relation like cardinality and extent.

The next chapter is dedicated to the planner and a novel monad for weighted backtracking that the planner is based on. The planner attempts to generate a plan based on the query graph that besides solving the query at hand, leaves in memory an optimal set of queries for the workload being run. In this chapter, the garbage collector is also discussed, that creates deletion operators as part of the plan such that the available budget is respected.

After that, we go into *Antisthenis*, a novel incremental computation system used by the query planner to determine, given a particular set of materialized relations, whether a relation materializable and the expected cost of materializing a relation. Antisthenis, besides reusing computations, takes advantage of the properties of the operations involved in the expression she is evaluating, like absorbing group elements, to heuristically avoid as much work as possible.

The final chapter about the FluidDB architecture describes the transpilation of plans generated by the planner to C++, as well as the supporting libraries that enable the evaluation of queries as C++ code, and the low level data organization of the database. The thesis closes with a chapter that describes our methods for benchmarking and some experimental results.

Lay Summary

TODO: lay summery

Acknowledgements

This thesis is dedicated to all the people that directly or indirectly supported me through my PhD studies, but also to all the people without whose support and hard work I would never have started this work to begin with.

It is customary to thank the supervisors first regardless of the circumstances but I would like to underline my gratitude to Stratis who really understood my idiosyncratic way of working and provided precisely the quality and quantity of guidance that I required. This is a rare quality and I cherish the opportunity I had to work with him. I would also like to express my gratitude to thank Boris Grot, not just in the general sense but because, especially during the first years of my thesis he went out of his way to make me feel included and gave value to my voice. For the similar reasons I would like to thank Christophe Dubach whose door was always open and his whiteboard always available, I deeply appreciate all the times he was willing to patiently listen to my crazy ideas.

I would like to thank my family and especially my parents, for being with me every step of the bumpy road that was my PhD, but mostly because without their hard work and sacrifices a PhD would not have been an option for me. I am privileged and humbled and I owe it primarily to them.

I would also like to thank Marina without whom I would definitely have given up in one of the many moments of desperation. Thank you for helping me with the impossible task of finishing a PhD during a pandemic.

Special thank you to my flatmate for two years Christos Velanis who got me safe and sound through some fairly dark times, and for the ideological steeling brought about by this cohabitation. Thanks to him and to Christos Maniatis for the heated philosophical debates that would typically end after 6am.

Last but not least I would like to thank my friends that were there through all of it, the

thick and the thin Antreas, Giannis, Aisha, Maria, Irini, Iordanis, Nelly, Christos, Pigi, Niki, The President, and Nikolas.

Declaration

I declare that this thesis was composed by myself, that the work contained herein is my own except where explicitly stated otherwise in the text, and that this work has not been submitted for any other degree or professional qualification except as specified.

Christos Perivolaropoulos

Contents

Abstract	iii
Lay Summary	v
Acknowledgements	vii
Declaration	ix
List of Figures	xv
1. Introduction	1
2. Background	11
2.1. Relational databases	12
2.2. Query optimization and planning	13
2.2.1. Logical and physical query optimization	15
2.2.2. Cost estimation	16
2.3. In-memory relational databases	17
2.3.1. Transpilation	19
2.3.2. Third party JIT compilers	19
2.3.3. Direct machine code generation	19
2.3.4. Custom execution engines	20
2.4. Intermediate result recycling	20
2.5. Conclusion	21

3. Logical planning	23
3.1. Relational Algebra	24
3.1.1. Expressions	26
3.2. QDAG: The unified query graph	28
3.2.1. Clusters and irreversible edges	32
3.3. Deterministic query processing	35
3.3.1. Query parsing	36
3.3.2. Query preprocessing	37
3.3.3. Possible joins	39
3.4. Query Normal Form (QNF)	42
3.4.1. QNF Structure	43
3.4.2. QNF Computation	49
3.4.3. Pragmatism of QNFs	52
3.4.4. Symbol correspondence	52
3.5. Query clusters internals	55
3.5.1. Cluster polymorphism	55
3.5.2. Populating the graph with clusters	57
3.6. Relation shapes	58
3.6.1. Shape propagators	59
3.6.2. Defaulting functor	61
3.7. Summary and conclusion	64
4. Physical planning	65
4.1. HCntT logic monad	66
4.1.1. Background	66
4.1.2. The HCntT monad	70
4.2. The planner	77
4.2.1. Garbage collection	81
4.2.2. Order of traversal	85
4.3. Conclusion	88
5. Antisthenis	89
5.1. Introduction	91
5.2. Background	95
5.2.1. Arrows and profunctors	95

5.2.2. Mealy arrows: processes that remember	99
5.2.3. Zippers	102
5.3. Antisthenis core	103
5.3.1. Processes	104
5.3.2. Epochs and coepochs	106
5.3.3. Antisthenis caps and bounds	108
5.3.4. Antisthenis zipper	110
5.3.5. The Cmds functor	114
5.3.6. Resetting processes	117
5.4. General operators	118
5.4.1. Basic cost estimation operators	120
5.5. Node machines	125
5.5.1. Antisthenis machine tape	125
5.5.2. Process stack	126
5.5.3. Cyclical machines	126
5.5.4. Calculating historical query costs	130
5.6. Conclusion	133
6. Execution engine	135
6.1. Codegen introduction	136
6.2. Storage memory management	136
6.3. Data layout	137
6.3.1. Initial data conversion	138
6.3.2. Pages	142
6.3.3. C++ row iterators	143
6.4. Physical planning	143
6.5. Operator implementations	144
6.5.1. Generated code structure	145
6.5.2. Select algorithm	150
6.5.3. Projection algorithm	151
6.5.4. The join operator	152
6.5.5. Aggregation and sort algorithms	156
6.6. Conclusion	156

Contents

7. Evaluation	157
7.1. A note on size estimation	164
7.2. Conclusion	166
8. Conclusions and future perspectives	167
8.1. Future perspectives	169
A. Appendix	171

List of Figures

1.1. An efficient logical plan for a single query.	3
1.2. An efficient logical plan for a template based workload.	4
1.3. A sequence of plans optimizing the workload amortized cost and involving reverse operations. These queries are the same save for the predicate p_i selecting <i>lineitem</i>	6
1.4. The materialized relations are marked with grey. Assuming the absense of null values and assuming set semantics $lineitem = \sigma_p lineitem \cup \sigma_{\neg p} lineitem$, FluidDB can find a plan to generate any of the relations in the graph from $\sigma_p lineitem$ and $\sigma_{\neg p} lineitem$	8
1.5. A sequence of plans optimizing the workload amortized cost and involving reverse operations.	9
2.1. Common architecture of a query optimizer.	22
3.1. A selection t-node may generate both output nodes, $\sigma_p A$ and $\sigma_{\neg p} A$, rendering it a partition operation.	30
3.2. Multiple selections can be composed to create larger networks.	31
3.3. Representing a reversible join as a single t-node.	33
3.4. The QDAG corresponding to a join is broken down in multiple t-nodes to facilitate greater flexibility.	34
3.5. The QDAG corresponding to a join is broken down in multiple t-nodes to facilitate greater flexibility.	35
3.6. There are four cluster types. The NClust is trivial, just a single node. The other 3 are BinClust , JoinClust , and UnClust representing binary, join, and unary operations respectively.	56

3.7.	A QDAG accumulated out of every node in the SSB TPC-H.	58
4.1.	Marker "branches" are injected in the heap of prioritized branches to possibly trigger the fallback branch. Each marker gives the scheduler the opportunity to run a particular fallback branch. The entry of the fallback branch entry references the last marker so that only that one actually trigger the fallback branch. A child branch that succeeds removes the fallback entry so the final marker also fails to trigger the fallback.	76
4.2.	starting at node O , the output of a join operation, we can derive three MetaOps that can materialize it. $MetaOp\{in : \langle I_l, I_r \rangle, out : \langle O \rangle\}$, $MetaOp\{in : \langle I_l, I_r \rangle, out : \langle O_l, O \rangle\}$, $MetaOp\{in : \langle I_l, I_r \rangle, out : \langle O, O_r \rangle\}$, $MetaOp\{in : \langle I_1, I_2 \rangle, out : \langle O_l, O, O_r \rangle\}$. Because trying all combinations of outputs explodes the search space we always go for the largest and then let the garbage collector deal with the possible repercussions. On the other hand to materialize I_l there is only one MetaOp that relates to this cluster $MetaOp\{in : \langle O, O_l \rangle, out : \langle I_l \rangle, interm : \langle I_l \times I_r \rangle\}$	80
5.1.	A computation DAG. Circular nodes represent operations and square nodes represent parameters that are subject to change.	92
5.2.	The zipper splits the sub-processes that cooperate to compute values for the owning process in four categories a) initials that have not yielded a currently valid value b) iterations that have not yielded a valid partial value c) coiterations that have yielded a full valid value d) the cursor that is the next sub-process to be evaluated. Depending on the evaluation strategy of the operator once the cursor process evaluates it is replaced with an init process or an iteration process, and depending on the value the cursor process yields it is pushed to the iterations or the coiterations of the zipper.	112

5.3.	A tree of different evolution paths a zipper may undergo during between the calling of a process and the accumulated bound exceeding the cap. p_0 , p_4 , p_3 and p_4 are subprocesses. The values returned, the specific iteration subset data structure are not presented for brevity and to shift focus to the movement of the subprocesses within the zipper. Nodes in this tree that have two children correspond to the ItInit constructor, nodes with one node correspond to the CmdIt or CmdInit constructor depending on which process set the process in the nex cursor is drawn from. When a concrete value is reached or the cap is exceeded the process stops and the free monad that wraps Cmnds takes the value Pure	117
5.4.	Node C is recomputed because the \perp value is predicated on A being in the computation stack.	127
7.1.	The foreign key links in a SSB-TPC-H	159
7.2.	The foreign key links in a traditional TPC-H schema	159
7.3.	The total page read/writes for each query a sequence of queries. The baseline is the query being run by FluidB directly without any materialized nodes. The workload bars represent the cost of each query in a workload being accumulated into the same QDAG.The budget allowed for this is the minimum budget within which FluidB can run each individual query (2300K pages).	160
7.4.	The total page read/writes for each query a sequence of queries. The baseline is the query being run by FluidB directly without any materialized nodes. The workload bars represent the cost of each query in a workload being accumulated into the same QDAG.The budget allowed for this is about triple the minimum budget within which FluidB can run each individual query (6100K pages). A badly timed GC run causes some individual queries to be slower than their coutnerparts from a workload being run in a tight budget (figure 7.3.	162

List of Figures

- 7.5. The total page read/writes for each query a sequence of queries. The baseline is the query being run by FluidB directly without any materialized nodes. The workload bars represent the cost of each query in a workload being accumulated into the same QDAG.. The budget allowed for this is about triple the minimum budget within which FluidB can run each individual query (6500K pages). 165

List of Listings

3.1. Types of algebraic expressions that compose the non-RA parts of the queries.	28
3.2. This parser returns a query modifier that is meant to be applied to a very simple product query generated by the from clause. The actual selection parser is complex and handles many different cases, but it is built up from simple fundamental blocks.	37
3.3. Pseudo-python description of finding all possible for clarity, it is abbreviated to omit sanity checking, memoization, some type conversions, etc. (This algorithm does not include memoization for brevity)	42
3.4. The QNF datastructure.	44
3.5. A QNF name may be an unnamed column of a relation, a named column of a primary table or a literal.	46
3.6. A QNF projection field is a collection of expressions that refer to QNF names. The particular structure of this collection is parametric. When the collection Identity the QNF query is essentially just a column. A normal QNF query would instantiate f to HashBag , an unordered multiset.	46
3.7. The QNF aggregation form of the projection field is similar to projection only, much like the QGroup constructor, it also includes a HashSet of expressions on which to group.	47
3.8. Selection name refers to a version of the current QNF that has all fields erased except the projection.	47
3.9. A key that uniquely identifies a QNF is made of three separate hashes, one for each part of the qnf structure so that they can be updated independently.	48

List of Listings

3.10. A fast and loose definition of equality between QNFs that depends on the quality of equality.	48
3.11. A very inefficient but correct equality between QNFs.	49
3.12. QNF computation monad provides non-determinism, caching, and error handling.	50
3.13. The internal QNF building functions provide some more information that was created during the generation of the QNF, precisely a name map relating column names to QNF names, a map relating input QNF names to output QNF names, and the top level operator with the names translated appropriately to input or output QNF names.	51
3.14. The qnf product.	52
3.15. The definition of the query forest. The query forest is hashed so that we can avoid traversing the same query forest repeatedly. The query forest is essentially a non-empty of queries with forests at their leafs.	57
3.16. A propagator matches a cluster with shapes at the edges to the same kind of cluster with the shapes synchronized.	60
3.17. The defaulting functor definition.	62
3.18. The join semilattice that is defined in terms of the	63
3.19. Promoting of defaulting functor happens during code generation when an n-node is materialized.	63
4.1. The logic monad typeclass	67
4.2. Using a simple list to drive non-determinism is implicitly equivalent to a DFS algorithm which in many useful cases does not terminate.	67
4.3. Interleaving (in this example >*<) is not <i>actually</i> fair in the sense that it does not give all the processes	68
4.4. Prioritise branches that we want to be executed first.	69
4.5. The HCntT monad transformer allows continuation based non-determinism that allows switching between branches.	70
4.6. We parameterize over heaps to allow the user to decide an efficient priority queue for the branches.	70
4.7. Disolution is the process of turning an HCntT computation into a ListT	71
4.8. The dissolution algorithm in pseudo-python	72
4.9. The implementation for Alternative is the same as the implementation for MonadPlus	72

4.10. The halt process yields updates the priority of the branch and yields execution to the scheduler.	73
4.11. This computation will evaluate to 1 because, while the computation <code>return 0</code> always succeeds the branch fails.	74
4.12. The nested scheduler runs a subprocess within a single branch. Once is built on top of that to make sure the process stops once a result is returned.	77
4.13. The monad that defines all the useful effects used by the planner. GCConf <code>t n</code> is an immutable, from the perspective of the planner, configuration that includes the QDAG, the node sizes, etc. GCState <code>t n</code> is state that is mutated and private to each branch of the planner like the materialized status of the nodes, the set of transitions registered so far and various caches. The PlanningError <code>t n</code> is a planner specific type of error. The entirety of the result of planning is accumulated in GCState so the result of backtracking is just unit <code>(())</code>	78
4.14. A MetaOp refers to input, output, and intermediate nodes that are involved in the set of operations it abstracts. Furthermore it contains a computation that registers and returns the transitions involved in the MetaOp	79
4.15. The different states that a node is allowed to be in.	81
4.16. A sketch of the garbage collector algorithm in pseudo-haskell. The <code><///></code> operator is supported by the HCntT	85
5.1. Haskell typeclasses related to the notion of Arrow	96
5.2. Laws for the typeclasses related to arrows.	97
5.3. Haskell definition of a mealy arrow.	100
5.4. A MealyArrow can take on the properties of other arrows by swapping out the function <code>→</code> type for a parametric one.	100
5.5. The MB <code>a b m</code> monad can be used as a convenience to implement mealy arrows using a coroutine-like interfaces.	101
5.6. An example the usage of the MB functor to generate mealy arrows.	101
5.7. Definition of the MB monad transformer.	102
5.8. The type of an Antisthenis process is a mealy arrow paired with a kleisli arrow.	104
5.9. The definition of the return value of an Antisthenis process. It may be a final result, an error indicating that a final result is non-computable, or a bound for the final value.	105

5.10. The type definition of a configuration. It is a tuple containing information that can be used to derive whether a value is valid.	106
5.11. The type of a naive function checking the validity of a value based on epoch and coepoch.	107
5.12. The final function for epoch and coepoch combination.	108
5.13. A sample implementatation of the function that transforms the configuration received by a parent process into one suitable for the child process. Checks if the parent process needs to be reset and uses the partial result to constrain the cap.	110
5.14. The definition of the zipper.	114
5.15. An internal representation of the process evolving the internal representation of a process: the zipper.	115
5.16. Definition of the commands functor that provides different branches of evolution for zipper.	115
5.17. A function that each atnistenis opator needs to implement that decides the traversal of the tree created by the different possible evolutions of zipper. The implementation may also deem that a good reset point has been discovered.	116
5.18. Operator specific types that need to be implemented by every operator.	118
5.19. Tags are phantom types that define the Antisthenis operations. The tags themselves may be parameterized using type parameter p. For the needs of FluidB we define operations for addition, subtraction, multiplication and boolean operations.	119
5.20. The type w fully defines the operator so combining subprocesses into a process is unambiguous.	119
5.21. An object of type Conv $w \ w'$ can act as an interface between parent processes of type w to subprocesses of type w'	120
5.22. Pseudocode for the algorithm of evaluating a process up to a threshold defined by the cap. For brevity and to avoid too much unnecessary detail we omit sanity checks and the reset handling code.	122
5.23. The algorithm for creating predicates.	129
5.24. Algorithm for handling the predicates.	130
5.25. Definition of the type used for capping the cost of historical queries.	131
5.26. Definition of the type used for both bounds and results of historical queries.	132

5.27. Comparison between bounds and between bound and cap are different. Between bounds we need to account for the computability metric. The cap on the other hand defines a three-dimensional bound () that the bound must fall within in order to not exceed it.	133
6.1. For standard FFI communication C++ structs that do not contain fancy constructors	138
6.2. The supplier row representation in the generated C++ code. The <code>fluidb_string</code> type is a constant size array of characters.	141
6.3. Algorithm to infer the padding of members according to the Itanium ABI.	141
6.4. Convert a bama file to a paged data file.	142
6.5.	143
6.6. Using a <code>RecordMap</code> to sort the records of a file.	143
6.7. The type level maybe	146
6.8. The shape of a generated predicate type.	147
6.9. A record transformation type defines objects with no internal state that are callable.	148
6.10. The selection operator. It is parameterized by the predicate and the primary and secondary output types. Enough information about these values is known at compile time such that the compiler can generate highly specialized code.	149
6.11. The C++ declaration of the select.	150
6.12. A block representing a particular operator.	150
6.13. The projection operator produced from the SQL query <code>select p_partkey, p_color from part.</code>	152
6.14. The equi-join algorithm first sorts in place the inputs w.r.t. the equal subtuples and then merges them.	155
6.15. Class declaration of the join operator	155
7.1. Abbreviated version of the plans of queries 4 to 7 of SSB TPC-H. This demonstrates how an unfortunately timed GC can cause cause FluidDB to make some bad decisions	163
A.1. The SQL code for the SSB queries	175

Introduction

Something about the
beginning of a journey.

(Author)

Chapter summary

- FluidDB is an in-memory RDBMs optimizes data layout for space efficiency w.r.t. the workload
- The main novelty relates to the introduction of reversible relational operations which affords a new perspective on query planning and view selection.
- She materializes all intermediate nodes and deletes garbage collects when she runs out of space.

With the advent of technologies that make access to information scalable and affordable, the mental and temporal gap between collection of data and their analysis grows rapidly. At least two of the biggest players in the tech industry, Facebook and Google, base their competitive advantage almost exclusively on vast amounts of information that they have collected and their capacity for such collection and processing. In cases like these the layout of the stored data is independent of the growing number of applications taking advantage of it.

A mantra in database community used to be that “storage is cheap” and while that is true a more complete version of the mantra might be that “slow storage is cheap”.

Databases have traditionally been dealing with the tradeoff between memory and time efficiency within monetary constraints especially with the use of intermediate result recycling technologies which employ sophisticated ways of choosing parts of computation to be stored for reuse. This model has had some impressive results in OLAP workloads. We adopt a slightly different view of the problem: executing query plans does not only leave some specific intermediate results as residue, but rather *transitions the entire storage state from one where the result of the query is not materialized to one where it is*. With such a notion of query planning what new dimensions open up in the design space of a query planner? The creation of FluiDB is an attempt to study some aspects of this question. In particular FluiDB is based on the following pillars:

- Introduce reversible query operations that allow for more sophisticated plans based on available materialized views.
- The planner itself involves a garbage collector that will delete materialized views or primary tables that can be materialized from the remaining relations.
- An incremental numeric evaluation system allows the planner to efficiently and repeatedly infer the cost of materializing queries under a rapidly changing inventory of materialized views.
- Query execution is based on code generation.

These concepts allow FluiDB to dynamically adapt the data layout to the workload in ways that would not be possible in traditional intermediate view recycling systems. To make this more concrete let's look at an example.

Consider the following query over the TPC-H dataset that computes the average discount per country:

```

1  select      n_name, avg(l_discount)
2  from        lineitem, customer, nation, order
3  where       l_orderkey = o_orderkey
4  and         c_custkey = o_custkey
5  and         c_nationkey = n_nationkey
6  and         l_shipdate > 10-11-2015
7  group by    n_name

```

An optimizer considering this query in isolation would come up with some plan resembling the following following plan following:

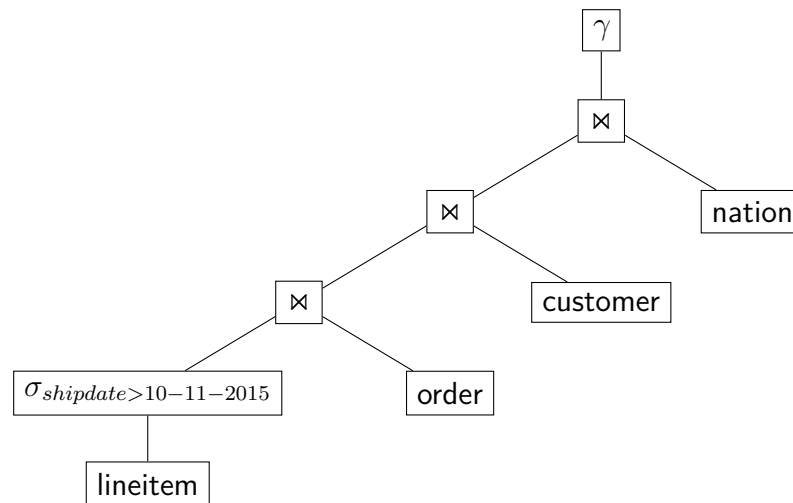


Figure 1.1.: An efficient logical plan for a single query.

Notice how the selection ($\sigma_{l_shipdate > 10-11-2015}$) is pushed all the way to the bottom of the tree because it is a cheap operation (worst case a full scan over the input) and can potentially shrink the input by a lot rendering the joins higher in the tree much cheaper.

Consider a workload that repeatedly runs a query generated from the following template:

```

1  select      n_name, avg(l_discount)
2  from        lineitem, customer, nation, order
3  where       l_orderkey = o_orderkey
4  and         c_custkey = o_custkey
5  and         c_nationkey = n_nationkey

```

```

6 and           l_shipdate > ${min_date}
7 group by      n_name

```

It is clear in this case that it would be beneficial if the workload is large enough for the cost of $lineitem \bowtie order \bowtie customer \bowtie nation$ to be amortized we would like to materialize the large join and only run $\gamma\sigma$ for each query attached as show in figure 1.2.

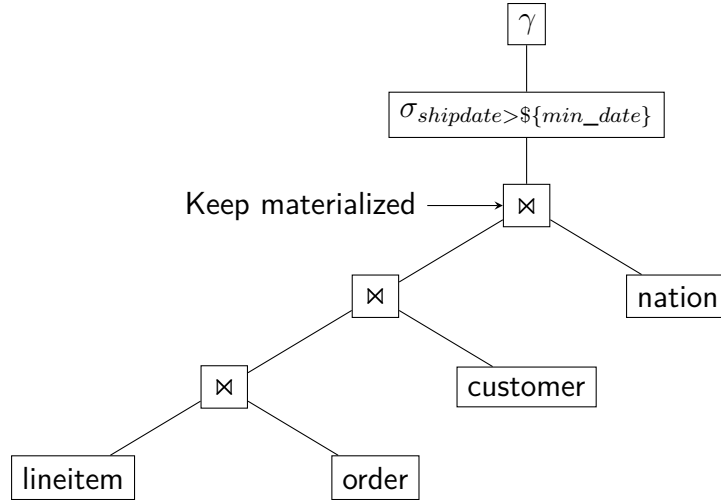


Figure 1.2.: An efficient logical plan for a template based workload.

This shift of focus from per-query optimization to considering the amortized cost of materialization of expensive relations is already a tall order for most RDBMSes that feature materialized views. A simple but effective approach to integrating incremental query materialization and the optimization processes, was presented in [76]. In their approach they maintain a *history pool* (a list of all the past queries) that is used to decide the benefit of materializing a sub-expression, and a *view pool* that keeps track of the materialized tables at every moment. Both these sets are taken into account during planning to produce a plan that will likely minimize the amortized cost of the workload. After the query is planned the sets are updated. A limitation of such an approach is that when dealing with relations like $lineitem \bowtie order$ in budget restricted settings, materialized view storage can quickly become a scarce resource.

There is an opportunity to reduce the effect of this problem by exploiting another common workload attribute: certain tables are frequently subsumed by the same intermediate result. In our example workload $lineitem$ is fully subsumed by $lineitem \bowtie order$, $lineitem \bowtie order \bowtie customer$, and $lineitem \bowtie order \bowtie customer \bowtie nation$. Similarly

order, *customer*, and *nation* are fully or partially subsumed by these relations. It seems wasteful to always keep the rows of all primary tables separately *and* concatenated in the rows *lineitem* ⋈ *order* ⋈ *customer* ⋈ *nation*.

In particular we can retrieve the rows of *lineitem* by projecting on the latter relation and deduplicating the resulting rows we can obtain any of the former. In this – contrived yet demonstrative – example a plan taking into account amortized costs could be:

Chapter 1. Introduction

```

1  Query  $\gamma_{\sigma_{p_0}}(lineitem \bowtie order \bowtie customer \bowtie nation)$  {
2     $Q_0 := \text{Materialize}[lineitem \bowtie order]$ 
3     $Q_1 := \text{Materialize}[customer \bowtie Q_0]$ 
4     $Q_2 := \text{Materialize}[nation \bowtie Q_1]$ 
5    # Not enough space to continue. Delete relations that we can rebuild.
6    GC {
7       $\text{Delete}[lineitem]$ 
8       $\text{Delete}[Q_0]$ 
9    }
10    $Q_3 := \text{Materialize}[\sigma_{p_0} Q_2]$ 
11    $Q_4 := \text{Materialize}[\gamma Q_3]$ 
12 }
13 Query  $\gamma_{\sigma_{p_1}}(lineitem \bowtie order \bowtie customer \bowtie nation)$  {
14   GC {
15      $\text{Delete}[Q_1]$ 
16   }
17    $Q_5 := \text{Materialize}[\sigma_{p_1} Q_2]$ 
18    $Q_6 := \text{Materialize}[\gamma Q_5]$ 
19 }
20 Query  $\gamma_{\sigma_{p_2}}(lineitem \bowtie order \bowtie customer \bowtie nation)$  {
21   GC {
22      $\text{Delete}[customer]$ 
23   }
24    $Q_7 := \text{Materialize}[\sigma_{p_2} Q_2]$ 
25    $Q_8 := \text{Materialize}[\gamma Q_7]$ 
26 }
27 ...
28 # Since the large join has all the columns of \li we should be able
29 # to create it by simply getting slicing and deduplicating
30 Query  $lineitem$  {
31    $lineitem := \text{Materialize}[\text{uniq}\{\pi_{cols(lineitem)} Q_2\}]$ 
32 }

```

Figure 1.3.: A sequence of plans optimizing the workload amortized cost and involving reverse operations. These queries are the same save for the predicate p_i selecting $lineitem$

By incorporating reverse relational operations where possible FluidB can indeed generate workload plans similar to the one described.

The solution we experiment with by implementing FluidB resembles the solution provided in [19] by Gou et.al for multi-quer. In their work they embed aggregations **group by** x_1, \dots, x_k into the \subseteq -lattice that arises from the powerset $P(\{x_1, \dots, x_k\})$. Thereby they encode the fact that **group by** A, B, C is subsumed, or can be computed by, either of **group by** A, B , **group by** A, C or **group by** B, C . Once the lattice is constructed a variant of the A^* path finding is used algorithm to search for the optimal aggregation plan. However they make no attempt to recycle tables, ie. garbage collect tables, whose data can be found in other relations, and narrows it's attention to aggregations.

From the aforementioned work we keep the basic notion of using a graph to represent relationships between queries and to derive the benefit of materializing a relation. We also use path finding techniques in that graph to create plans. However we introduce a more complex and ad-hoc hierarchy of relations to account for subsumptions the entire relational algebra, rather than just aggregations. That is very similar to AND-OR dags as found in [15]. The relations we express in that graph are bidirectional, so rather than only finding paths towards the goal and deleting relations when they are no longer needed via GC plan fragments. We simultaneously plan for moving towards the goal query and performing "backwards" operations to enrich the space of possible plans. We clarify this with an example which demonstrates a slightly simplified version of our system's functionality.

In figure 1.4 we provide representation of a subsumption graph that our RDBMS might create after witnessing selections on *lineitem*. For brevity let $p := shipdate > 10 - 11 - 2015$, $q := quantity < 24$ and $r := discount < .06$. The RDBMS has encountered $\sigma_p(lineitem)$, $\sigma_q(lineitem)$ and $\sigma_q\sigma_r(lineitem)$. For the purposes of this example we assume that according to the workload the planner has determined that the priority in terms of usefulness of each qury involved is in descending order of usefulness:

- $\sigma_q lineitem$
- $\sigma_{p \wedge r} lineitem$
- $\sigma_p lineitem$
- $\sigma_{\neg p} lineitem$

...

- *lineitem*

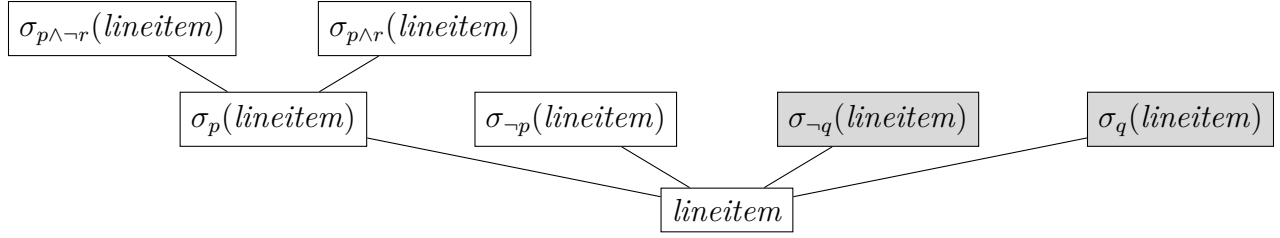


Figure 1.4.: The materialized relations are marked with grey. Assuming the absence of null values and assuming set semantics $lineitem = \sigma_p lineitem \cup \sigma_{\neg p} lineitem$, FluidDB can find a plan to generate any of the relations in the graph from $\sigma_p lineitem$ and $\sigma_{\neg p} lineitem$.

The query that we are planning is $\sigma_{quantity < 24} \sigma_{discount < .06}(lineitem)$, which is denoted in the figure as $\sigma_{q \wedge r}(lineitem)$. Our total size budget is 2.5.

Following the edges in the graph, to plan $\sigma_{q \wedge r}(lineitem)$ we need $\sigma_q(lineitem)$ and for that we need $lineitem$ and is least beneficial. So first we plan the union

$$\sigma_{\neg p}(lineitem) \cup \sigma_p(lineitem) \rightarrow lineitem$$

Then we need $\sigma_q(lineitem)$ but we are now using 2 units of space and adding .6 more would exceed our space budget of 2.5. $lineitem$ is the least beneficial of our materialized views but it is required for our next step, ie. creating $\sigma_q(lineitem)$. $\sigma_{\neg p}(lineitem)$ is deleted since it's derivable from $lineitem$. Then $\sigma_q(lineitem)$ is created and now we are using 2.1 units of space.

$$lineitem \rightarrow \sigma_q(lineitem)$$

Finally we need to create the final relation $\sigma_{q \wedge r}(lineitem)$. However it's space requirement is .5 and we would be exceeding our budget. We can't delete $lineitem$ even though it is our least beneficial table because we would have no way of recreating it.

Here we have two options. One is to delete $\sigma_{\neg p}(lineitem)$, our most beneficial table. The other, which is the one that the system should opt for, is to backtrack. When we created $\sigma_q(lineitem)$, instead we create both $\sigma_q(lineitem)$ and $\sigma_{\neg q}(lineitem)$.

$$lineitem \rightarrow \{\sigma_q(lineitem), \sigma_{\neg q}(lineitem)\}$$

Once both those tables are materialized we can safely delete *lineitem*. Now our space usage is 1.1 and we can safely

$$\sigma_q(lineitem) \rightarrow \sigma_{q \wedge r}(lineitem)$$

Which is the requested query and, in summary, the final can be represented as:

```

1  Inventory {
2     $Q_0 := \sigma_{\neg p}(lineitem)$ 
3     $Q_1 := \sigma_p(lineitem)$ 
4    ...
5  }
6  Query  $\sigma_{q \wedge r}lineitem$  {
7     $lineitem := \text{Materialize}[Q_0 \cup Q_1]$ 
8    GC { Delete $[Q_0]$  }
9     $Q_2, Q_3 := \text{Materialize}[\{\sigma_q(lineitem), \sigma_{\neg q}(lineitem)\}]$ 
10   GC { Delete $[lineitem]$  }
11    $Q_4 := \text{Materialize}[\sigma_r Q_2]$ 
12 }

```

Figure 1.5.: A sequence of plans optimizing the workload amortized cost and involving reverse operations.

This radical approach to adaptation of data to the workload is particularly useful in contexts where (fast) memory is scarce compared to the data, particularly in-memory databases. In that vain we lean heavily on executing plan by transpiling the plan to C++ code.

This thesis describes in detail how we implemented FluidB, a system targeted at performing this kind of reasoning to planning and executing queries. In chapter 2 we provide a brief overview of the state of the art and common practices that pertain to query processing; in chapter 3 we describe how queries are processed and stored at a logical level; in chapter 4 we describe how FluidB performs planning and garbage collection to come up with a concrete physical plan for each query; in chapter 5 we

Chapter 1. Introduction

go over an incremental evaluation system we built to facilitate the requirements of the physical planner; in chapter 6 we go over the algorithms and particular techniques that FluidB employs to transpile the physical plan into C++ code; finally at chapter 7 we describe our experimental evaluation of FluidB on the benchmark SSB-TPC-H. We present a conclusion and some future directions in the final chapter 8

Background

I am not the sower of chance,
the forger of new life

(K. B.)

Chapter summary

- Relational databases are question answering systems that deal with information organized in tables or relations.
- Query optimization and planning revolves around finding an efficient algorithm for answering a query.
- In-memory databases keep all the data in main memory so they are closer to the processor
- many in-memory databases employ some variant of code generation to execute the query plans.
- Intermediate result recycling is the practice of reusing computation between queries.

FluidDB is a system that generally focuses on relational query, optimization and planning. This chapter aims to give the reader some idea about where FluidDB fits in the design space and the historical context in which it was developed. First, we outline a very high level overview of the *query language and operators* involved in database management, as well as the overall architecture of such systems. Following that we will focus on the *query planning* subsystems of relational query databases and some traditional approaches to query evaluation. Afterwards, we will focus particularly on *in-memory relational databases* and the trade-offs that govern their design. Finally, we look into systems that utilize *intermediate result recycling* and how they solve the problem of automatically selecting and maintaining intermediate results for workloads.

2.1. Relational databases

Databases are more than just a method of accessing data. In their essence they are machines for question answering. The typical database works in a perpetual loop of reading queries and coming up with answers based on a set of data points. There are two important aspects that every database needs to define in order to delineate its operational semantics:

- The language in which the queries are expressed in
- The representation of data in terms of which the queries are expressed and the results are presented.

The oldest, most studied, and most common model is the relational model which defines queries in terms of *relational algebra* and organizes data in *tables* or relations.

A relation is typically an unordered set of tuples $\{d_1, d_2, \dots, d_k\}$ where each element d_i represents an attribute. The core relational algebra is an extension of the algebra of sets (that defines operators of set union \cup , intersection \cap , product \times , and difference $-$) that includes the operators of joins \bowtie , projection π and selection σ . Upon this foundation relational databases typically define extra operators to increase the expressive power of the language like aggregations γ , semijoins \ltimes , sorting, limiting, etc.

A typical relational database processes operates in a cascading fashion. It initially receives a query in a textual form. While relational algebra is the language that underpins the operation of a relational database, it is rarely the language in which users interact

with it. Instead, most relational databases expect queries written in a variant of SQL, a query language that is parsed into a tree of relational algebra operators.

This tree is processed by the query optimizer which rewrites it into a representation that is efficient to be executed by taking advantage of the mathematical properties of RA and often gathering statistics about the underlying data itself. This process leads to the formation of a relational algebra expression called the *logical plan*. It is a high level description of a sequence of operations required to produce a result. A logical plan only includes the denotational semantics of these operations, making no assumptions about their implementation.

Each RA operator is typically implemented by several different algorithms, each being efficient and even possible in some situations but not in others. For example, a join can be implemented by nested loops or with merge join among others. While the former algorithm is general and can implement any join, it is very inefficient. On the other hand a merge join is much more efficient but can only implement joins of the form $\bowtie_{a=b}$ (equijoins or natural joins). Furthermore, a *logical plan* typically contains only implied information about the scheduling of the operators, for example, the logical plan $(\pi A) \bowtie (\sigma B)$ implies that the join can not begin being evaluated before the projection and selection but the latter can be evaluated in any order, or even simultaneously via pipelining. All these details about the execution of the query are resolved by the *physical planner*. The physical planner accepts a logical plan and emits an unambiguous algorithm that will produce the result of the query, called the *physical plan*.

The final step is actually executing the physical plan which is handled by the *query execution engine*. The execution engine simply evaluates the physical plan on top of the data, making use of the correct algorithms, auxiliary structures like indexes managing memory, handling page-level caching, etc.

2.2. Query optimization and planning

Query optimization and planning generally refers to the processes that take place between reading a query and executing a physical plan. The concerns of query optimization are the correctness of the final result and the efficiency of the plan generated, usually in terms of time, but also in terms of space.

Finding an optimal plan is in the general case NP-complete [12], but query optimizers can do a good job at finding good plans using heuristics. The selection and organization

of these heuristics as well as query cost estimation are the main problems that make a query optimizer nontrivial. In particular one could separate the tasks of a query optimizer into three broad categories:

- Query rewriting
- Query plan enumeration
- Size and cost estimation (cost model)

At a very high level, the architecture of a query optimizer is demonstrated in figure 2.1 and is broadly comprised of the following components:

- The *parser* which receives the query in textual form and produces a logical plan in the form of an abstract syntax tree.
- The *deterministic optimizer* or *logical planner* that rewrites the logical plan applying optimizations that are based solely on the general properties of the relational algebra.
- The *physical planner* that transforms a logical plan into a physical plan that can be unambiguously executed by the execution engine to produce a result. The physical planner typically has at least some information about the underlying data that the plan will operate on like estimations about the statistics of the values or the cardinalities of the relations.
- The information about the data being manipulated by the plan is inferred by the *cost estimator*. It uses a cost model to predict the cost of plans and the cardinality of relations by taking into account the provenance of relations as well as physical properties of the data like the presence of indexes, the data layout, etc.

These subsystems are presented here as separate for simplicity and because in many database systems they are clearly delineated, but it is also common that they blur into each other. For example, in some RDBMSes the physical and logical planner are merged into one [7, 37, 44]. In fact, it is increasingly common for systems to intersperse query planning with query execution, adapting the optimization strategy [4] to concrete information about the intermediate results evaluated rather than purely relying on estimations

and predictions [21, 27, 50, 56]. The degree to which these subsystems are separate is a major concern in the design space of query optimizers.

Another important aspect to be considered, and which is indeed important in the design of FluidDB, is the number of queries considered at a time during optimization and the way in which they are considered. The optimizer usually considers one query at a time and maintains little or no state between executions. Although it has found little adoption in main stream databases, multiquery optimization has been researched extensively [14, 40, 53, 67]. What is more common in recent years is recycling intermediate results. RDBMSes that incorporate this technique materialize and cache intermediate relations, reusing them when they appear as sub-queries in later queries [28, 39, 42].

The final query processing concern that is relevant to the design of FluidDB regards the traversal and pruning of the search space. As mentioned, query optimization is generally NP-complete, so the viable options DBMS designers are left with are randomized algorithms, ML approaches, and heuristics. Virtually all systems implement heuristics entirely or to some degree, while recently more and more also incorporate randomized algorithms [33] and machine learning [59, 66].

2.2.1. Logical and physical query optimization

The *logical planner* transforms a syntax tree in the form of relational algebraic expressions, where the operators only contain information about the semantic meaning of operations and none relating to the algorithms that will eventually be executed on the input data. It is for all intents and purposes a rewrite engine for relational algebra. Typical transformations performed by the logical query optimizer are

- Predicate normalization to conjunctive normal form, e.g. $(a_0 = 1 \wedge a_1 = b_1) \vee b_2 = 3 \hookrightarrow (a_0 = 1 \vee b_2 = 3) \wedge (a_1 = b_1 \vee b_2 = 3)$
- Predicate push-down e.g. $A \bowtie_{a_0=1 \wedge p} B \hookrightarrow (\sigma_{a_0=1} A) \bowtie_p B$.
- Cartesian product to joins $\sigma_p(A \times B) \hookrightarrow A \bowtie_p B$
- Searching the join ordering space.

The *physical query planner*, on the other hand, will specialize the logical operators deciding on the particular algorithm that should be used. A physical query planner therefore requires low level information about the query that relates to the indexes available, possible ordering of the data, materialized views, etc.

Either of the planners needs to enumerate the plans under consideration while traversing the search space. There are two major approaches to this:

- The *top-down* approach, where the planner establishes the top level operator and branches searching the children, backtracking as necessary. This was the approach of the descendants of the *volcano optimizer* [6] where they implement a "search engine" that uses a branch and bound approach to optimization with caching.
- The *bottom-up* approach where the planner builds and connects fragments of a plan accumulating those fragments into a larger plan in a traditional dynamic programming approach. [34, 61]

An important optimizer, on which most modern optimizers are still being based, is the *cascades optimizer* [7]. In short, cascades keeps track of *groups* of equivalent query expressions and uses those as the fundamental atom that it manipulates. For example instead of keeping track of $A \bowtie (B \bowtie C)$ and $(A \bowtie B) \bowtie C$ separately they would be part of the same query group. It then uses a global hash table (the "memo" structure) to match the best plan that corresponds to each group.

Many database systems use optimizers similar to cascades (like the Microsoft SQL server, Postgres, MemSQL[49], and Greenplum – now Orca [45] to name a few). Notable among the descendants of cascades is Apache Calcite [55], a framework for implementing query planning used by a number of commercial and research databases [62].

2.2.2. Cost estimation

Probably the hardest aspect of the planner design is the cost estimation algorithm, which is required by the optimizers, especially the physical planner. These are required in order to select plans and to navigate the search space. A good cost model can help basic planners make decent plans and a bad one can cause sophisticated planners make horrible plans [48]. As cost estimation, we generally refer to a number of different related procedures that are broadly the *estimation of the cost* of an arbitrary plan and include the prediction of the *cardinality* of not yet materialized relations.

It seems that the most important challenges involved in the design of a cost model relate to the fact that we are fundamentally operating with scarce and highly uncertain information. Especially relating to cardinality estimation, uncertainty and bad predictions propagate and make make it extremely hard for the planners to make correct decisions.

Consider, for example a join. Some joins are similar to Cartesian products, producing large output tables, and some joins are more similar to lookups producing only a few rows. A cost model that confuses the kind of join will make very bad predictions w.r.t. the cost of any relational algebra expression that uses the said join, causing an optimizer to select an expensive plan or avoid a cheap one.

Cardinality estimation most commonly takes the form of selectivity estimation, i.e. what percentage of tuples from the input make it to the output of a selection or a join. This is sensitive to the statistical properties of the values and the selection/join predicates.

The most common approach is to maintain pre-computed statistics related to the primary tables. This yields better plans than keeping no statistics at all, but they are very hard to maintain and their effectiveness becomes very limited for complex queries. Another approach employed is to delay the decisions of the optimizer and essentially merge the planning and plan execution process. This does not completely eliminate the problem, but it allows the system to have more up to date and precise data that would allow it to correct course early in the event of exceptionally bad estimations. One family of techniques that is becoming popular and was initially used in IBM DB2 [16], is caching the statistics of previously computed relations and using that data to make better predictions in the future. More recently, this takes the form of employing machine learning to learn from past cardinalities [60]. This approach is still rarely used due to the engineering cost of keeping the collected information up to date.

To give the reader a more well-rounded intuition of the state of cost estimation besides cardinality estimation, it is worth mentioning that, to estimate the cost of an operation, Postgres uses magic configurable variables to weigh IO with CPU evaluation time, while DB2 runs micro-benchmarks on the production system to make these estimations.

Cost estimation is beyond the current scope of FluiDB and we only implement the most naive cardinality estimation, but it is worth mentioning how more mature systems approach the issue.

2.3. In-memory relational databases

In-memory databases are databases where all data lives in main memory. The design of in-memory databases is different from the design of a disk-backed database in a number of respects. To name a few:

Chapter 2. Background

- Page buffers have little use and caching of data in general has very different goals. While in disk-backed databases caches are mainly used to avoid disk IO, in in-memory databases they focus on reuse of computation.
- Concurrency control is much simpler as storage synchronization concerns are almost entirely eliminated.
- In disk-backed databases, only a small percentage of time is spend on actual computation [57]. Much of the non-compute latency is directly linked to the persistent storage.

On the other hand, there are concerns that are specific to the *lack of backing database storage*. To name a few:

- Should the system rely on record IDs like in a persistent database or can it use direct pointers to records?
- Error prone software can bring the data to an irrecoverable state.
- Query execution algorithms have fundamentally different characteristics. While in persistent databases IO dominates the runtime the bottlenecks for an in-memory database are much more complex and they can include things like locking, cache misses, predicate evaluation, and data movements.
- As main memory is a much less abundant resource than persistent storage, in-memory databases are often distributed, making network a major concern.

One increasingly common technique to address many of these issues, which is also used by FluiDB, is *code generation*. Since workloads for in-memory databases are typically CPU bound, there are major gains in performance to be had by specializing the code being executed. A typical interpreter based plan execution engine makes heavy use of virtual function calls and conditionals inside tight loops, which kills performance on virtually all modern architectures. The value proposition of code generation is to inline or hard-code the virtual functions and erase the conditionals at runtime to reduce the number of operations and make better use of hardware optimization.

We identify 4 different approaches in the literature to solving this problem:

2.3.1. Transpilation

Transpilation of a physical plan to a systems language like C or C++ and then fed to an off-the-shelf compiler [29] is expensive but it generates highly efficient code more easily debuggable execution plans. The most notable complete database systems that used this technique were Microsoft Hekaton that generates C code from SQL queries and older versions of MemSQL. For FluiDB, we reuse a lot of techniques introduced in HIQUE [29] to translate physical plans to template-heavy C++ code, making the assumption that query compilation will be much faster than the query runtime.

2.3.2. Third party JIT compilers

JIT compilation has received a lot of attention in the compiler community, especially in the context of the JVM. A number of database systems have taken advantage of this to speed up query execution. Most notable of these are SPARK [46], which generates Scala AST which is then converted to JVM byte code, and Neo4j that directly generates bytecode out of the queries.

Systems like Peloton [52] and recently Postgres [54] compile query plans to LLVM IR code that is then passed to the LLVM compiler to generate high performance machine code. Notable among the systems that use this approach is HyPer [75] which addresses the problem of query compilation overhead with an adaptive execution approach: they built an IR interpreter that starts running the query while the compiler does proper optimization and compilation of the LLVM IR program being interpreted. Cheap queries are thus completed reasonably fast, while in the case of complex queries the interpreted program is seamlessly replaced by the compiled one once the LLVM compiler finishes generating machine code.

2.3.3. Direct machine code generation

Some databases do not reuse any compiler or JITing VM, but rather directly generate highly specific machine code out of the physical plan. The first system to attempt that, like most techniques used today, was System-R, which originally compiled SQL statements directly to machine code by stitching together code fragments from a "fragment library" [2] but it was quickly deprecated due to the large engineering effort required. Oracle also includes similar to some of their databases and MemSQL express their plans in a custom

language called MPL (MemSQL Plan Language) for which they have a custom compiler that translates it directly to machine code.

2.3.4. Custom execution engines

The final category is per database code generation rather than per-query code generation. Volcano/EXODUS [6] and more recently SageDB [73] generate an optimizer and execution engine that is specific to the database schema but not to the queries.

2.4. Intermediate result recycling

A materialized view is a relation defined by a query that is persistently stored. A view that is not stored is said to be virtual. *View selection* is the process of selecting an appropriate set of materialized views to improve the performance of a workload [36]. Automated materialized view selection or intermediate view recycling has been on the radar of the database research community for a while now. A few approaches to this problem have to do with AND/OR directed acyclic graphs [9], modeling the problem as a state optimization [11], and lattices to represent data cube operations (i.e. multiple aggregations over the same relation) [70].

A related problem is *multiquery optimization (MQO)* [11] that attempts to plan multiple queries simultaneously. An efficient solution using AND/OR DAGs was proposed by Roy in [14] where they insert queries and their sub-queries in a graph and attempt to evaluate a plan by traversing that graph from multiple sources in a fashion similar to volcano optimization. Building on that Kathuria et. al. [51] present an approximation algorithm that runs in time quadratic to the number of common subexpressions and provides theoretical guarantees on the quality of the solution obtained.

Researchers have further looked at opportunistically reusing intermediate results that would be materialized. Both in traditional databases [28, 39] and on non-database contexts like MapReduce [35], and there have been attempts to unify the planner with the materialized view selection engine [43]. Notable in this field is Nectar [26] which is also not an RDBMS which automatically compresses rarely used data into programs that generate that data. Nectar focuses on sharing computation and data as much as possible.

Work such as MRShare [31] tries to bridge the gap between intermediate result re-

cycling and MQO by automatically grouping queries in a workload in such a way that computation can be maximally shared.

2.5. Conclusion

We presented an overview of the various techniques and challenges of problems that are directly related to the design of FluidDB and also some adjacent ones. FluidDB draws inspiration from all these systems to explore a new approach to managing in-memory data being queried. It subscribes to the relational model of data and queries and parallels many techniques introduced for intermediate result recycling and mutly query optimization. As an in-memory database it employs code generation and takes it to the furthest extreme generating highly generic C++ like the HIQUE system, fully compromizing on the performance of the C++ compiler.

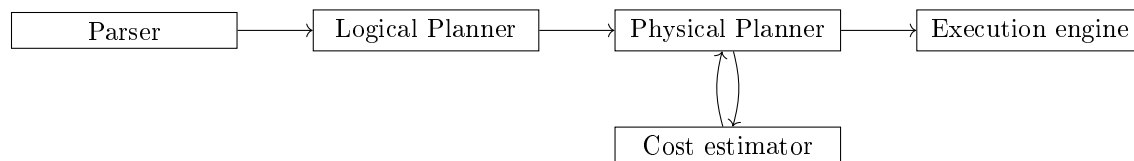


Figure 2.1.: Common architecture of a query optimizer.

Logical planning

Is there no end to this
accursed forest? Somebody
must climb a tree and have a
look round. The only way is to
choose the tallest tree that
overhangs the path.

(*Thorin*)

Chapter summary

- FluiDB expresses queries in relational algebra and defines reverse operations for most operations.
- Query plan trees are accumulated in a bipartite DAG (QDAG) where the types of nodes are tables and bidirectional operations.
- Queries are normalized to hashable objects (QNF) that abstract a wide range of rewrites. This allows efficient merging query plan trees into the QDAG.
- Sub-graphs of the QDAG (clusters) represent higher level operations and are used as *propagators* to infer information of the n-nodes like cardinality, column type, and uniqueness of subtuples.

FluidDB is focused on query optimization and planning when the queries in a workload being served can have common subqueries. To exploit these similarities between queries, FluidDB transforms them to bipartite graphs of subqueries and bidirectional operators, which she incrementally merges into a global query DAG and she uses to enumerate the possible plans. In this chapter, we will dive into how FluidDB processes and merges these graphs. and into a detailed description of the graphs themselves.

We will start by looking into what kinds of queries FluidDB can understand and how she processes them. Then we will look at how she relates the queries to each other to construct a tightly interrelated story of the workflow by looking at the basic query graph, the heart of FluidDB.

For manipulating a graph of queries, it is important that FluidDB can distinguish between semantically different queries and to be able to tell when they are the equivalent. We will look at how FluidDB normalizes queries to efficiently determine this equivalence relation.

Then we will see how FluidDB organizes the nodes in the query graphs into clusters representing single operations, how these clusters are formed within the query graph, and how those are used to maintain consistency among information about intermediate results.

Finally, we conclude by describing how the query processing infrastructure interfaces with the rest of the components: the planner, Antisthenis and the code generator, described in later chapters.

3.1. Relational Algebra

Much like most other RDBMSs, FluidDB deals with queries expressed in terms of a *relational algebra*. Here we define the FluidDB relational algebra semantics, that is, the traditional relational algebra extended by unary operators for ordering $s_e(X)$, limiting $l_i(X)$, anti-semijoin-ing \Join and aggregating $\gamma_e(X)$, as listed in table [3.1](#).

Description	RA	SQL
Aggregation	$\gamma_e(X)$	<code>select * from X group by e</code>
Ordering	$s_e(X)$	<code>select * from X order by e</code>
Limit	$l_i(X)$	<code>select * from X limit i</code>
Join	$X \bowtie_p Y$	<code>select * from X, Y where p</code>
Anti-semijoin	$X \ltimes_p Y$	
Union	$X \cup Y$	
Selection	$\sigma_p(X)$	<code>select * from X where p</code>
Projection	$\pi_r(X)$	<code>select r from X</code>

Table 3.1.: Correspondence between FluidB’s relational algebra expressions and SQL.

Before going any further, it is important to clarify some assumptions made to accommodate the design of FluidB. Particularly that

$$\sigma_p A \cup \sigma_{\neg p} A \equiv A$$

This means that there no **NULL** values are allowed since a predicate over a **NULL** value always evaluates to **False**, breaking the double negation theorem. This allows the definition of a bidirectional variant of the select (σ) operator.

Further, we assume set semantics for all relations, meaning that a valid relation needs to have at least one unique sub-tuple. The user is not burdened with the responsibility of enforcing this, it is rather enforced by a pre-processing stage on projections and aggregations, exposing extra columns than the ones specified by the user such that the resulting relation has at least one sub-tuple that is unique for each row (see section 3.3 on query pre-processing).

The semantics of the sorting operation are preserved by making a small concession to the set semantics: The *only* time when records of a relation are assumed by the planner to be ordered is for expressions of the form $s_e(X)$ and only for the duration between it running and the next operation. In practice, this means that the FluidB planner generally treats s like an identity operation.

Besides these caveats all operators known to FluidB behave as expected. The first of the extensions adopted is the antijoin operator \ltimes_p . It is equivalent to $A \ltimes_p B \equiv A \setminus \pi_A(A \bowtie_p B)$ and it is primarily useful for creating a bidirectional join by capturing the rows left over from a join (see subsection 3.2 on the QDAG). The important property we

will be interested in is:

$$A \Join_p B \cup \bar{\pi}_{cols(A)}(A \Join_p B) \equiv A$$

$\bar{\pi}_{cols(A)}$ means "group by the unique columns of A " As syntactic sugar we also define:

$$A \Join_p B := A \Join_{\neg p} B$$

3.1.1. Expressions

Relational algebra expressions contain expressions that operate at the field level. These appear in projection column definitions, selection predicates, aggregation column definitions, etc. These expressions are organized in up to 4 nested layers based on the kinds of operators:

- The top layer is purely Boolean operations **Prop** that can be rewritten in Cartesian normal form.
- The intermediate layer **Rel** describes relations between values (\equiv , **like**, $<$, \leq , etc). This layer is not nested, i.e. we do not allow relations between relations.
- The lowest layer **Expr** is expressions of any kind that may return any value. The above layers can be expressed in terms of this layer to support the ternary operator **if** <expr> **then** <expr> **else** <expr>.
- The symbol layer describes the literal and symbolic atoms involved in the query. Its type is parametric in most parts of FluiDB to allow us to tag the symbols in the expression without breaking compatibility.

Predicates of selections and joins (joins are always θ joins at the type level) have type **Prop** (**Rel** (**Expr** e)) for some type e and projections have type **QProj** [(e , **Expr** e)]. A special layer is provided for aggregations **Aggr** that makes it possible to disallow nested aggregation functions at the type level making the aggregation operator **QGroup** [(e , **Expr** (**Aggr** (**Expr** e))))] [**Expr** e], where the first argument of the constructor is the aggregation functions and their names, and the second argument is the expressions on which the grouping should happen.

The detailed definitions of the algebraic expressions are presented in listing 3.1.

```

1  -- | Algebraic type of a proposition
2  data Prop a = P2 BPOp (Prop a) (Prop a)
3    | P1 UPOp (Prop a)
4    | P0 a
5  data UPOp = PNot
6  data BPOp = PAnd | POr
7
8  -- | Algebraic type of an expression
9  data Expr s = E2 BEOp (Expr s) (Expr s)
10   | E1 UEOp (Expr s)
11   | E0 s
12  data UEOp =
13    EFun ElemFunction
14    | ENot
15    | EAbs
16    | ESig
17    | ENeg
18  data ElemFunction
19    = ExtractDay -- erased after parsing
20    | ExtractMonth -- erased after parsing
21    | ExtractYear -- erased after parsing
22    | Prefix Int
23    | Suffix Int
24    | SubSeq Int Int
25    | AssertLength Int
26  data BEOp =
27    -- Boolean (truthy)
28    EEq | ELike | EAnd | EOr | ENEq
29    -- Numeric
30    | EAdd | ESub | EMul | EDiv
31
32  -- | Algebraic type of a relation.
33  data Rel a = R2 BROp a a
34  data BROp

```

```
35  = REq
36  | RGt
37  | RLt
38  | RGe
39  | RLe
40  | RLike
41  | RSubstring
42
43  -- | Algebraic expression of an aggregation function
44  data AggrFunction = AggrSum
45  | AggrCount
46  | AggrAvg
47  | AggrMin
48  | AggrMax
49  | AggrFirst
```

Listing 3.1.: Types of algebraic expressions that compose the non-RA parts of the queries.

3.2. QDAG: The unified query graph

The operation of FluidDB revolves around an AND-OR graph of relations and algebraic operations, not unlike the ones commonly used in multiquery optimization literature. This graph encodes all plans that the query optimizer and planner will enumerate by means of traversal. Said graph is a bipartite graph where *t-nodes* (similar to AND nodes in other literature) represent relational operations and *n-nodes* (similar to OR nodes elsewhere) represent relations that will act as intermediate results. A set of extensions is introduced to the traditional AND-OR DAG to facilitate the operation of FluidDB. In the interest of motivating those extensions, we will take a very high level view of the operations on the QDAG before describing them in detail.

There are multiple ways of composing relational operations to materialize a certain relation. Each one is a relational algebraic representation, we commonly refer to them as logical plans and they are represented by trees similar to ASTs. These trees are made of two kinds of nodes: operation nodes (t-nodes), and relation nodes (n-nodes). For

example, the query $A \bowtie B \bowtie C$ may be represented as $(A \bowtie B) \bowtie C$ or $A \bowtie (B \bowtie C)$ due to associativity of \bowtie . Those two expressions represent two tree forms where the join operators \bowtie are represented as t-nodes and each of A , B , C , $A \bowtie B$, $B \bowtie C$ and $A \bowtie B \bowtie C$ are represented as n-nodes.

While t-nodes can be triggered bidirectionally, they are not symmetric: a subset of n-nodes are on the output side and another on the input side. Each n-node can be materialized by *triggering* one of the input t-nodes, by running a connected operation to which said n-node is on the output. It is clear then where the name OR-node found in other literature originates: any of the input nodes may participate to a plan involving the OR-node (n-node for us). Conversely, each t-node can be triggered \iff all of its input n-nodes are materialized. Because *all* input nodes must be materialized, t-nodes are elsewhere referred to as AND-nodes. We depart from the standard nomenclature to avoid confusion in the presence of the extensions introduced by FluidB. The planning of a query amounts to searching for a sequence of t-nodes to be triggered that would result in the corresponding node being materialized.

The QDAG of a query is created by unifying all logical plans that would derive it by means of merging n-nodes that are equivalent. For example the mentioned trees $(A \bowtie B) \bowtie C$ and $A \bowtie (B \bowtie C)$ shared the nodes $A \bowtie B \bowtie C$, A , B and C . A global QDAG is maintained by unifying the QDAG of each query encountered into the global QDAG, again, by merging equivalent n-nodes.

From a QDAG we can derive logical plans for queries. The planner operates in terms of the global QDAG. It traverses the QDAG searching for a plan (i.e., a sequence of t-node triggers) that will cheaply materialize the requested node(s). This way, historical nodes can be naturally considered. The problem of reusing materialized intermediate results is thereby conceptualized as a graph traversal problem.

This technique is heavily inspired by multiquery optimization (MQO) approaches like [15] where they use the AND-OR DAG to solve multiple queries simultaneously, optimizing computation reuse.

QDAGs are different from AND-OR DAGs in a subtle but important way: each t-node may have multiple outputs and can be partially bidirectional. In other words, a t-node, as well as multiple inputs, may have multiple outputs, any number of which may be materialized by it being triggered. A *reverse trigger* of said t-node would materialize inputs provided that all outputs are materialized. Where AND-nodes are similar to functions that when evaluated produce an output given some input, t-nodes correspond to

bidirectional relations between sets of input and output relations.

This not only creates opportunities in terms of what plans are possible, but also opens the gate to a whole new regime for how garbage collection is involved in the process of query planning. Any n-node may potentially be garbage collected, even primary tables, as long as the garbage collector can prove that it is reconstructable by triggering or reverse triggering t-nodes. The storage budget is thereby taken into account not only when selecting or rejecting materialized views, but also while planning the query solution itself.

To make this process clearer, as an example, consider a DAG that contains just a selection in figure 3.1 and a more complete example of a QDAG involving multiple selection is presented in figure 3.2

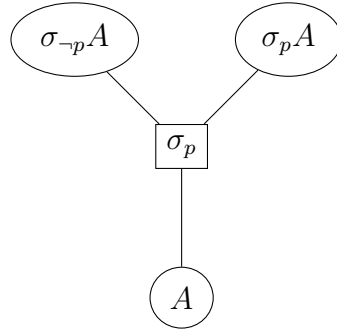


Figure 3.1.: A selection t-node may generate both output nodes, $\sigma_p A$ and $\sigma_{\neg p} A$, rendering it a partition operation.

In the figures throughout this work we will denote n-nodes as circular and t-nodes as squares or rectangles. When materializing $\sigma_p(A)$ the planner has the option of also materializing $\sigma_{\neg p}(A)$. If deemed beneficial A can be safely deleted as the t-node can be reverse triggered to produce it. A reverse trigger amounts to the union $\sigma_p(A) \cup \sigma_{\neg p}(A)$. It is generally assumed that any row not selected by p is necessarily selected by $\neg p$. To remind the reader of our RA assumptions, this means that we do not support missing or **NULL** values as they break Boolean double negation. Furthermore, relations are taken to be unordered sets of tuples.

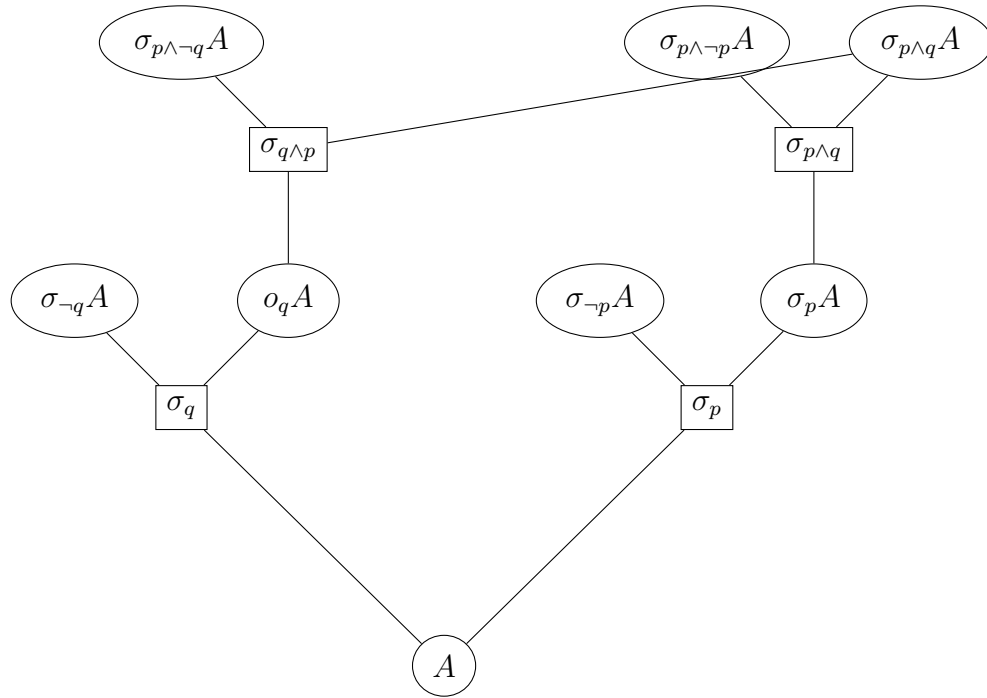


Figure 3.2.: Multiple selections can be composed to create larger networks.

In figure 3.2 if $\sigma_{\neg p}A$, $\sigma_{q \wedge p}A$, $\sigma_{p \wedge \neg q}A$ are all materialized, then A is materializable as their union. This way $\sigma_{p \wedge q}A$ is readily available with no space overhead. Instead, the price paid is that in the event that A is required, possibly in a different context, constructing it requires the overhead of unioning all its parts. Furthermore, there are two options for materializing $\sigma_p A$: either by selecting over A , or via the union $\sigma_{q \wedge p}A \cup \sigma_{\neg q \wedge p}A$. Table 3.2 demonstrates the correspondence between relational algebra operators and the equivalence that enables their reverse

Operator	Equivalence
Aggregation ($\gamma_e(X)$)	\perp
Sorting $s_e(X)$	\perp
Limit ($l_i(X)$)	$l_i A \cup d_i A \equiv A$
Join ($X \bowtie_p Y$) / Anti-semijoin ($X \ltimes_p Y$)	$\bar{\pi}_{cols(A)}(A \bowtie B) \cup (A \ltimes B) \equiv A$
Selection ($\sigma_p(X)$)	$\sigma_p A \cup \sigma_{\neg p} A \equiv A$
Projection ($\pi_r(X)$)	$\pi_r A \bowtie \pi_{\bar{r}} A \equiv A$

Table 3.2.: Each operation that may end up in the final result and the equivalence that ensures reversibility. \perp values indicate that the operator is not generally reversible.

A concept closer to the functionality of t-nodes and n-nodes than AND and OR nodes, and in fact our very inspiration for this approach, is *propagators* [23]. In short, propagators are hyperedges in a hypergraph, the nodes of which contain partial information about a value. Information about a value can be accumulated from its neighbors. Each node can obtain at least as much information about its value as the most precise neighbor. The example provided in the paper is one of measuring the height of a building by means of a barometer. One approach is timing the barometer's free fall from the top of the building, another is by measuring the length of the shadow of the barometer and comparing with the length of the shadow of the building at the same time of day. One may even use the barometer in its intended function and derive the height of by comparing the difference in air pressure at the roof and at the ground floor. Each of the mentioned methods will give very noisy or partial information about the magnitude in question. In the propagator paradigm, each of these measurements might be represented as a node in a hypernet. A propagator edge connecting these nodes would combine information from all measurements and improve each of the nodes' accuracy in their perception of the underlying value. A t-node is similar to a propagator in that it does not have a strict direction of information flow and can generate information on either side given information on the other.

3.2.1. Clusters and irreversible edges

Generally, there is no one-to-one correspondence between t-nodes and relational operators. Usually that is the case but there is an exception: the join operator. As we made clear, a t-node can generate any number of its inputs by reverse triggering it as long as

all of the outputs are materialized. Then one approach to implementing a join operator might be:

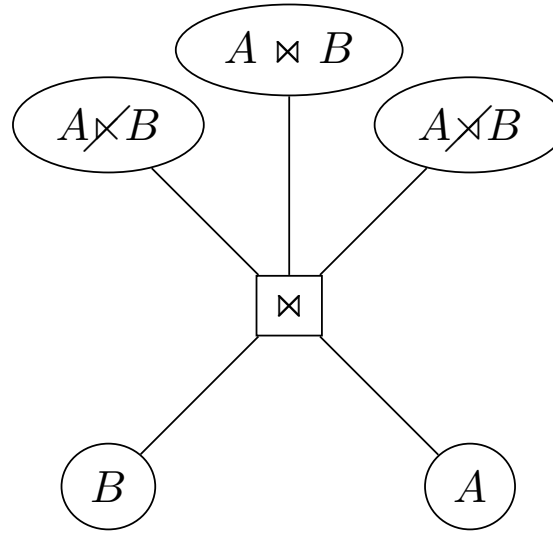


Figure 3.3.: Representing a reversible join as a single t-node.

The \Join operator is the antijoin operator, $A \Join B$ includes the tuples in A that do not participate in $A \bowtie B$. It is the complement of the semijoin \bowtie . Therefore

$$\bar{\pi}_{cols(A)}(A \bowtie B) \cup (A \Join B) \equiv A$$

$\bar{\pi}_{cols(A)}$ means project on the columns of A group by the unique columns of A . Therefore, by having materialized $\{A \Join B, A \bowtie B, A \Join B\}$ we can create A, B and vice versa. We could, however be slightly more flexible than that, as we might want to materialize $A \Join B$ but not $A \Join B$ to the end of ensuring the materializability of only A and not B .

In this case neither the entire input nor the entire output are materialized and yet all tables involved are materializable. To overcome the limitation that this poses, we make the graph a bit more complex and a bit more flexible by breaking the join into two stages (see figure 3.4). First, we separate the join operands (A and B) into the part that will be used in the join ($A \bowtie B$ and $B \bowtie A$ respectively), and the part that will not ($A \Join B$ and $B \Join A$ respectively) and then we join only the useful parts $(A \bowtie B) \bowtie (B \bowtie A) \equiv A \bowtie B$. Thus we fundamentally break the 1-1 correspondence between t-nodes and RA operators. Rather operators correspond primarily to *clusters* of nodes. As we will see in detail in 3.4 the join cluster is the only multi-t-node cluster.

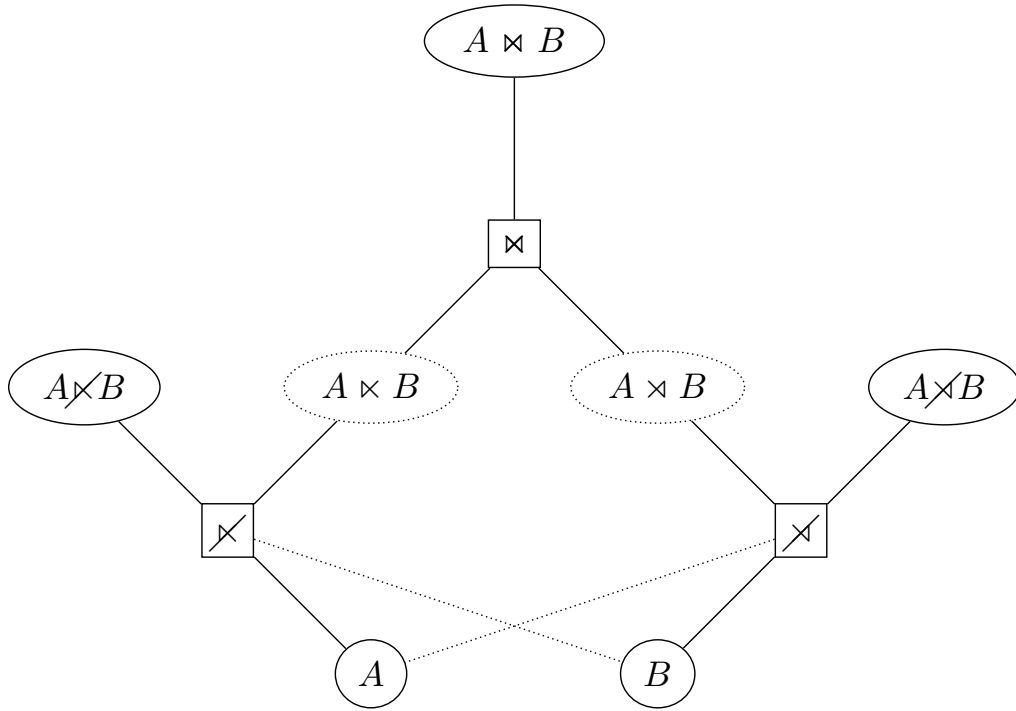


Figure 3.4.: The QDAG corresponding to a join is broken down in multiple t-nodes to facilitate greater flexibility.

The left antijoin node partitions A into $A \times B$ and $A \times B$, which is left semijoin in our null-less context. Similarly for the right antijoin node which partitions B to $A \times B$ and $B \times A$. Joining the semijoin relations $A \times B$ and $B \times A$ is equivalent to joining A and B .

In figure 3.4 we used an antijoin t-node. From the semantics of the antijoin $A \times B$ it is clear that no tuples from B are included in either $A \times B$ or $A \times B$ as B acts as a kind of predicate for selecting over A . This means that even though B is an input, there is no natural way of creating B from the t-node's output. For these cases we introduce *irreversible edges*. Remember that edges are always connections between n-nodes and t-nodes. An irreversible input edge (a connection from an n-node to a t-node) can only be used as an input to the t-node and its n-node can *not* be materialized by reverse triggering the t-node.

Irreversible edges appear in the inputs of anti-join t-nodes but also on the input of operations that cannot be information preserving, notably aggregations. The output of an aggregation can not in general be supplemented with extra information to reconstruct the input, unless the entire input table is replicated. We therefore give up the reversibility of

such operations completely and mark the input edge as irreversible. This is generally not a big problem as on the one hand the results of aggregations tend to be small compared to the input relations, and in the case of anti-join t-nodes, which are only found in join clusters, for every anti-join node $A \not\bowtie B$ that is unable to generate B by reverse triggering, there is a $A \bowtie B$ that will happily generate B albeit refusing to generate A .

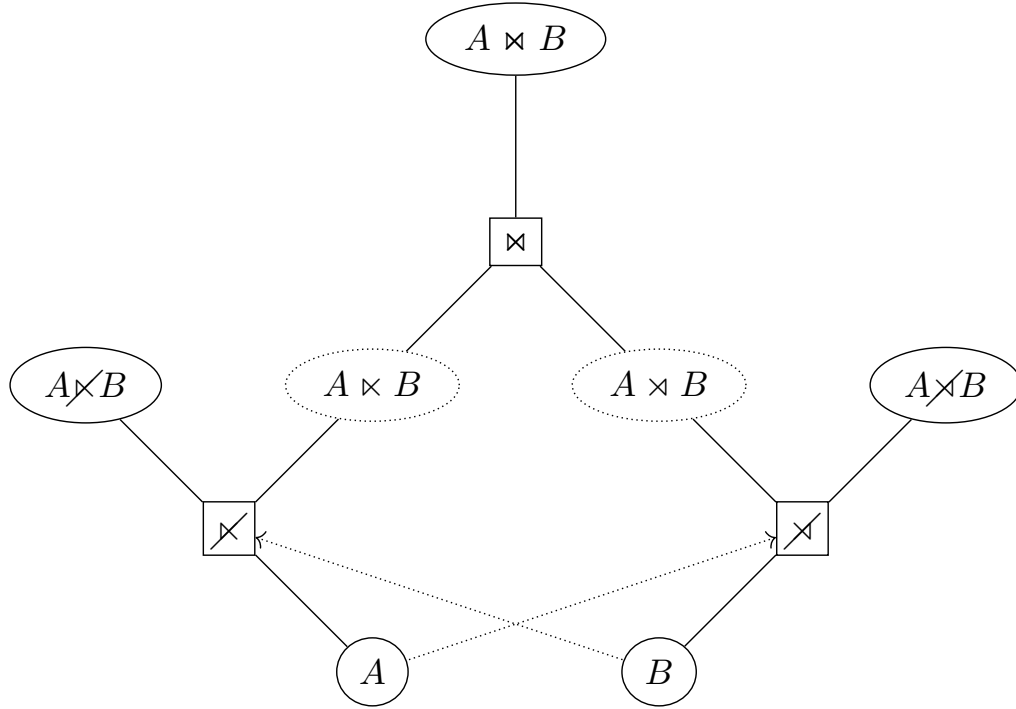


Figure 3.5.: The QDAG corresponding to a join is broken down in multiple t-nodes to facilitate greater flexibility.

3.3. Deterministic query processing

In this section, we will look at the initial steps of what happens when FluidDB receives a query. Initially, the query is accepted in textual form as SQL. Then it is sanitized to adhere to the properties we defined for relational algebra and some trivial optimizations are applied to remove artifacts of the parsing process. Finally, the query is translated to a forest of possible join orderings. These are rewrite-based optimizations that are deterministic and do not take into account the rest of the workload.

3.3.1. Query parsing

We use a custom library for parsing SQL queries into relational algebra based on parser combinator monads [74]. Our parser handles both the parsing and decorrelation of the queries. In brief, parser combinators is a functional interface to combining parsers to make larger, more complex parsers. A parser p has the form $p \ a$ which means that if the parser is used to parse the matching text, it will return a result of type a . For this framework to work, two combinators are required:

- sequence (monadic bind) $(\gg) : p \ a \rightarrow (a \rightarrow p \ b) \rightarrow p \ b$ which sequentially tries two combinators and the latter one may use the result of the former.
- the 'alternative' parser combinator $\langle \rangle : p \ a \rightarrow p \ a \rightarrow p \ a$ which tries one parser and if that fails tries another.
- **MonadGen** $e \ p$ can generate a value of type $p \ e$. Useful for creating names for unnamed expressions in projections, groups.

Simple parsers like the `parseInt` parser combinators are then combined to create more complex parsers (listing 3.2). This approach is fundamentally not dissimilar to writing yacc files but it is much more flexible as parser combinators function as a domain-specific language that inherits all the power of the host language, in our case Haskell.

```

1  parseInt :: MonadParse p => p Integer
2  parseInt = do
3    sign <- fmap (maybe 1 (const $ negate 1)) $ parseMaybe $ char '-'
4    void $ parseMaybe $ char '-'
5    (* sign) . read <$> readWhile isDigit
6
7  parseSelectProj
8    :: (MonadGen e p, MonadParse p)
9    => p e -> p (Query e s -> Query e s)
10 parseSelectProj eM = do
11   word "select"
12   fmap (Q1 . QProj) $ sep1 (word ",") $ do
13     ex <- parseExpr (E0 <$> eM)
14     parseMaybe (word "as" >> eM) >= \case

```

```

15     Just e → return (e,ex)
16     Nothing → do
17     e ← maybe pop (return . mkSymbol) $ asSymbol ex
18     return (e,ex)

```

Listing 3.2.: This parser returns a query modifier that is meant to be applied to a very simple product query generated by the **from** clause. The actual selection parser is complex and handles many different cases, but it is built up from simple fundamental blocks.

3.3.2. Query preprocessing

Immediately after parsing, the queries have several potential inconsistencies with the notion of relational algebra that we described in section 3.1. For that reason, we implement an intermediate stage where we rewrite the query to comply with various properties of the RA, and which remove redundancies. In brief, those are the following:

- Projections and selections are squashed. For example $\sigma_p \sigma_q$ becomes $\sigma_{p \wedge q}$ and $\pi_m \pi_{m'}$ becomes $\pi_{m \circ m'}$.
- Date intervals are translated to dates. This is because the physical planner does not handle dates specially so all dates and date intervals are converted to timestamps and second intervals. For example the expression `date '2020-12-01' - interval '90' day` is immediately turned into `date 2021-03-01`. Note that an expression like `date '2020-12-01' + shipping_time` is not supported as it is not date arithmetic between literals.
- Optimization of the **like** operator: `<VAR> like <STRING_LITERAL>` expressions where `<STRING_LITERAL>` begins or ends with `%` are turned into `suffix_of` and `prefix_of` expressions. For example, the expression `a like "%.txt"` would be turned into `".txt".suffix_of(a)`.
- Substring expressions are turned into 0-indexed ones from 1-indexed as is the case in SQL.
- It is asserted that each intermediate relation has uniquely named columns. For example, there are necessarily name conflicts in $A \bowtie A$.

- It is asserted that sorting appears higher in the query AST than a selection, a join or an aggregation. All three of these operators affect the ordering of their inputs (see chapter 4 on code generation), and sort is the only operation that asserts tuple ordering on its output. Therefore $s(A \bowtie B)$ is valid but $s(A) \bowtie B$ is not.

Once the preliminary sanitization and simple optimization steps are done we mark each instance of the primary table symbols with the columns and their types. This is simply information derived by the database schema. Then we do some processing and annotation on symbols referring to symbols:

- First *remap unique sub-tuples*. Each intermediate result, and each relation that we want to be able to reason about needs to have uniquely identifiable rows. For each relation in the QDAG we are keeping track of a set of columns, the unique sub-tuple, the combination of which is unique among each tuples of the same table. It is fairly straightforward to keep track of the unique subtuple for most operations, e.g., the unique sub-tuple of a Cartesian product or join is generally the concatenation of the unique sub-tuples of the input tables, the unique sub-tuple in the case a selection is the same as the unique sub-tuple of the input, in aggregations the unique sub-tuple is the same as the columns on which we aggregate, etc. The main challenge is projections and, in a similar manner aggregations, where it is possible that the operation does not project on the entire unique sub-tuple of the input. To mitigate this, we perform a pre-processing step to rewrite the projections to expose the entire unique sub-tuple. It is worth mentioning that there are cases where there are more than one combinations of fields that form a unique sub-tuple. We opt to keep track of all of them in order to have more flexibility in this pre-processing step: ideally we want to modify the projection as little as possible. In other words we want to project as few extra columns as possible to keep the output table size as small as possible. As a simple example the query $\pi_{p_color} part$ would become $\pi_{p_color, p_partkey} part$.
- Annotate each symbol referring to a column with the corresponding primary table if there is one (named columns of projections do not correspond to a primary table). This is useful both in determining possible joins (as we will see in section 3.3.3) and in determining the types of columns. This inference is sometimes trivial, for example, in the case of references to columns of primary tables, or when annotating literal values, but when references to columns of projections it is more complex.

For example, in the query $\sigma_{p(a)}\pi_{a \rightarrow f(X)}$, we need to infer the type of a . FluidDB will apply some very simple heuristics to fake proper type inference. This is usually enough to come up with accurate types as SQL's type system is basic.

- Projection sharing is determined. As we will see in detail in section 6.5.3 where we will be discussing the implementation details of the projection algorithm and its reverse, it is important that at least one unique tuple of a projection is shared between a projection and its complement. At this stage, it is decided which unique subtuple is to be shared and it is recorded in the RA expression tree.

3.3.3. Possible joins

Join ordering is one of the most detrimental optimizations of a query. It is the selection of the associative order of joins. The relation $A \bowtie B \bowtie C$ can be expressed as $(A \bowtie B) \bowtie C$ or as $A \bowtie (B \bowtie C)$. Individual joins can have a wide range of behaviors. From being similar to Cartesian products, being prohibitively costly, and exploding the cardinality of the result; to being similar to lookups, very cheap, and with a very small result. Optimizers that process queries in isolation can aggressively push down selections and prune the search space of join combinations when they encounter expensive intermediate results. FluidDB, however, needs to keep track even of more expensive intermediate results as the high cost may be amortized in the long run.

For this reason, at this stage we enumerate all possible join permutations building a graph of all possible joins.

- All joins are turned into selection-products, so $A \bowtie B$ becomes $\sigma(A \times B)$
- selections are pulled up and merged, $(\sigma_p \sigma_q A)$ becomes $\sigma_{p \wedge q}$
- products are represented as sets of relational algebraic expressions.

Then all possible permutations of joins are created. For example $\sigma_{a=b \wedge b=c}(A \times B \times C)$ becomes $\{(A \bowtie_{a=b} B) \bowtie_{b=c} C, A \bowtie_{a=b} (B \bowtie_{b=c} C)\}$.

Note that we do not include queries that are only different w.r.t. commutativity, i.e., we don't include both $A \bowtie B$ and $B \bowtie A$. This part of FluidDB depends on the consumer of the data structure to disambiguate between them. Another caveat is that since Cartesian products are very large and are almost never used in the actual plans, an option that disallows product intermediate results is enabled by default. Therefore,

in our example we prune relations like $\sigma_{a=b}((A \times C) \bowtie_{b=c} B)$. Furthermore, allowing Cartesian products exponentially explodes the space of possible plans to a degree that makes the resulting data structure virtually unusable.

Once selections are merged, we find product trees, i.e., contiguous parts of the tree that only contain Cartesian products, and combine them into multi-sets of subtrees such that $\sigma(C \times (A \times B) \times D)$ becomes $\sigma(\{A, B, C, D\})$.

Then selection propositions are normalized to conjunctive normal form and into sets. For example $\sigma_{p_1 \vee (p_2 \wedge p_3)} A$ becomes $\sigma_{(p_1 \vee p_2) \wedge (p_1 \vee p_3)} A$, and in fact the conjunctive terms are also represented as sets of subterms, so the final form would be better expressed as $\sigma_{\{p_1 \vee p_2, p_1 \vee p_3\}} A$. This leaves us with a tree of terms of the form $Q[\sigma_{\{p_0, p_1, \dots, p_l\}} \{Q_0, Q_1, \dots, Q_k\}]$ where each term p_i refers to columns in one or more Q_j relations and $Q[\cdot]$ is a query tree with \cdot as a subtree. We will refer to such a selection-product pair as an SP-term.

To turn this representation into a forest, we will turn each SP-term into a set of regular join trees, the Cartesian product of which will yielded the possible joins. We summarize the process as follows:

- We iterate over all pairs of sub-multisets into which the multiset $\{Q_0, Q_1, \dots, Q_k\}$ can be split. We follow each of the following steps for each of the resulting sub-multisets.
- We split the p terms in three parts depending on which sub-multiset of relations their symbols refer to. One set are terms where the columns mentioned refer *only* to the left sub-multiset (left p -terms), another set are terms where columns mentioned refer *only* to the right sub-multiset, and the rest (right p -terms) must necessarily have references to columns of *both* sub-multisets. We call the latter category *connector p -terms*.
- We turn each side of the split into an SP-term by associating it with the left and right p -term set. We repeat the process of generating join trees out of each SP-term created and connect them in an all-to-all fashion by creating theta joins using the conjunction of the connector p -terms. In case the the set of connector p -terms is empty, that would mean that the connection of the join trees must be done with a Cartesian product. As mentioned above, we want to avoid inserting Cartesian products in our graph so we reject those cases outright. Optionally, we may even want to reject all non-natural joins altogether, unless that would mean not coming up with any plans at all.

The fundamental forest structure of literally collecting each logical plan individually is highly redundant. Therefore, we abstract away the redundancy by representing the forest as a *bipartite tree*.

- One kind of node represents a plan fragment, or a *sub-tree*
- The other kind represents a forest of semantically equivalent subtrees earmarked with a unique hash for them. The fact that the subtrees are semantically equivalent means that they essentially evaluate to the same data. They therefore have the same normal form as discussed in section 3.4 on Query Normal Form (QNF). Finding the normal form hash of one would mean the normal form hash of the rest. This will prove useful when inserting this structure in the graph. We refer to these as *sub-forests*.

In other words, the relational algebra tree is interleaved with OR-nodes of different representations of the subquery. For example $\gamma(\sigma_{a=b \wedge b=c}(A \times B \times C))$ becomes $\gamma\{A \bowtie_{a=b} B) \bowtie_{b=c} C, A \bowtie_{a=b} (B \bowtie_{b=c} C)\}$. This way we mitigate the combinatorial explosion of interleaved nesting joins with other operations.

The algorithm for building the forest out of the pre-processed query is described in listing 3.3.

```

1  def possible_joins(q):
2      # Split a query that is equivalent to
3      # sel(p1 and p2 and p3 and ... , prod(q1, q2, q3, ... ))
4      # into
5      # [p1,p2,p3, ..], [q1,q2,q3, ... ]
6      and_props, subqueries = as_prod(q)
7      # If there are no top level joins traverse the query
8      if len(and_props) == 0:
9          return q.map_children(possible_joins)
10     # All possible partitions where the partitions are non-empty
11     for lqs, rqs in possible_partitions(subqueries):
12         # has_free_vars(p, [q1,q2,q3 ... ]) checks if all the atoms in p
13         # are either literals or columns in one of the q1,q2, ...
14         lps = filter(lambda p: not has_free_vars(p, lqs), and_props)
15         rps = filter(lambda p: not has_free_vars(p, rqs), and_props)

```

```
16     # The props that connect the partitions into a join.
17     jps = filter(lambda p: has_free_vars(p,lqs)
18                  and has_free_vars(p,rqs),
19                  and_props)
20     # If there are no predicates connecting it is product. We
21     # disallow products.
22     if len(jps) == 0: continue
23     # A valid join to be considered while planning!
24     lq = sel(lps,prod(rqs))
25     rq = sel(rps,prod(rqs))
26     yield join(jps,possible_joins(lq), possible_joins(rq))
```

Listing 3.3.: Pseudo-python description of finding all possible for clarity, it is abbreviated to omit sanity checking, memoization, some type conversions, etc. (This algorithm does not include memoization for brevity)

It is worth stressing that, many of the subtrees generated this way will be equivalent. As analyzed in [63] and discussed in section 3.4, in pure functional settings like our, DAGs can have explosive complexity during traversal. Even worse in this case even the space complexity can indeed be needlessly exponential. For that reason we maintain a hash map matching hashes with the corresponding matching subtrees. Before inserting a newly created forest into our structure we perform a lookup in the hash map. If we find something we throw away the newly created forest for the optimizer runtime's garbage collector to recycle, and we use the one in the hash map instead. Otherwise, we insert our forest both in the hashmap and in the tree. This way we deduplicate our data structure to some degree. Furthermore, in the same vain as [63] we earmark all our forests with their hash so we don't have to traverse them again to compute it. The result of this is that in total we completely traverse the forest once, in order to compute the hashes of each sub-forest and then we use those hashes to test sub-forests for equality.

3.4. Query Normal Form (QNF)

It is imperative that the QDAG has as little redundancy as possible, two subqueries that describe the the same relation should be represented as the same node. Equivalence of RA expressions is infamously NP-complete [1] but we make a best effort to normalize

queries such that they are hashable and comparable via simple structural equality. The core of normalization is transforming the queries to Cartesian normal form

$$\zeta_{x_1}(A) \bowtie_{p_1} \zeta_{x_2}(B) \bowtie_{p_2} \zeta_{x_3}(\sigma_{p_3}(C)) \rightarrow \zeta_{x_1;x_2;x_3}(\sigma_{p_1 \wedge p_2 \wedge p_3}(A \times B \times C))$$

where $\zeta \in \{\pi, \gamma\}$.

3.4.1. QNF Structure

FluidDB's Query Normal Forms are highly polymorphic data structures separate from RA data structures. The general QNF form can encapsulate one of several special forms, all of which strive to abstract away as many valid RA rewrites as possible, and ideally achieve a one-to-one correspondence between a QNF and a table. While the QNF in its general form (listing 3.4) is specialized to represent several different forms internally to the QNF subsystem, outside of it there are two main objects:

- The normal form of a relational algebra query
- A column of a normal-form query.

```

1  data QNFQuerySPDCF sel_f prod_f dbg_f col_f f e s =
2    QNFQuery
3    { qnfColumns :: col_f (QNFProj f e s) (QNFAggr f e s)
4      -- ^ What columns kept. Each of these is a column on one of the
5      -- QNFProd.
6      , qnfSel    :: sel_f (QNFSel e s)
7      -- ^ And conjunction of selection. The QNFNames here lack a qnfProd
8      -- and qnfSel because they are names of this QNFQuery itself so to
9      -- avoid a recursive structure that we need to keep up to date we
10     -- just promise that they will be empty.
11     , qnfProd    :: prod_f (QNFProd e s)
12     -- ^ QNFProd e s refers to multiple rewrites of a single
13     -- relation. Essentially {Query (..) (QNFQuery e s)}. Each element
14     -- of this set needs to have the same set of QNFQuery's. The binary
15     -- operators allowed here MUST expose plans from ONE side. Otherwise
16     -- the qnfColumns will be ambiguous. In practice only Product/Join

```

```

17      -- expose both sides but if there are more in the future know that
18      -- qnf machinery will break.
19      ,qnOrigDEBUG' :: dbg_f (Query e s)
20      ,qnfHash :: QNFKey e s
21      -- ^ Cache the hash. Hashing is done separately for each field so we
22      -- don't recompute hashes of unchanged fields after each operation.
23  }

```

Listing 3.4.: The QNF datastructure.

In section 3.4.2 on building QNFs, we will look at several specializations that are useful as intermediate. In this section we will focus on each of the fields of the **QNFQuery** constructor.

QNF Product collection

The `qnfProd :: prod_f (QNFProd e s)` field contains the unordered collection of sub-queries $\{A_1, \dots, A_n\}$ in the equivalent $\pi\sigma(A_1 \times \dots \times A_n)$, typically as a hash-set (via the instantiation of `prod_f` as **HashSet**). Each of the A_i terms may be one of the following:

- A primary table
- A QNF with an opposite kind `qnfCol` (if Q is a projection an A_i may be an aggregation)
- A set of equivalent RA expressions that are not reducible to the $\pi\sigma(A_1 \times \dots \times A_n)$, for example $\{A \bowtie B, B \bowtie A\}$ or $\{l_N(A)\}$

Focusing on the latter case, there are certain equivalent rewrites that are not easily expressible in the QNF. For that reason instead of perpetually extending the QNF system we provide this expensive method of simply enumerating equivalent queries in RA form. These RA queries have at the leaves of the expressions either QNFs or primary tables. The only restriction that we place on these queries is that the schema of the inputs is simply forwarded to the output without modification of the individual fields so that the `cnfCol` vector can refer directly to them.

QNF Columns, projections and aggregations

The `qnfColumns :: col_f (QNFProj f e s) (QNFAggr f e s)` field indicates which columns are exposed by the query. There are two kinds of column sets: aggregated columns and projected columns. Depending on the context `col_f` is instantiated to a type that allows the column to be *always* aggregation, *always* projection, or either of the two.

Whenever a `QNFQuery` is found outside of the QNF sub-module `col_f` will be instantiated as `Either`, which is to say that a query may be an aggregation (`Right`) or a projection (`Left`). A query with no explicit projection operand at its root is translated to QNF by enumerating the columns into this field. In all of these storing them in a container of parametric type `f`, which is the main distinguishing factor between QNF columns and QNF queries: a query projection exposes an unordered bag of columns through the projection/aggregation field (`f` is instantiated to `HashBag`) while a QNF column exposes *exactly one* column (`f` is instantiated to `Identity`).

We extend the notion of a column to to define names (see listing 3.5), the atoms of predicates and other expressions (see section 3.1.1). A name in the context of QNFs can be one of three different kinds:

- **A column of a nonprimary relation.** This relation *must* be a member of the `cnfProduct` collection of the QNF structure where the name appears. Because there may be more than one equivalent query in the product collection, we index the columns with an integer.
- **A column of a primary table.** These columns have a name and are associated with a table (`s`) which needs to be in the product collection like in the case of nonprimary columns. It is indexed much like the nonprimary relation.
- **A literal value**

```

1 data QNFNameSPC sel_f prod_f col_f e s =
2   PrimaryCol e s Int
3   | Column (QNFColSPC sel_f prod_f col_f e s) Int
4   -- ^ a a column c that refers to a:
5   --
6   -- QNF[P[c],S[c],Prod[Query[a ~ QNF]]]
7   | NonSymbolName e

```

```

8  -- ^ A literal.
9  -- A column is a QNF that has one element in the projection set
10 -- (Identity) and no debug query.
11
12 type QNFColSPC s p c = QNFQuerySPDCF s p EmptyF c Identity

```

Listing 3.5.: A QNF name may be an unnamed column of a relation, a named column of a primary table or a literal.

In listing 3.6 we present the definition of the projection field of the QNF structure. It departs from the definition we examined in section 3.1 where we described the RA in two important ways, a) no names are provided for the columns and b) the columns are unordered, depending on the definition of `f`. We maintain, however, the definition of each column as an `Expr`. FluiDB makes no attempt to normalize the expressions, so equivalent expressions that differ syntactically will produce different QNFs.

```

1 type QNFProj f e s = f (Expr (QNFName e s))

```

Listing 3.6.: A QNF projection field is a collection of expressions that refer to QNF names. The particular structure of this collection is parametric. When the collection `Identity` the QNF query is essentially just a column. A normal QNF query would instantiate `f` to `HashBag`, an unordered multiset.

As mentioned, the columns in the QNF name, both primary and of general relations, are indexed in the case that there are more than one equivalent. This is an edge case that can cause problems because different indexing can cause equivalent nodes to have different QNF representations. For example, the query `select * from A as A1, A as A2 where A1.a = A2.a` is symmetric in all respects and therefore the way QNF names representing columns `A1.a` and `A2.a` are indexed makes no difference provided that the equality expression `A1.a = A2.a` is ordered consistently.

Consider, however, the query `select * from A as A1, A as A2 where A1.a = A2.a + 1` which is not symmetrical. Indexing `A1` and `A2` differently does not change the semantics of the query, but it changes the QNF representation. This is not commonly a problem, a rare edge case rather, but it is important to keep in mind. In an attempt to mitigate this *all* different combinations of indexing are emitted by the QNF builder (see section 3.4.2).

From the perspective of QNFs, aggregations (defined in listing 3.7) are very similar to

projections. However, they only operate on the non-recursive **Aggr** algebra, in contrast to the **QGroup** unary RA operator that incorporates **Expr** (**Aggr** (**Expr** e)) typed expressions.

QNF moves the inner and outer **Expr** to one level up and one level down respectively, so an aggregation $\gamma_{l \mapsto f(g(h(x)))} S$ would be rewritten to the normalized $\pi_{l \mapsto f(l_1)} \gamma_{l_1 \mapsto g(l_0)} \pi_{l_0 \mapsto h(x)} S$.

```
1 type QNFAggr f e s =
2   (f (Aggr (QNFName e s)), HS.HashSet (Expr (QNFName e s)))
```

Listing 3.7.: The QNF aggregation form of the projection field is similar to projection only, much like the **QGroup** constructor, it also includes a **HashSet** of expressions on which to group.

QNF Selection

The `qnfSel :: sel_f (QNFSel e s)` field contains an expression of the selection predicate (see listing 3.8). With very few exceptions `sel_f` is instantiated as an unordered hash containing the terms of the Cartesian normal form of the selection. For example a `qnfSel` with value $\{p_0, p_1, p_2\}$ corresponds to the selection predicate $p_0 \wedge p_1 \wedge p_2$.

```
1 type QNFSel e s = Prop (Rel (Expr (QNFSelName e s)))
```

Listing 3.8.: Selection name refers to a version of the current QNF that has all fields erased except the projection.

We went into detail in the previous section how the atoms of the predicates in QNFs, when they refer to columns, take the form of QNF structures themselves. Further, we asserted that the columns in the expressions of projection/aggregations refer to the product. Similarly, selection names refer to columns of the product. This means that in the QNF of the expression $\sigma_{f(a)}(Q_0 \times Q_2 \times \dots)$ the atom a is encoded as a column of one of the Q_i subqueries.

This way we can update the fields of the QNF without needing to keep the columns in sync with the columns referenced in `qnfSel`. This eliminates the need for both indexing the columns and for dealing with primary columns in **QNFSelName**. This has to do with the way selections are being created: the RA expression $\sigma_p A$ is translated by first translating A to a named QNF. Then p refers directly to the columns of the query.

QNF Misc fields and hashing

There are two more fields in the generic QNF datastructure we described. `qnfOrigDEBUG' :: dbg_f (Query e s)` is the more straightforward one. It simply keeps track of the query used to create the QNF for debugging purposes and is disregarded by both the equality and hashing operations. Typically `dbg_f` is instantiated to `EmptyF`, an empty container to save on memory.

The other field is `qnfHash :: QNFKey` (see listing 3.9). Because the QNF is a highly self referential data structure equality in a purely functional environment like Haskell traversing it for the purposes of equality or is a costly operation. For that reason, we depend on hashes to detect equivalence, and we cache those hashes to avoid recomputing them, leaning on the purity of Haskell. As the QNF has three distinct parts that are being often updated during building of a QNF `qnfHash` is represented as a tuple of their respective hashes. This way, when one of the fields is updated, we do not have to recompute the hash of each of the others.

```
1 type QNFKey e s = (Int,Int,Int)
```

Listing 3.9.: A key that uniquely identifies a QNF is made of three separate hashes, one for each part of the qnf structure so that they can be updated independently.

Since QNFs are rarely used as anything other than a token for relating RA form queries FluidDB saves on memory and use `QNFKey` in place of the `QNFQuery` structure (see listing 3.10). On the other hand, in a more pedantic but correct mode, we can avoid depending on the quality of the hash function and define equality as recursively checking between all types (see listing 3.11). In practice, we have never found the former method to yield a false positive for equality.

```
1 instance Eq (QNFQuerySPDCF sel_f prod_f d col_f f e s) where
2   x == y = qnfHash x == qnfHash y
```

Listing 3.10.: A fast and loose definition of equality between QNFs that depends on the quality of equality.

```
1 instance Eq (QNFQuerySPDCF sel_f prod_f d col_f f e s) where
2   x == y = qnfProd x == qnfProd y
3   && qnfSel x == qnfSel y
4   && qnfColumns x == qnfColumns y
```


Listing 3.11.: A very inefficient but correct equality between QNFs.

3.4.2. QNF Computation

As it became clear while describe the parts of a general QNF in many cases transformation algorithm from RA to QNF can be easily inferred. In this subsection we will look at a few parts of the QNF building process that are maybe less than obvious.

QNFs are built bottom up from the leaves to the top. Building the QNF we apply the operations incrementally to a named QNFs (NQNF). Applying the operator σ_p to a QNF A we need to be able to relate the symbols appearing in p with the columns of the QNF A . To mitigate this define a *named QNF* which is simply a tuple of a QNF along with a mapping between symbols and columns of that QNF. For example, the name map in the NQNF derived from the query $\pi_{n \mapsto f(a)} A$ would map the symbol n to a column with the elements $\{f(a) | a \in A\}$.

The entire process of computing QNFs works within a monad stack with the following features:

- A **ListT** monad allows for non-deterministic computation in order to account for the different possible column index assignments described in subsection 3.4.1 on QNF projections/aggregation.
- Since the resulting QNF is a highly self referential structure, it is to be expected that a lot of the computation might be duplicated. The cache object is also shared between different QNF building processes that we know are likely to have a lot of overlap, for example, when inserting nodes into the graph without shared caching FluidB would be repeatedly computing the QNF of the same RA subqueries. This takes the form of simply memoizing the interface functions for building the QNF.
- The process of building a QNF can fail when processing malformed RA forms.

These characteristics are encoded in the monad within which we build the QNF (listing 3.12)

```
1 -- | The monad in which the computation happens. Not all transformations
2 -- require non-determinism.
```

```

3 type QNFBuild0 e s = StateT (QNFCache e s) (Either (QNFEError e s))
4 type QNFBuild e s = ListT (QNFBuild0 e s)

```

Listing 3.12.: QNF computation monad provides non-determinism, caching, and error handling.

Every QNF building function accepts one or two `NQNFQuery` forms and an operator unary or binary relational algebra operator to be applied to the query. It returns an `NQNFResult` (listing 3.13) contains

- A named QNF which is the main payload of the QNF building function.
- A mapping between input and output QNF columns. When there is no direct mapping (like in the case of aggregations), this is an empty list.
- The RA operator that was provided as input with the expression atoms translated to QNF names referring to the input or output columns. The way the names in the operator expressions are translated to `QNFNames` is specific to each operator. For example, in the selection operator, all atoms are directly translated using the name map of the input `NQNFQuery`. When the operator is a projection $\pi_{\{n_0 \mapsto K_0, n_1 \mapsto K_1, \dots\}}$, the builder function takes as input a QNF Q_{IN} and the parameters of the projection $\{n_0 \mapsto K_0, n_1 \mapsto K_1, \dots\}$, and outputs a `NQNFResult` containing the QNF Q_{OUT} . In the projection operator output, the left hand side atoms of the projection $\{n_0, n_1, \dots\}$ are columns of Q_{OUT} and the atoms in the right hand side expressions Q_0, Q_1, \dots are columns of Q_{IN} . This is for convenience assuming that the RA query that was just normalized is a logical plan. The top level operator can then be used to annotate QDAG elements that relate to the normalized query. The importance of this is analyzed in depth in section on symbol correspondence.

```

1 data NQNFResultDF d f op e s =
2   NQNFResult
3   { nqnResNQNF :: NQNFQueryDCF d Either f e s
4     , nqnResOrig :: op
5     -- Map the names exposed by the input qnf to names on the output
6     -- qnf. Unless qnfResOp is QProj, QGroup each input name corresponds
7     -- to 0 or 1 name in the output. Thus if qnfResOp is QProj or

```

```

8  -- qnfResInOutNames has no meaning and will be `const Nothing`
9  ,qnfResInOutNames :: [(e,QNFCol e s),(e,QNFCol e s))]
10 }

```

Listing 3.13.: The internal QNF building functions provide some more information that was created during the generation of the QNF, precisely a name map relating column names to QNF names, a map relating input QNF names to output QNF names, and the top level operator with the names translated appropriately to input or output QNF names.

The details of most QNF functions are fairly mundane and easy to infer with the possible exception of products (listing ??). When calculating the the product of two NQNFs Q_1 and Q_2 we first find conflicts, where columns of Q_1 can be found in the product collection of Q_2 . If none are found, we can simply merge all of their fields and recalculate the hashes. If there is at least one conflicting pair of **QNFNames** (c_1, c_2) , non deterministically increment the index of c_1 or c_2 . For example, if we want to compute the QNF of the product $\sigma_{a_0=1}(A_0) \times \sigma_{a_0=2}(A_0)$ we get in QNF-like form $\sigma_{a_0=1 \wedge a_1=2}(A_0 \times A_1)$ and also $\sigma_{a_1=1 \wedge a_0=2}(A_0 \times A_1)$.

```

1  def qnf_product_code(nqnf1, nqnf2):
2      # The named qnf is a pair of a map names→columns and a QNF.
3      nm1, qnf1 = nqnf1
4      nm2, qnf2 = nqnf2
5      # Find relations in the product set that could cause name
6      # conflicts.
7      conflicts = [c for c in qnf1.qnfProd if c in qnf2.qnfProd]
8      for c in conflicts:
9          # Python does not support non-determinism but with ListT we
10         # can do it in haskell. The `nondet` function non
11         # deterministically selects one of the arguments and returns
12         # it. The result of the overall computation is a stream of all
13         # possible results.
14         #
15         # The `increment_indices` function returns the same named CNF
16         # with all columns' indeces that corresond of the argument CNFProd
17         # incremented.
18         ncnf1, ncnf2 = nondet(

```

```

19         (increment_indices(ncnf1, c),ncnf2)
20         (ncnf1,increment_indicse(ncnf2, c)))
21     return ncnf1 + ncnf2

```

Listing 3.14.: The qnf product.

3.4.3. Pragmatism of QNFs

As mentioned, query equivalence is an NP-complete problem. QNF is only a heuristic that does fail to identify many cases of equivalent RA queries. For example, our design of QNFs does not capture the equivalence $\sigma_p A \cup \sigma_{\neg p} A = A$. It captures enough rewrite rules, however to be useful in most common cases. In particular, it covers the following RA properties:

- $\sigma_p \sigma_q A = \sigma_{p \wedge q} A$
- $\pi_m \pi_n A = \pi_{m \circ n} A$
- $A \bowtie B = B \bowtie A$
- $(A \bowtie (B \bowtie C)) = (A \bowtie B) \bowtie C$
- $\sigma_\theta(A \times B) = A \bowtie_\theta B$
- $\pi_m \gamma_{n,g} A = \gamma_{m \circ n, g} A$
- $\gamma_{m,g} \pi_n A = \gamma_{m \circ n, g \circ n} A$
- $A \Join B = B \Join A^*$
- $A \cup B = B \cup A^*$

3.4.4. Symbol correspondence

We need to normalize the names in each operator, but we also require that there is no naming ambiguity at any level. Why, then, is it necessary to normalize the symbols if they are unambiguous in each context? If we were only dealing with one query, we could declare that we trust that the user would be consistent when naming their columns (**select**

col **as** name **where** ...), and that the broader problem of RA equivalence is beyond the scope of this thesis and move on to more central issues for FluidB.

Normalizing the queries that correspond to n-nodes is of little use if we do not also normalize the symbols that refer to the columns of each n-node. To produce the code implementing an operator we need to be able to establish a stable correspondence between

- the symbols that refer to the columns of the input tables the symbols that refer columns that refer to the output the symbols in
- the predicates and expressions of the RA operators that t-nodes correspond to.

This is a problem unique to FluidB because most systems will produce variants of QDAGs that are specific to the query being planned. FluidB maintains a QDAG that must maintain its consistency between different queries.

The first step before inserting an operator in the graph is to normalize the symbol name. We consider that two symbols are equal only if they refer to the same column. This is done by means of converting the query into normal form as described previously.

This is required to resolve the possible ambiguity of symbols with the same name but different scope when expressions are moved around. It abstracts away the symbol names but creates some complications. For example, in the query $\sigma_{a=x \vee a=z}(\gamma(\sigma_{a<y}(A)))$ the a in $a = x$ and the a in $a < y$ are not equal symbols as they refer to different columns. The former refers to a column in the table $\gamma(\sigma_{a<y}(A))$ and the latter to the table A . They both refer to some transformation of a column labelled a . On the other hand a in $a = z$ and a in $a = x$ are indeed equivalent as they refer to the same column. To clear up the confusion we will say that symbols are *strictly-equivalent* if they refer to the same column but that they are *provenance-equivalent* when they refer to different columns but they are transformations of the same column, i.e., they would be represented by the same in a normal relational algebraic setting.

In particular, each materialized n-node must have a precise schema that represents the order of columns. The identities of these columns must be agnostic to the particular names decided by the user, the decorrelating subsystem, or any other query processor. *Otherwise it would be impossible for two clusters to share an input/output n-node would both be able to interact with it only in the case that they happen to be in total agreement about the column names of the relation that the n-node represents.*

Thus, to interact with the nodes inserted in the graph previously and to allow future queries to match the nodes we insert we split it's symbols' representation in two notions.

- A local notion of a symbol is subject to the names given by the circumstances of its creation (user-defined names, generated symbols, etc)
- a global notion that is a symbol representation dependent only in the column itself.

The local symbol is useful to find provenance equality should it be needed (being cautious when merging symbols based on them sharing a local notion that the global notion maintains correctness, i.e., it refers to the right column). It is clear that local notions of symbols are sensitive to the name chosen, the global is not.

Using the global notion of symbols creates a problem with the correspondence of the query schema of the output with the schemata on the input side of a cluster (see section 3.2 on the QDAG). In other words, if we throw away the names of the columns in the inserted operations in favor of column unique ids, how would the code generator know how to arrange the elements of each input tuple it encounters to generate the output tuple?

To mitigate that, we use the internal operations of the QNF query builder. To build a cluster, we use as parameters a NQNF along with an n-node for each input and output relation, and the operator with column references in terms of local notions of symbols. As described in detail in section 3.4.2 on building QNF, the result of building the QNF that corresponds to an n-node also returns the root operator with globalized symbols, as well as a correspondence between globalized symbols at the input with globalized symbols at the output side of operator.

Each cluster is marked with all relevant information required by the code generator. But is that enough? Not quite. As we have alluded to, each relation we deal with, and therefore each n-node, must be associated with a schema which contains information about:

- Which columns are available
- What is the type of each column
- Which columns are constant values
- What sets of columns form unique sub-tuples
- What is the order in which the columns are stored.

Thereby each node is coupled with a structure containing all this information (see section 3.6)

3.5. Query clusters internals

Clusters are the highest level building blocks of the graph. This section describes of the semantic role of the cluster as a set of connected n-nodes and T-nodes.

Clusters represent *logical operations* and they are the connection between the basic graph and the relational algebra semantics.

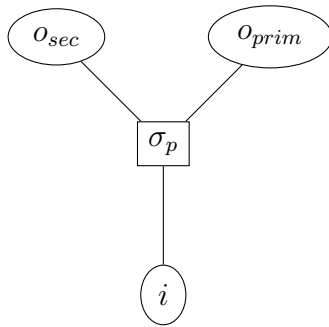
3.5.1. Cluster polymorphism

As alluded to when describing the QDAG in section 3.2, clusters come in different shapes depending on what the underlying operator is. In this section we will take a closer look at the structure of the different cluster types. Each type of cluster has input, output and intermediate slots, and each of the slots typically contains an n-node reference and an RA operator that relates the node to the cluster. They also contain slots to accomodate the t-nodes connecting said n-nodes. In dealing with clusters it is assumed that input and output n-nodes contain values that may be shared by other clusters. In contrast, intermediate n-nodes and t-node slots are always completely private to each cluster.

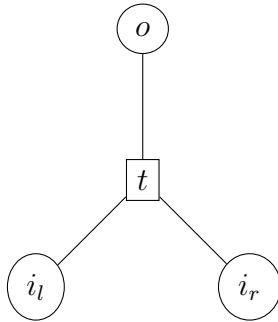
We define are three different kinds of clusters and a unified type **AnyCluster**:

- **NClust** is a cluster containing exactly one slot which corresponds to a primary table.
- **UnClust** a cluster corresponding to a unary operation. It has two slots in the output position, one slot in input position, no intermediate n-nodes and a single t-node slot (figure 3.6a).
- **BinClust** is a cluster corresponding to a non-join binary RA operation. It has one output slot, two input n-node slots and a t-node slot (figure 3.6b).
- **JoinClust** augments the **BinClust** with two intermediate n-nodes, two extra output n-nodes, and two input t-nodes to facilitate the antijoin results of a FluidB join operation (figure 3.6).

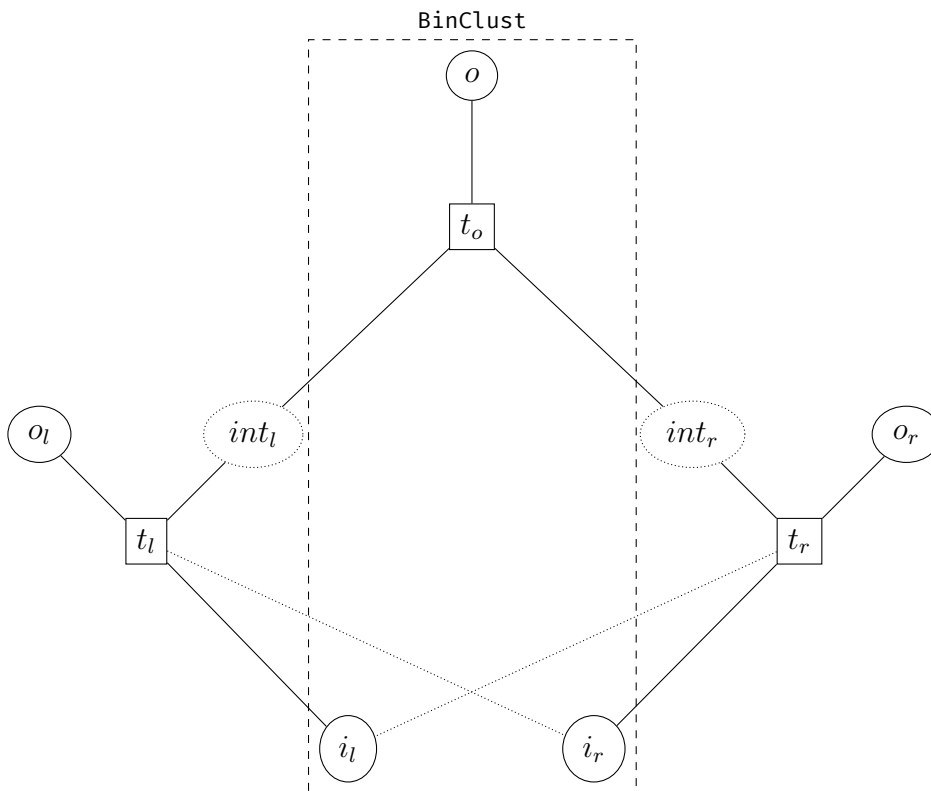
Every cluster must always be hashable and the hash is cached inside the cluster itself.



(a) The structure of a **UnClust**.



(b) The structure of a **BinClust**.



(c) The structure of a **JoinClust** contains a **BinClust** and includes two intermediate and two complementary output nodes to facilitate data related to antijoins.

56
Figure 3.6.: There are four cluster types. The **NClust** is trivial, just a single node. The other 3 are **BinClust**, **JoinClust**, and **UnClust** representing binary, join, and unary operations respectively.

3.5.2. Populating the graph with clusters

Once the query is transformed into a forest of possible join orderings as described in section 3.3.3 on enumeration of possible joins, the query is transformed into a forest of plans we can insert it into the tree (listing 3.15). This is an iterative process where each operation in each of the trees in the forest is translated to a cluster which is then inserted into the QDAG. The process of inserting a sub-forest or subtree into the QDAG has the effect of extending the QDAG by the connected clusters corresponding to the operators at the nodes of the forest, and also returns an n-node that corresponds to the subforest or sub-tree relation. However, since the insertion function is idempotent, i.e., inserting the same forest in the QDAG more than once has no effect either on the returned value or on the resulting QDAG, we memoize the function. Since the same function recursively inserts all sub-trees and sub-forests we get for "free" a mechanism of not traversing the same tree twice during insertion.

```

1 data QueryForest e s =
2   QueryForest
3   { qfHash :: Int
4     , qfQueries --
5       :: Either --
6         [(Query e (QueryForest e s))] -- An actual forest
7         (s, RelationShape e s) -- A final leaf
8   }

```

Listing 3.15.: The definition of the query forest. The query forest is hashed so that we can avoid traversing the same query forest repeatedly. The query forest is essentially a non-empty of queries with forests at their leafs.

The process of traversing the forest and inserting the clusters in the graph is fairly straightforward. The actual creation of the clusters, however, is slightly more convoluted. First, a cluster of a type corresponding to the query operation at the forest node is created (see section 3.5.1 on cluster polymorphism). That cluster initially contains QNF forms at its edges. The QNF forms are then matched against existing QNFs in the QDAG and the missing ones are generated. If no new nodes were generated and the matched nodes all correspond to the same existing cluster, we infer that the cluster is already in the QDAG and the process finishes. Otherwise, we need to create and insert a cluster and a corresponding propagator (see section 3.6.1). Finally all new nodes, QNFs, and the

cluster are registered and cross matched so they are easily cross referenced. We memoize the entire process using the [qfHash](#) to avoid re-traversing the same path.

As an example of what a more realistic QDAG looks like, first run through all 12 queries of the SSB TPC-H benchmark (see chapter 7) produces the graph presented in figure 3.7.



Figure 3.7.: A QDAG accumulated out of every node in the SSB TPC-H.

3.6. Relation shapes

In the section on QNFs we discussed the need to abstract away information that pertains to the particular structure of a relation in favor of semantic information that is rewrite and data independent. Relation shapes cover the opposite need: they encode the particular cardinality and shape of a relation at a specific point in time.

A relation shape is a datastructure that contains information about the shape of the result of query.

```

1 data RelationShape' e' =
2   RelationShape
3   { rsSchema :: [(e',ColumnProps)]
4     ,rsUnique :: NEL.NonEmpty (NEL.NonEmpty e')
5     ,rsSize :: RelationSize
6   }
7

```

```

8  type RelationShape e s = RelationShape' (ShapeSym e s)
9
10 data ShapeSym e s =
11   ShapeSym { planSymQnfName :: QNFName e s
12             , planSymQnfOriginal :: e
13             }
14 deriving (Show, Generic)

```

A shape contains all combinations of columns that comprise a unique subtuple. This is useful primarily to determine foreign key joins. The expression

$$Q_1 := A \bowtie_{bId_A=id_B \wedge cId_A=id_C} (B \times C)$$

behaves like a foreign key join while the join when id_B is the only element of the unique subtuple of B and id_C is the only element of the unique subtuple of C .

$$Q_2 := A \bowtie_{bId_A=id_B} (B \times C)$$

does not because for the expression $B \times C$ neither id_B nor id_C are unique on their own. Only their combination is. Therefore, the cardinality $|Q_1|$ is at most the size of A because each row of A matches at most one row of $B \times C$.

So the unique subtuples are computed using the following rules:

- In general the set of unique subtuples of $A \bowtie B$ is computed by concatenating all combinations of the unique subtuples of A and B .
- The unique subtuples of a selection, semijoin, antijoin, and limit are the same as the unique subtuples of the underlying relation.
- The unique subtuples of a projection are the unique subtuples of the input that are entirely exposed. The reader is reminded that FluidB's notion of RA requires that there is at least one such unique subtuple.

3.6.1. Shape propagators

Shape propagators are based on the the idea of the propagator as first introduced in [24] and analyzed more in depth in [64] and mentioned in section 3.2. Very briefly, a

propagator network is a hypergraph of cells containing mutable state and hyperedges as multi-directional functions that, when triggered, update the state of the connected cells so they are consistent with each other. The state of each cell represents some information about a value (like a bound or a probability distribution), and the propagator takes into account the state of each cell and propagates that information to the other cells (for example tightening their bound or finding a linear combination of alternative values based on the certainty of each). Kmett [72] formalized the notion of "information about a value" drawing from [38] encoding it as a lattice with \top meaning contradiction, and \perp meaning that there is no information about the value. Then the new values of cells are being joined (\vee) with the old ones and the value is updated to reflect the combination of information. The example provided by Kmett is each cell value being embedded in a 3 level lattice where the bottom level is **Nothing**, the middle level is **Just** x and the top level being a Haskell error (ironically denoted as \perp in literature unrelated to propagators).

As the name of our propagator implementation suggests, our idea is to correspond each cluster to a propagator to form a propagator network that computes the shapes of the relations corresponding to n -nodes.

Intuitively, we model a propagator as a partial function that changes the values at the edges of a cluster. Each cluster is equipped with a propagator which, when triggered synchronizes the values at the edges of the cluster. The partiality of the function stems from the fact that it might detect irreconcilable inconsistencies between the cells (see listing 3.16). This conception of the propagator departs slightly from the conception described in the paper in that it clearly separates the *cluster* as a datastructure that holds a fixed number of interdependent values, and the function that updates the values in that cluster.

```

1  type ACPropagator a e s t n =
2      EndoE e s (PropCluster a NodeRef e s t n)
3  type EndoE e s x = x → Either (AShowStr e s) x
4  type ShapeCluster f e s t n =
5      PropCluster (RelationShape e s) f e s t n
6  type PropCluster a f e s t n =
7      AnyCluster' (ShapeSym e s) (WMetaD (Defaulting a) f) t n
8  newtype WMetaD a f b = WMetaD { unMetaD :: (a,f b)}

```

Listing 3.16.: A propagator matches a cluster with shapes at the edges to the same kind of cluster with the shapes synchronized.

Given a cluster with n -nodes at the edges, we lookup the shape of each node and add that to the cluster edge. Once a cluster is triggered and the values synchronized those values are checked against the old ones and for the ones that were updated we find the clusters in which they participate and run the same process. The convergence of this process is justified in [38].

3.6.2. Defaulting functor

As mentioned, the cells of a propagator involve more structure than raw values. Indeed, FluidDB wraps these values in a functor that we dub as *defaulting functor (DF)* to reason about the consistency of the relation shape. In particular we require that the cell can handle cases where

- The shape is not computed yet, i.e., there is no information about the value.
- The shape is inferred via the propagator network and is therefore subject to change and
- The shape stored in the cell corresponds to a materialized node and other parts of the FluidDB's internal state depend on it.

The latter case is salient. Different logical plans can lead to different, yet semantically equivalent as per the FluidDB RA.

The DF may hold up to two values (see listing 3.17).

- An empty DF means we have no information about the value.
- A single value (the "default" value) corresponds to a derived shape that was computed using the propagator network. The default value is updated when propagators are triggered.
- A full value exists alongside a default value and is constant through its lifetime as it matches the value of some state that is external to the propagator network, in our particular application it represents the shape of a materialized relation. Until that relation is deleted, and the full propagator is demoted to a defaulting one, the registered shape is not allowed to change.

```

1 data Defaulting a =
2   DefaultingEmpty
3   | DefaultingDef a
4   | DefaultingFull a a

```

Listing 3.17.: The defaulting functor definition.

Since instances of the DF are meant to accommodate propagator values, adhering to Kmett's approach to cell value maintenance, we define a semilattice over the DF and define it in terms of a monoid (listing 3.18).

Semigroups are not commutative, so we give meaning to the order. The left operand of the semigroup operator is the newer value and the right one is the older value. Therefore, the DF does not require all the properties of the semilattice which is required to be commutative and associative. In general, the semigroup implemented by the underlying type of the DF is taken to mean *better* \diamond *backup*, i.e. keep as much of the left operand as possible unless the right operand is found to be "better" for some definition of better. For example, when calculating the cardinality we include a value of certainty between $[0, 1]$. A more certain right-hand cardinality is preferred over a less certain left-hand side.

Bringing our attention to the implementation of the **Def** \diamond **Full** combination in the semigroup definition 3.18 it is worth noting that it inverts the order of the combination of the default values. The reason we do this is to give priority to the full valued (materialized) nodes when propagating to their immediate neighbors. We need this when committing an operator to the code generation.

```

1 instance Semigroup a => Semigroup (Defaulting a) where
2   DefaultingEmpty <math>\diamond</math> a = a
3   a <math>\diamond</math> DefaultingEmpty = a
4   DefaultingDef a <math>\diamond</math> DefaultingDef a' = DefaultingDef $ a <math>\diamond</math> a'
5   DefaultingDef a <math>\diamond</math> DefaultingFull a' b' =
6     DefaultingFull (a <math>\diamond</math> a') b'
7   DefaultingFull a b <math>\diamond</math> DefaultingDef a' =
8     DefaultingFull (a <math>\diamond</math> a') b
9   DefaultingFull a b <math>\diamond</math> DefaultingFull a' _ =
10    DefaultingFull (a' <math>\diamond</math> a) b

```

```

11
12 instance Semigroup a => Monoid (Defaulting a) where
13   mempty = DefaultingEmpty

```

Listing 3.18.: The join semilattice that is defined in terms of the

When committing to a query plan we are also committing to full values for the defaulting values in the propagator cells. Every operator of the plan corresponds to a relation shape propagator. When committing an operator to the plan by translating it into C++ code, we promote all the defaulting values at the edges. It is clear that the generated code expects a certain data layout and therefore the *full parts of the values must be structurally consistent with each other*.

As we went over in the previous sections, however, the default values are consistent with other default values of the same cluster w.r.t. the amount of information they hold, but not w.r.t. the order of the columns in the relation shape. Therefore, before promoting the default functors to full functors by duplicating the default value to fill both fields of the **DefaultingFull** constructor, we trigger the propagator to enforce full consistency between the **DefaultingDef** values.

The full part and the default part of the the defaulting value are not necessarily structurally synchronized. We *internally synchronise* a cell before triggering the propagator to make sure we maintain the consistency of the newly promoted cells with the already existing full cells.

It is notable that there is no formal guarantee that the propagator will be able to create structurally consistent cells. In that case the propagator will fail. Handling of this failure gracefully is beyond the current scope of this work.

When an n-node is to be materialized, the corresponding defaulting functor containing the value is *promoted*. Promotion of propagators is copying the default value of a non-full propagator to the full state (listing 3.19).

```

1 promoteDefaulting :: Defaulting a -> Defaulting a
2 promoteDefaulting = \case
3   DefaultingEmpty      -> DefaultingEmpty
4   DefaultingDef x       -> DefaultingFull x x
5   d@(DefaultingFull _ _) -> d

```

Listing 3.19.: Promoting of defaulting functor happens during code generation when an n-node is materialized.

3.7. Summary and conclusion

In this chapter, we detailed how FluidDB processes queries. From parsing to generating a graph of queries (QDAG) and inferring the shapes of the relations in that graph. We also discussed different conceptions of relational algebra that are used in FluidDB (normal RA and QNF) and reverse operations. This will serve as a solid bedrock on which we can build the query planner and on top of the code generator.

There are a few shortcomings to this model that we have not yet addressed that we plan to implement in later incarnations of FluidDB. The most important ones being:

- As things stand, the QDAG will keep growing infinitely as new queries arrive and there is no obvious way to prune it while guaranteeing the materializability of all the nodes.
- There is no obvious path to supporting updates in the primary tables. There has been a lot of work on materialized view maintenance but most of it assumes that the primary tables are already materialized.

Physical planning

As you were undoubtedly gathering, the anomaly is systemic, creating fluctuations in even the most simplistic equations.

(*The Architect*)

Chapter summary

- We introduce a monad-based functional method of weighted backtracking search that supports *once* and *fallback*, two cut-like operations.
- The FluidDB query planner follows an A^* -like method for searching plans plans are searched in a top-down fashion, guided by a cost model that takes into account the cost historical queries.
- The query planner can make use of *forward* or *reverse* triggering of operators.
- The garbage collector clears up space on-demand maintaining the materializability of all the nodes in the graph.

This chapter goes in detail about the architecture of the logical planing of a query. Query planning is based on an A*-like search algorithm for through the space of partial plans. We initially describe the fundamental computational structures that facilitate this search and their unique properties that allow us the necessary flexibility. Then we discuss the basic backtracking algorithm and the way the graph is traversed. We move on to discussing the garbage collection that generates plan fragments that allow the plan to free up space in the underlying storage. Finally, we talk about the cost model that guides the evolution of the plan.

4.1. HCntT logic monad

The FluidB planner is designed in terms of a backtracking search algorithm that searches in the space of subplans for an optimal plan. Because the algorithm involves a lot of complicated heuristics it is important that a powerful underlying framework is deployed that matches the purely functional infrastructure in which FluidB is implemented. In particular, we require a monad that supports both weighted search and soft cuts. Because none of the solutions we found in the literature support both we developed the **HCntT** backtracking monad that is described in this section.

4.1.1. Background

Monads for backtracking in a functional context have been proposed in various incarnations. The most common monad that encapsulates **ListT** $m\ a = \text{forall } b . (a \rightarrow m\ b \rightarrow m\ b) \rightarrow m\ b \rightarrow m\ b$, also dubbed the Church encoding of lists or an application of the Cayley theorem on list. It was first proposed (in an untyped format in scheme) in [3] where they explicitly use it for implementing a logic programming framewrk back in 1987, but most authors in the field cite [13] as the seminal work on the concept where the authors seem to have independently re-discovered an application of the concept in 2000, and also [71] where the authors focused on a notion of fairness in backtracking. As shown in [65] this representation has $O(n^2)$ BFS complexity.

The other common representation of a list monad transformer, is the more straightforward

```
1 newtype ListT m a = ListT (Maybe (a, ListT m a))
```

Which mirrors precisely the idea behind **MonadLogic** from [71]. As demonstrated in listing 4.1 the **MonadLogic** able to **cons** (dubbed **reflect**) and **uncons** (dubbed **msplit**).

```

1 class MonadPlus m => MonadLogic m where
2   msplit :: m a -> m (Maybe (a,m a))
3   -- ... other methods are defined in terms of msplit
4
5   reflect :: Maybe (a,m a) -> m a
6   reflect Nothing = mzero
7   reflect Just (a,as) = pure a <|> as
8
9   -- msplit ==> reflect = return

```

Listing 4.1.: The logic monad typeclass

In [71], on top of this framework the authors build an **interleave** and a monadic **bind** (**>>=**) operator that works with it to make computation be slightly more "fair". In this context fair means that each of the branches gets equal computation resources. Interleaving is in a certain sense fair between two computations, but when more than two are involved, half the time goes to the first computation, out of the half that is left, half (1/4 of the total) goes to the second, out of the 1/4 left, half (1/8) goes to the third and so on. While this powerseries that describes the way the way resources are allocated to branches may seem arbitrary, in practice it may mean the difference between terminating and non-terminating computations. For example, the code in listing 4.2 can be translated to something like the code in listing 4.3 which does terminates due to the fact that while it does not give the same change to all the branches it does not completely starve any of them.

```

1 nonTerm :: [(Int,Int,Int)]
2 nonTerm = do
3   (a,b,c) <- (,,) <$> genNaturals <*> genNaturals <*> genNaturals
4   guard $ a + b - c == 10
5   return (a,b,c)

```

Listing 4.2.: Using a simple list to drive non-determinism is implicitly equivalent to a DFS algorithm which in many useful cases does not terminate.

```

1  interleaveTest :: [(Int,Int,Int)]
2  interleaveTest = runLogic @[] $ do
3    (a,b,c) ← genNaturals >*< genNaturals >*< genNaturals
4    guard $ a + b - c == 10
5    return (a,b,c)

```

Listing 4.3.: Interleaving (in this example `>*<`) is not *actually* fair in the sense that it does not give all the processes

Can we do better than unnecessarily avoiding non-termination? Sometimes we can with *weighted search*. Weighted search refers to the backtracking search where branches are weighted, or *prioritized*, and branches with higher priority are scheduled before branches with lower priority. A sketch of the API is demonstrated in listing 4.4. The backtracking monad implements the `halt` operation (called `tell` in [65], presumably to echo the `MonadWriter` interface) which accepts a value indicating the priority of the current branch and yields control to the scheduler. The scheduler then passes control to the branch with the highest priority. The value passed to `halt` must implement a monoid such that the priority of a branch is the concatenation of all values passed to `halt` in that branch up to that point.

```

1  -- | Non-deterministically return an integer but before that pass
2  -- control to the scheduler which prioritizes branches based on a
3  -- total sum which it tries to minimize.
4  stream :: (HeapKey h ~ Sum Int, IsHeap h, Monad m) => HCntT h r m Int
5  stream = go 0 where
6    go i = do
7      halt $ Sum i
8      return i <|> go (i+1)
9
10 -- | Non-deterministically choose three numbers and only keep the
11 -- branches that satisfy a + b - c == 10. Avoid diverging by preferring
12 -- branches for which the sum of the numbers is minimum.
13 test2 :: IO [(Int,Int,Int)]
14 test2 = takeListT 5 $ dissolve @ (CHep (NonNeg Int)) $ once $ do
15   (a,b,c) ← (,,) <$> stream <*> stream <*> stream
16   guard $ a + b - c == 10

```

```

17     return (a,b,c,x)
18 -- > test
19 -- [(5,5,0),(6,4,0),(5,6,1),(6,5,1),(6,6,2)]

```

Listing 4.4.: Prioritise branches that we want to be executed first.

None of the work we could find easily implements all the required features simultaneously so we implement yet another backtracking monad transformer, **HCntT**. With that in mind, for the FluiDB planner we require that our logic framework supports the following features:

Weighted search Not all plans that match our criteria, i.e. that solve the query within the space are equally admissible. We need to find as good plans as possible and we do not want the planner to spend time looking into plans that are unlikely to be efficient. For that reason neither breadth-first nor depth-first traversals are ideal for our purpose. We need a robust way to search in a weighted manner.

Soft-cut/fallback As we will see in more detail in the next section, the planner is initially optimistic about being able to materialize a node until it hits the budget limit. If the node can not be materialized within the available budget, the branch fails and a new branch tries running the garbage collector first and then trying to materialize the node. If the budget turns out to be large enough FluiDB should completely disregard the latter branch. The reason is that the later a GC pass happens the more options it will have and therefore the better job is likely to do w.r.t. deleting useless relations. To achieve that we implement the operator `<||>` (pronounced *fallback*) which is similar to prolog's soft-cut. As we will discuss in detail later, unlike other similar approaches our fallback operator refers to the *continuation* of the left-hand operand rather than the content of the operand itself.

Once In the context of non-weighted search it is fairly straightforward to implement an operator that demands that a sub-computation yields no more than one value (prolog's `once`). Simply run the entire computation in-place requesting one result and if it doesn't fail return that result. This approach can still work with weighted but we would like the `halt` calls inside the computation to have global effect for the scheduler. In other words we want the scheduler to be able to interleave this branch with other branches while preserving the semantics of **once**.

4.1.2. The **HCntT** monad

The **HCntT** monad uses delimited continuations and, similarly to [65] it traverses the search space in a best first manner by running the highest priority continuation available (see listing 4.5).

```

1  type HCntT h r m = ContT (HRes h m r)
2    (ReaderT (HeapKey v)
3     (StateT (CompState h r m) m))
4
5  newtype HRes heap m r = HRes (heap (Branch heap r m), [r])
6  type Branch h r m = ReaderT (HeapKey v)
7    (StateT (CompState h r m) m) (HRes h m r)

```

Listing 4.5.: The **HCntT** monad transformer allows continuation based non-determinism that allows switching between branches.

We define a recursive relationship between the monad within which the computation takes place –**HCntT**– and the value returned to the scheduler at each step –**HRes**. In particular **HRes** contains a priority queue (heap) of continuations and a set of final results. The heap type is parametric to allow flexibility with respect to how to best handle the particular key types (see listing 4.6). The heap is required to be stable, i.e. items with the same key are returned in the order they were inserted, and the **HeapKey empty** must be higher priority than all other **HeapKeys**.

```

1  class (forall v . Monoid (h v), Functor h, Monoid (HeapKey h), Ord (HeapKey h))
2    => IsHeap h where
3    type HeapKey h :: *
4
5    popHeap :: h v -> Maybe ((HeapKey h, v), h v)
6    singletonHeap :: HeapKey h -> v -> h v
7    maxKeyHeap :: h v -> Maybe (HeapKey h)

```

Listing 4.6.: We parameterize over heaps to allow the user to decide an efficient priority queue for the branches.

The **HCntT** needs to be *dissolved* to run the actual search. Dissolution (listing 4.7) is the process of turning an **HCntT** value to a **ListT** value. The **ListT** object will lazily produce

just as many results as are required, and of course just as many effects as required, and not more. This indicates that the process of building computations is composable but not incremental. The price of the `<///>` combinator that we described is that, unlike in other non-continuation based monad, once we start drawing results from the computation we can not apply further constraints on the resulting objects, at least not in terms of the **HCntT** object. Contrast that with the case of **ListT** where we can draw the first couple of results and then use the rest in a different computation.

```
1 dissolve :: (IsHeap h, Monad m) => HCntT h r m r -> ListT m r
```

Listing 4.7.: Disolution is the process of turning an **HCntT** computation into a **ListT**.

The way dissolution works then (listing 4.8), is to first commit to the computation constructed by applying the continuation and obtaining an **HRes** object. If that object contains concrete results they are yielded one by one into the resulting **ListT**. Then, if the heap is not empty, the highest priority branch is popped and scheduled to run until it yields a new **HRes**. The concrete results are yielded into **ListT** and the new heap is combined with the old one. Scheduling a branch entails running the reader layer of the monad transformer using its previous value. This is how the values provided to halt by the computation are accumulated within the branch. thereby incrementally updating the priority value.

As mentioned when describing the heap, a valid heap for **HCntT** must be stable. While this is a weighted search, it remains a question of how the options that have the same priority should be ordered. To avoid complicating things, up to this point we have been implying that continuations are accumulated in a single heap. But with just one heap, appending the newly produced heaps to the left would make for a depth first traversal of the *same-priority search space* while appending on the right would result in a breath first traversal. We want the operators to have the flexibility to decide on which side each branch should be appended. For that reason the **HRes** actually contains two heaps: one appended to the left (DFS), and one to the right (BFS). This is of particular importance to the correct implementtion of `<///>` as we will see shortly.

```
1 def dissolve(branches):
2     while len(branches) > 0:
3         priority,best_branch = branches.pop()
4         (sub_branches_left,sub_branches_right),results = best_branch.run()
```

```

5         branches = sub_branches_left + branches + sub_branches_right
6         for r in results:
7             yield r

```

Listing 4.8.: The dissolution algorithm in pseudo-python

In the following we will briefly describe the implementations of various combinators of **HCntT**.

Alternative/MonadPlus

The most fundamental combinator for any backtracking monad is the one spawning branches, the implementation of **Alternative** or, its equivalently in our case, **MonadPlus**. These typeclasses require that we can compose non-deterministic values ($a \triangleleft b$ or $mplus\ a\ b$ that non deterministically takes the value of a or b) and that we can make branches fail (`empty` or `mzero`). The implementation for **HCntT** is fairly straightforward (listing 4.9): for `empty` or `mzero`, which indicates failure of a branch, simply disregards the continuation and returns an empty **HRes**. The `mplus` or \triangleleft simply returns an **HRes** with only the two alternative branches having maximum priority (i.e. `mempty :: HeapKey h`) and bounded to the current continuation. Both those branches are pushed from the left so that they are scheduled before other same-priority branches.

```

1 instance (IsHeap h, Monad m) => Alternative (HCntT h r m) where
2     ContT m <|> ContT m' = ContT $ \f -> return
3         $ HRes (Tup2 (h (m f) <|> h (m' f)) mempty, [])
4     where
5         h = singletonHeap mempty
6     empty = ContT $ const $ return emptyHRes

```

Listing 4.9.: The implementation for **Alternative** is the same as the implementation for **MonadPlus**.

Halt: passing control to the scheduler

We mentioned the `halt HCntT` process. Because we want to be able to transform the **HCntT** monad we define the class of **MonadHalt** that can support this operation. Most common monad transformers of **HCntT** (like **ReaderT**, **StateT**, etc) can trivially support

`halt`. The implementation of `halt` for `HCntT` itself is demonstrated in listing 4.10. The provided priority is offset by the previous priority of the branch using the `ReaderT` monad transformer.

```

1 class Monad m => MonadHalt m where
2   type HaltKey m :: *
3   halt :: HaltKey m -> m ()
4
5 instance (Monad m, IsHeap h) => MonadHalt (HCntT h r m) where
6   type HaltKey (HCntT h r m) = HeapKey h
7   halt v = ContT $ \nxt -> do
8     v' <- asks (v <◇>)
9     return $ HRes (Tup2 (singletonHeap v $ nxt ()) mempty, [])

```

Listing 4.10.: The halt process yields updates the priority of the branch and yields execution to the scheduler.

Soft-cut/fallback

The `</>` (pronounced fallback), `left </> right` runs the left hand side operand (primary) and if no values are produced in the *entire* computation based on that, only then does it try to evaluate the right hand side (fallback). `HCntT` does not guarantee that the right hand side will be scheduled immediately after the primary branch completely fails, but it does guarantee that it will be scheduled before it moves on to a new priority value. In other words it is only guaranteed that the priority of the fallback branch will tie the last failing branch of the left hand side in terms of priority, but if there are other branches in the queue that tie, there is no guarantee of how those will be scheduled.

There are two challenges that our particular feature set imposes to implementing this:

- We want `</>` to operate on the entire operation, unlike [71] that makes the decision on whether to run the fallback solely based on whether the left hand side returns a value. For example in listing 4.11 the none of the branches of the fallback operator fail individually but the entire branch fails for one of the two.
- In the context of a weighted search, control needs to be able to escape a branch that passes through the left hand side operand before it is exhausted. Therefore

we can not simply dissolve the left hand side and reconstruct it. However, then how would the operator know when the branches are exhausted?

```
1  comp = do
2    -- The left computation succeeds on its own but the scheduler will
3    -- move on to the right branch because the branch fails later on.
4    x ← return 0 <//> return 1
5    guard (x > 0)
6    return x
```

Listing 4.11.: This computation will evaluate to 1 because, while the computation `return 0` always succeeds the branch fails.

We solve these problems with the use of a special kind of branch we call a *marker*, and a global state that keeps track of all markers in the branch heap. Since branches are arbitrary processes we equip them with access to global state (local to the computation), a lookup table (**CompState**) full of fallback branches (see figure 4.1). The main idea is that a fallback branch corresponds to an entry in the lookup table. This entry is modified accordingly by the children of the right hand side. At the time of the cut, we push a marker into the heap with priority lower than the child branches so that when it is scheduled, it will perform some checks based on the lookup table entry and schedule the fallback function if all the children branches have finished without yielding a concrete result. At any time there is exactly one valid marker per active fallback and it is denoted as such by the fallback entry in the lookup table. The rest of the markers are invalid and should be equivalent to noops.

More precisely about the internals of the marker processes, each one refers to a location in the lookup table via its closure. Also each marker is uniquely identifiable so each valid entry in the lookup table references one marker, the lowest priority one that corresponds to it. When a marker is scheduled, it looks up the fallback branch entry in the table and checks if the entry also refers to that marker. There are 3 possible scenarios that may play out at this point:

- The fallback entry in the table has been invalidated by a branch that yielded a result. In this case, the marker is invalid and just returns.
- The fallback entry is valid but does not correspond to the scheduled marker. This means due to the left hand side branch spawning low priority subbranches, another

marker has been inserted to (possibly) trigger the fallback at some time in the future. The current marker is invalid

- The fallback location is valid and corresponds to the scheduled marker. This means that the fallback process needs to be run and removed from the table.

What is the lifecycle of the fallback entries? For a visual aid to the following see figure 4.1: When an expression **A** `< // >` **B** appears we create a new fallback entry in the lookup table containing **B** and recursively "infect" all spawned sub-branches of **A** to perform the following actions immediately after they generate new branches and results:

- If there is at least one valid concrete result, invalidate the fallback in the lookup table and stop infecting child branches with the currently described hook.
- If the fallback is invalid it means there have already been valid results that rendered the fallback obsolete. Stop infecting sub-branches.
- If none of the above happened, check the *priority* of the last marker corresponding to the fallback (registered in the lookup table entry).
 - If it is strictly lower than the lowest priority subbranch do nothing because there is a well-placed marker to handle it.
 - Otherwise, create a new marker of the same priority as the lowest priority subbranch and put it in the right-append heap.

This way we know it will be scheduled *after* the last branch relating to the fallback.

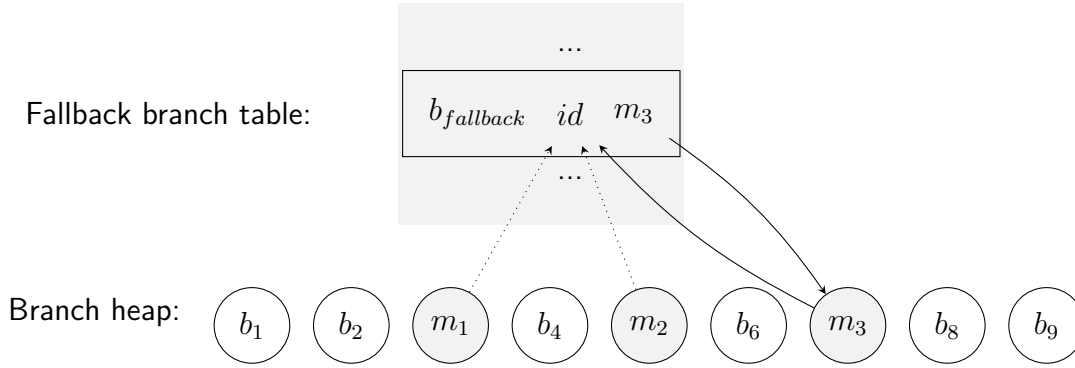


Figure 4.1.: Marker "branches" are injected in the heap of prioritized branches to possibly trigger the fallback branch. Each marker gives the scheduler the opportunity to run a particular fallback branch. The entry of the fallback branch entry references the last marker so that only that one actually trigger the fallback branch. A child branch that succeeds removes the fallback entry so the final marker also fails to trigger the fallback.

There are a few optimizations that can be implemented to avoid too many lookups in the fallback table in the case of deep fallback nesting that take advantage of the fact that new markers for outer fallbacks may only be created when new markers for inner fallbacks are created. However, for FluidDB we did not find **CntT** to be a performance bottleneck so we leave it for a future iteration.

Once

The **once** operator runs the argument computation and stops once it returns the first result. FluidDB uses **once** to run the garbage collector without exploding the search space, as the GC may have to search an exponentially large space of plan fragments. We prune that space by settling for the first result the planner can find and carefully crafting the GC such that it tries more robust plans first (see section 4.2.1).

The **once** operator is built on top of a concept we call a *nested scheduler* which is exactly what it sounds like. The nested scheduler is built from an **HCntT** value (the *nested computation*) and needs to know what to do in case of a success and in case of complete failure (no more branches to run), implemented in the general form as **nested** (see figure ??).

The nested scheduler then is a process that *always* returns a single subbranch which

has the priority of the highest priority subbranch of the nested computation. When scheduled by the outer scheduler, the inner scheduler internally schedules the next branch of the nested process. When the process yields results or completely fails, the corresponding hook is run. `once` implements these hooks as "return the result and stop" and as "propagate the failure" respectively. Since the implementation of `once` is fairly small it is provided along with the types in listing 4.12.

```

1  nested
2      :: forall h m r a .
3      (Monad m, IsHeap h)
4      => (Tup2 (h (Branch h r m)) -> r -> [r] -> HRes h m r)
5      -> HRes h m r
6      -> HCntT h r m a
7      -> HCntT h r m a
8  nested success fail c = ...
9
10 once :: forall h m r a . (Monad m, IsHeap h) => HCntT h r m a -> HCntT h r m a
11 once = nested (\_h r _rs -> HRes (Tup2 mempty mempty, [r])) emptyHRes

```

Listing 4.12.: The nested scheduler runs a subprocess within a single branch. `Once` is built on top of that to make sure the process stops once a result is returned.

4.2. The planner

The planner is the subsystem of FluidB that given the state of the QDAG and a target node to be materialized, produces a logical plan that will materialize the query. Our notion of a logical plan is slightly more specific than what is commonly considered a logical plan, i.e. the RA representation of the query. In our case it is a sequence of *transitions* that are to be transpiled to C++ by the code generator. There are three kinds of transitions:

- t-node trigger that assumes the input n-nodes are materialized and produces a subset of the output nodes.
- t-node reverse trigger that assumes that the output nodes are materialized

- n-node deletion

The planner operates by backtracking using **CntT** monad. Each branch maintains some branch-internal effects that are reified as monad transformers on top of the **CntT** defining the **PlanT** monad transformer (listing 4.13).

```

1 type PlanT t n m =
2   StateT
3     (GCState t n)
4     (ReaderT (GCCConfig t n)
5       (ExceptT (PlanningError t n)
6         (HCntT PlanHeap (()) m)))

```

Listing 4.13.: The monad that defines all the useful effects used by the planner. **GCCConfig** *t n* is an immutable, from the perspective of the planner, configuration that includes the QDAG, the node sizes, etc. **GCState** *t n* is state that is mutated and private to each branch of the planner like the materialized status of the nodes, the set of transitions registered so far and various caches. The **PlanningError** *t n* is a planner specific type of error. The entirety of the result of planning is accumulated in **GCState** so the result of backtracking is just unit **(())**.

It is important to clarify the way we use the term "materialized node" (**Mat** as opposed to "not materialized" – **NoMat**) from the planner's perspective. A mapping of node states is passed to the planner at the beginning of the planning process. This mapping indicates which nodes are initially materialized. Then as transitions get registered update the mapping is updated to reflect the effect that the plan so far would have on the materialized relation set in storage. Since the planner operates via backtracking (using the **CntT** monad described), each branch maintains its own mapping of node states and therefore considers a different set of materialized nodes.

We maintain the partial plan (sequence of transitions) as part of the state of each of the planners branches. *Registering a transition* means that a transition is added to the partial plan.

The main loop of the planer is a recursive process of asserting that nodes are materialized. If the node is already materialized the assertion succeeds, if not the planner attempts to trigger, or reverse trigger, a neighbouring t-node that would lead to the node being materialized, asserting first that the input nodes of that transition are materialized. In practice, due to intermediate n-nodes being auxiliary and not corresponding to real

nodes (see section 3.5 on clusters), we know that t-nodes are organized in sequences that need to be triggered entirely or not at all. For example (figure 4.2) the intermediate \bowtie nodes are not actually materialized as part of the join operator. This means that a plan triggering \bowtie t-nodes but not the corresponding \bowtie t-nodes is invalid. In other words, within a plan all intermediate nodes must be used as input exactly once for every time they are materialized. To ensure this we organize the t-nodes into **MetaOps** (listing 4.14) that can atomically be triggered or reverse triggered. **MetaOps** also abstract the distinction between triggering and reverse triggering as they expose a set of input, a set of output nodes and a process that registers the correct transitions once the **MetaOp** is triggered. The process of asserting an n-node being materialized then becomes the process of non-deterministically selecting a **MetaOp** with the node in question in the output set. Before actually splicing the **MetaOp** computation we assert that the input set is materialized.

```

1 data MetaOp t n = MetaOp {
2   metaOpIn      :: NodeSet n,
3   metaOpOut     :: NodeSet n,
4   metaOpInterm  :: NodeSet n,
5   metaOpPlan    :: forall m' . Monad m' => PlanT t n m' [Transition t n]
6 }

```

Listing 4.14.: A **MetaOp** refers to input, output, and intermediate nodes that are involved in the set of operations it abstracts. Furthermore it contains a computation that registers and returns the transitions involved in the **MetaOp**.

Two questions should arise from the above description: a) how does the planner avoid cycles where it recursively tries to materialize the parents and then the children? and b) how does it know not to materialize a node twice? Both those problems are addressed by refining the possible states of the n-nodes (listing 4.15). Rather than being just **Mat** or **NoMat**. We also disambiguate between **Initial** and **Concrete** nodes. Nodes in the **Initial** state are allowed to change their materialization status. In contrast **Concrete** nodes have their state fixed with very few exceptions. Then, when a node is asserted to be materialized its status is first checked and following scenaria are possible:

- the node's status found to be **Concrete** and **NoMat** and the branch immediately fails as it is not allowed to be set to materialized.

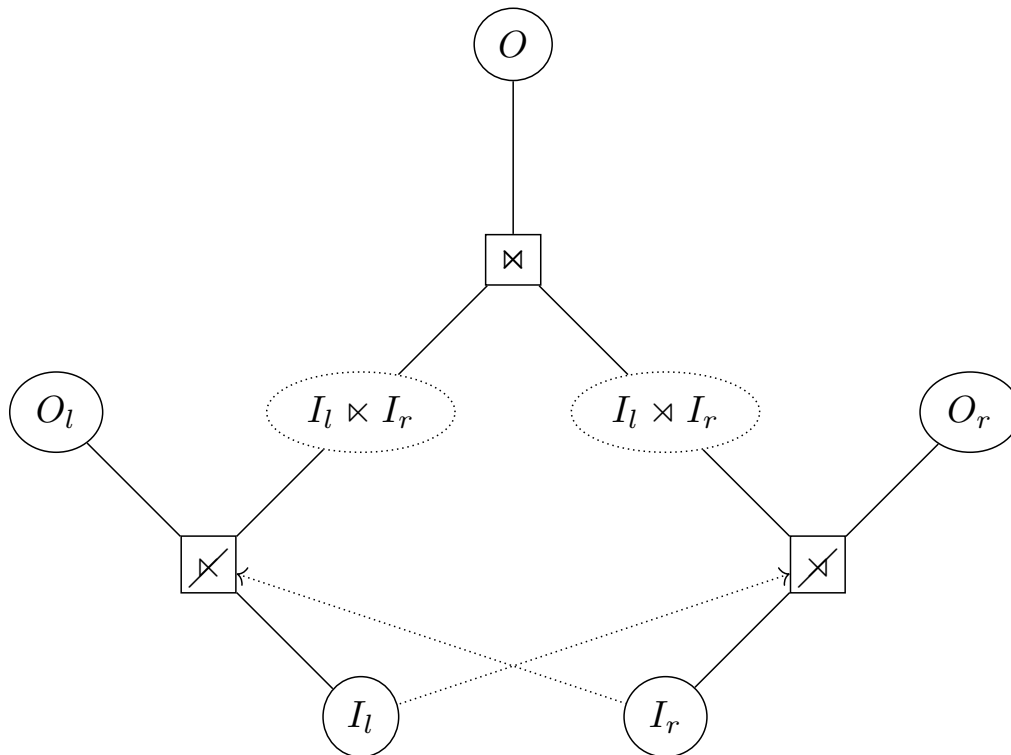


Figure 4.2.: starting at node O , the output of a join operation, we can derive three MetaOps that can materialize it. $MetaOp\{in : \langle I_l, I_r \rangle, out : \langle O \rangle\}$, $MetaOp\{in : \langle I_l, I_r \rangle, out : \langle O_l, O \rangle\}$, $MetaOp\{in : \langle I_l, I_r \rangle, out : \langle O, O_r \rangle\}$, $MetaOp\{in : \langle I_1, I_2 \rangle, out : \langle O_l, O, O_r \rangle\}$. Because trying all combinations of outputs explodes the search space we always go for the largest and then let the garbage collector deal with the possible reprecussions. On the other hand to materialize I_l there is only one **MetaOp** that relates to this cluster $MetaOp\{in : \langle O, O_l \rangle, out : \langle I_l \rangle, interm : \langle I_l \times I_r \rangle\}$

- if it is **Concrete** and **Mat** the assertion simply succeeds
- If it is **Initial** and **Mat** it is turned into **Concrete** and **Mat** and the assertion succeeds.
- if the node's status is **Initial** and also **NoMat**, two things need to happen: first the node is set **Concrete** and **NoMat** in order to avoid cycles, and the aforementioned process of finding a **MetaOp** is carried out recursively in order to materialize the node.

```

1 data IsMat = Mat | NoMat
2 data NodeState =
3   -- Concretified states are allowed to change only within the same
4   -- lexical scope.
5   Concrete IsMat
6   -- Initial states are subject to change when encountered as a
7   -- neighbour or by the GC.
8   | Initial IsMat

```

Listing 4.15.: The different states that a node is allowed to be in.

Once an n-node's is materialized and we start materializing its siblings, it is important that the garbage collector (which we will be discussing in section 4.2.1) does not delete said n-node. To avoid this problem, items materialized are set to **Concrete Mat** until all their siblings, that are inputs to the same **MetaOp**, are materialized. Once the **MetaOp** is triggered and its output nodes are set to the **Concrete Mat** state the input nodes are set to **Initial Mat** as it is now safe for the garbage collector to remove them.

It is worth highlighting the importance of each node in the depset being completely materialized before the process of materializing the next one commences, i.e., that no two sibling n-nodes are simultaneously in the process of being materialized. The reason is that at any time a single trail of parent nodes is marked as **Concrete** and **NoMat**, otherwise we can't be sure whether a **Concrete NoMat** node that renders a **MetaOp** non-triggerable is actually in the process of becoming materialized, meaning that when that process succeeds the **MetaOp** under consideration will be triggerable.

4.2.1. Garbage collection

One of the fundamental design decisions of FluiDB is that she is committed to materializing all intermediate nodes and keeping them around for as long as possible. In one

hand the plan selected is crafted so that the intermediate results are maximally useful for the overall workload, on the other, when the available storage runs out, the garbage collector mechanism selects the least useful nodes to be deleted, creating free required for solving the query. In this section we focus on the latter.

The garbage collector is triggered right before a **MetaOp** is to be triggered. It calculates the space required for materializing the output nodes and the space available to decide whether inserting a plan fragment generated by a GC run is required. Then it selects a subset of materialized nodes to delete in order to make room for the new results. The selection process has two hard constraints:

- The nodes being deleted must be *deletable*, i.e. the node's state is **Initial** and if the node is deleted that it can be reconstructed from the remaining materialized nodes.
- The nodes being deleted must not be required for the continuation of the current plan. We call these nodes *protected*. For example, say we are planning for the query $A \bowtie B$ and we have materialized A already but while materializing B the GC is triggered. A should not be in the repertoire of nodes the GC can delete under any circumstances.

Selecting a subset of nodes is not an easy problem so we follow a simple heuristic and leave a proper solution for future work.

Large nodes are costly to create and therefore the larger the node, the less inclined the planner is to create it, and the more useful it is likely to be. With that in mind the heuristic we follow is for the garbage collector to prioritise deleting small nodes over deleting larger ones.

The first order of business for the GC is to find the set of deletable nodes and sort them by size. It tries to delete them one by one starting from the smaller ones and working its way up to the larger ones. Every time a node is deleted it is possible that other nodes that were previously established to be deletable lose that property. For example if all nodes A , $\sigma_p A$, and $\sigma_{\neg p} A$ are **Initial** and materialized, each one of them can be materialized from the others and therefore all of them are marked as deletable. If the GC deletes $\sigma_p A$ first the other two are no longer materializable as they depended on A to be materializable.

Nodes that are established as non-deletable in this way are switched from **Initial** to **Concrete** to avoid re-calculating their materializability.

If after this process not enough space was created, the GC resorts to starting a new *planner epoch*. The function of a planner epoch change is for the planner to revisit assumptions that it has made in the form of setting n-nodes to **Concrete Mat**, i.e. when a new epoch begins the unprotected n-nodes materialized in the current run become available for deletion. When a new planner epoch begins all the nodes states and transitions that were recorded since the last planner epoch started are stashed into the epoch stack. The epoch stack resides in the branch-local state (**GCState**) and contains a map of nodes to their states at the end of the epoch and a sequence of transitions. When a new epoch is pushed into the stack a new mapping of states is created according to the following rules:

- The materialization status does not change between epochs: All materialized nodes from the previous epoch are materialized in the new one materialized and all non-materialized nodes are still not materialized.
- All non-protected **Concrete** nodes become **Initial** nodes.

The sequence of transitions for the new epoch is empty. The reason we keep separate lists of transitions is, as we will see in more detail in the code generation chapter, that the subset of output nodes that a physical operation generates is established at the time of physical planning based on the materialized nodes according to the epoch corresponding to the transition.

An epoch contains all the important information about the state of the planner. For this reason when an epoch is inserted in the stack it is checked for equality against the other epochs in the stack. If two equivalent epochs are inserted in the stack the branch fails permanently as it means that there is a cycle.

This enables another optimization, the *free-to-delete* nodes. As mentioned previously the planner prioritizes the version of **MetaOps** that materializes all the outputs for example it will prioritize the branch that follows $MetaOp\{in = \langle A \rangle, out = \langle \sigma_p A, \sigma_{-p} A \rangle\}$ over the one that follows $MetaOp\{in = \langle A \rangle, out = \langle \sigma_p A \rangle\}$. It is rarely the case, however, that all the outputs are required. In the example just mentioned then the former example, triggering the former **MetaOp** and then garbage collecting $\sigma_{-p} A$ without it being used is never desirable. For that reason we mark the node $\sigma_{-p} A$ as *free-to-delete* as soon as it is created based on the fact that it is not required. When the garbage collector deletes free-to-delete nodes it does so without registering a deletion transition. This way the

physical planner will generate a plan where $\sigma_{\neg p}A$ is never created in the first place. All nodes lose their free-to-delete status upon the creation of a new epoch.

The comprehensive algorithm for the GC is presented in listing 4.16

```

1  gc reqSize = do
2    -- Try the current epoch and if that fails retry with a new epoch.
3    return () <\/> newEpoch
4    -- find the deletable nodes and sort them by size
5    deleteables ← sortOnM getNodeSize ≡≡ filterM isDeletable ≡≡ getAllNodes
6    -- Try deleting each node and stop deleting when amassing enough
7    -- free pages.
8    forM_ deleteables $ \n → do
9      freePgs ← getFreePages
10     when (freePgs < reqSize) $ tryDelete n <\/> markAsConcrete n
11
12  tryDelete node = do
13    guardM isDeletable
14    isFreeToDel ← getIsFreeToDel node
15    unless isFreeToDel $ register $ DelNode node
16    setStatus node (Initial NoMat)
17
18  isDeletable n = do
19    st ← getNodeState n
20    case st of
21      Initial Mat → do
22        setNodeState n (Initial NoMat)
23        ret ← isMaterializable n
24        setNodeState n (Initial Mat)
25        return ret
26      _ → return False
27
28  markAsConcrete = do
29    st ← getNodeState n
30    case st of
31      Initial m → setNodeState (Concrete m)

```

32 `_ → return ()`

Listing 4.16.: A sketch of the garbage collector algorithm in pseudo-haskell. The `<///>` operator is supported by the **HCntT**.

The garbage collector may have a final trick up their sleeve when all else fails: neighbor materialization. When there exists a dependency set of a non-deletable node that takes up less size than the node itself, the transition, the GC can attempt to trigger the corresponding **MetaOp** in order to generate the dependency set and render the node deletable. For example the θ -join node $A \bowtie_{\theta} B$ is likely to take up more space than A and B combined. This makes the search space explode and is a generally low-yield strategy, therefore it is by default disabled.

It is hopefully clear at this point that the process of garbage collection involves a huge search space. We sacrifice a some of our plan repertoire by running the entire process wrapped in a **once** operator that we described in the previous section.

Finally, a word about protected nodes. Nodes are protected within the context of a branch from the time they are established as part of the input set of a **MetaOp** we are making triggerable until the time said **MetaOp** is actually triggered during the normal planning (i.e. not the GC). Because a node may be the input of more than one simultaneously considered **MetaOps**, protection of a node is not a boolean value that is set when the node is encountered as **MetaOp** input and unset when the **MetaOp** is triggered, but rather a natural number variable that is incremented and decremented respectively. When that value is zero, the node is not considered to be protected.

4.2.2. Order of traversal

We mentioned that the **MetaOps** are selected non-deterministically. In this section we will go in depth on the FluidDB planner's strategy on the order in which it considers its options using the **CntT** framework, making heavy use of the **halt** operator to decide branch priorities. A branch's priority is dependent on the particular frontier at the time of a halt and is determined by four factors:

- The cost of the **MetaOps** that the branch is in the process of making triggerable.
- The cost of the **MetaOps** that have already been triggered.
- The sum of the estimated costs of each node in the frontier.

- A weighted sum of the *stochastic cost* of historical queries based on the current frontier.

The final two factors are values, on which we will focus in this section, are computed often and are dependent on the set of materialized nodes. Since the set of materialized nodes does not change in a completely arbitrary way from computation to computation a naive approach to calculating those values would involve a lot of duplicate work. For this reason we use Antisthenis which will be expanded in a chapter 5 to minimize the amount of work required.

Let's start with the sum of estimated costs. We have a very rough way of estimating costs. We incrementally estimate the cost of each node of the frontier using the following formula:

$$c_{nomat}(n) = \min_{op \in \text{metaops with } n \text{ output}} \left\{ cost(op) + \sum_{i \in inputs(op)} c(i) \right\} \quad (4.1)$$

$$c_{mat}(n) = 0 \quad (4.2)$$

Where the cost $c(n)$ of n-node n is $c_{mat}(n)$ if the node is materialized and $c_{nomat}(n)$ otherwise. In plain english the cost is recursively estimated as cheapest combination **MataOp** plus the cost of materializing the input nodes.

This approach has several problems like the fact that nodes that are used more than once, are added more than once and that it does not take at all into account the budget constraints. It is, however, a good enough heuristic for prioritizing the branch.

The final factor, which takes into account the historical queries is a bit more tricky. The planner keeps track of the last couple of queries requested and the corresponding n-nodes, and tries to estimate how beneficial following the particular branch would be in the event where a query similar to those is requested in the future. We estimate that by summing up a notion of cost for those queries.

A naive approach would be to just use the the same algorithm as we did for the frontier nodes. However, this approach is prone to getting stuck behind materialized nodes. Past queries, especially recent ones are likely to still be materialized and therefore their cost is likely to be 0. Even if the GC has gotten around to deleting them, their near dependencies are likely to block the cost estimator go far. This makes the estimation quite bad for two reasons:

- Nodes found to be materialized when prioritizing branches may be deleted by a GC pass.
- We don't care about the cost of the particular nodes, but rather about *similar* nodes. For that reason we want to avoid depending too heavily on a particular materialized node being materialized.

A realistic way of calculating the expected cost that an n-node will have at a future point in time, which we very informally and heuristically attempt to approximate, would be to instead of considering the cost of materialized nodes to be zero, to calculate the likelihood that a materialized node will still be materialized when we encounter it again. This is related not only to an estimation of how many page-writes separate the moment of cost estimation and the actual plan the cost of which is being estimated, but also all the decisions that the garbage collector will make in that time. After FluiDB's budget is exhausted for the first time, for every page write a page needs to be garbage collected. We considered a couple of options for stochastically modelling the behavior of the GC like assuming that it chooses random pages or random tables to delete, but we could find none that was a pragmatic approximation and computationally viable.

For this reason we decided to follow a pragmatic approach to the problem and simply assume that for every materialized node there is a constant probability that it will not still be materialized when we need it, scaling the cost c_{nomat} of that n-node by that factor to get the stochastic cost. So the cost formula h for the nodes now is:

$$h_{nomat}(n) = \min_{op \in \text{metaops with } n \text{ output}} \left\{ h(op) + \sum_{i \in \text{inputs}(op)} cost(i) \right\}$$

$$h_{mat}(n) = \lambda \cdot h_{nomat}(n)$$

Where $\lambda \in (0, 1)$ is the estimated probability that the the node is no longer materialized when needed. As before materialized nodes have cost c_{mat} and non-materialized nodes have cost c_{nomat}

4.3. Conclusion

In this chapter we saw two distinct concepts: a weighted backtracking framework and a system implemented on top of that framework, namely the FluidDB's physical planner. These two combined enable FluidDB to search for quality plans within the QDAG in a flexible way. Some aspects of the planner, like the GC algorithm, may seem fairly naive. Since the planner is fundamentally a heuristic-based system, the design space is limited only by the imagination of the engineer designing it. The aim of this part of the work, and by extension of this chapter, is to explore ideas relating to the physical planner balancing breadth and depth.

Antisthenis

I can see a horse, and I can see
a man; but humanity and
horsehood I cannot see.

(Antisthenes)

You have the eye which sees a
horse and a man; but the eye
which can see horsehood and
manhood you have not.

(Plato)

Chapter summary

- *Antisthenis* is a framework for building systems for incremental evaluation of algebraic expressions using heuristics to optimize the order in which subexpressions are evaluated and to calculate some classes of self referential computations.
- *Mealy arrows* are a construct to implement adaptable computations that can produce a value and a new, semantically equivalent computation. Each iteration adapts to the previous caching values and changing the evaluation plan to be more efficient. Antisthenis computations are implemented as network of mealy arrows that compose larger adaptive (incremental and partial) computations.
- Machines are named Antisthenis processes that can be referred to by other machines. Antisthenis provides tools for handling some classes of self referential computations.
- We implement Antisthenis operations to cover requirements of the FluidDB physical planner:
 - calculate min/sum of natural numbers – used to compute *cost estimations*
 - boolean functions – used determine *node materializability*
 - calculate min/sum of numbers annotated with certainty metrics – used to calculate the estimated cost of materialization of nodes given nodes that *might* be materialized.

This chapter goes in depth describing the various parts that compose Antisthenis, a framework for building incremental computation systems we developed to assist the FluidDB planner in incrementally and efficiently computing costs and checking if nodes in the QDAG are materializable during the evolution of the inventory.

The chapter starts with an introduction about the motivation for developing Antisthenis as a collection of synchronously communicating processes performing incremental computation. Next we look at our conception of a Mealy arrow that enables the construction of Antisthenis processes. In the following section we describe the core of Antisthenis, which is the structure of a single Antisthenis process. From there we move on to describe some specific operation implementations that are useful for FluidDB and then describe the systems that facilitate intercommunication and high level organization of Antisthenis processes. Finally, we conclude with some caveats and some future work that can be done to make Antisthenis easier to apply to solutions other than FluidDB.

5.1. Introduction

The order in which subexpressions of an algebraic expression are evaluated can be important not only for the performance of the expression evaluation but, in the context of infinite expressions, for the very termination of the evaluation itself.

Expressions are typically modelled as abstract syntax trees or computation graphs (DAGs), and evaluation is commonly modelled as a reduction of that tree or DAG. Depending on the evaluation strategy each operand is either fully or partially evaluated. Typically, the strategy for evaluating a subexpression operand, however, is oblivious to the parent expression and a partially evaluated subexpression (thunk) is opaque from the perspective of the parent expression. Breaking this assumption provides some opportunities for incremental computation and efficient evaluation in the presence of non-total arguments.

We built *Antisthenis* around this principle to achieve three important goals:

- To *incrementally evaluate* an expression, i.e. to only re-evaluate the parts of it that depend on parameters that were updated with respect to the last time the expression was evaluated.
- Take advantage of absorbing elements and other operator-specific strategies for *early stopping* the evaluation.

- Often sub-expressions are not total, i.e. they may not be computable. A common reason for that is the expression being self-referential. These errors may be detrimental to the evaluation the top level expression, but often, as we will see, they are not.

Incremental computation is a tool for *efficiently* evaluating expressions of variables that change over time. The area has received some attention recently [32, 41]. An incremental evaluator will typically accept a directed acyclic graph (DAG) of interdependent computations like the following one

$$A = a + B + C + D$$

$$B = C \times b$$

$$C = D + c$$

$$D = 0$$

Which can be represented as a dag in figure 5.1.

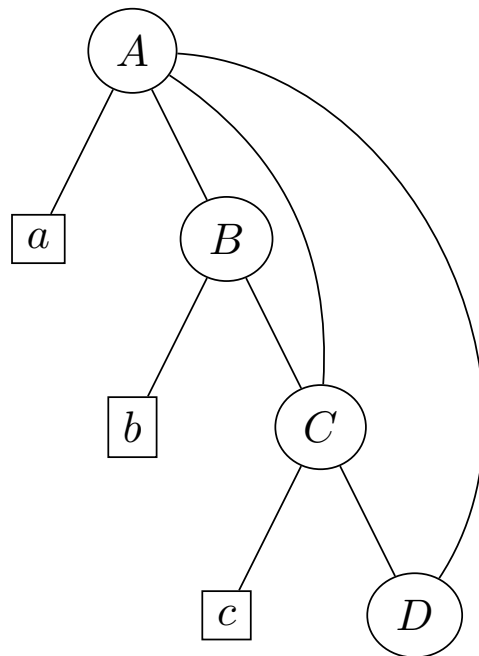


Figure 5.1.: A computation DAG. Circular nodes represent operations and square nodes represent parameters that are subject to change.

Where A , B , C , and D are expression nodes in the DAG and a , b , and c are parameters

that change over time. A non-incremental evaluation system would evaluate all four nodes every time a value were requested, regardless of whether and which parameters have changed. Incremental evaluation systems keep track of which variables change and only re-evaluate the nodes of the DAG that have been updated. In our example if only the parameter b has been updated since the last time A was evaluated, an evaluation of A would require only A and B to be re-evaluated.

Incremental evaluation systems employ a wide range of tricks from simply memoizing the functions [5] to more complex strategies like [41] where changes in parameters are propagated in the form of “dirtying” DAG nodes and more explicitly breaking the computation into reusable parts.

To our knowledge none of the existing approaches to incremental computation take into account the presence of either absorbing elements or non-computable subtrees (i.e. recursive values) in their scheduling the evaluation order. Scheduling is either done naively or left up to the user. However, it is a hard requirement of the FluidDB query planner for the incremental computing framework to automatically exploit the properties of the computation at hand in order to incrementally compute node materializability and query costs in the presence of different materialized query sets. Furthermore, in all the approaches we are aware of the network of sub-expressions was assumed to be well-behaved, in that all nodes can return a valid value and there are no cyclical references.

The main principle behind the architecture of Antisthenis is to manipulate the order in which sub-expressions are evaluated in order to exploit properties of the parent operator and avoid fully evaluating each of the sub-terms. For example when estimating the cost of a database query we take it to be the minimum cost of each of the possible logical plans. It is unlikely that one would need to come up with a full estimation of each plan, especially the more expensive ones, before realising the lowest cost. Instead, we would like to accumulate a lower bound for the cost of each plan.

Another prime example of an opportunity pruning the expression tree is the exploitation of absorbing elements. In the general case, in a non-lazy evaluation strategy, to compute the value of a function $f(A, B, C)$ where A , B and C are expressions, we would need to fully evaluate *all* three arguments first. Since the absorbing element of multiplication in the reals is 0, for $f(a, b, c) := a \times b \times c$, a result of 0 for the first argument A would render the evaluation of the rest of the arguments redundant.

As alluded to, a lazy evaluation strategy paired with a “Sufficiently Advanced Compiler”, can take steps towards mitigating this problem, but it can only go so far. Consider

the following example:

$$A = B \times C \times D$$

$$B = \sum_i i$$

$$C = 10 - 10$$

$$D = \sum_i i$$

Here, when evaluating A , a traditional evaluation system would look at the first argument first, namely B . B is very expensive to evaluate and very soon a human would figure out that it will never evaluate to zero. C , however, does evaluate to zero, and once that is evaluated the evaluation of A is complete. In a complex set of automatically generated such expressions it is important that the evaluator does not keep trying to evaluate a complex expression when evaluating another simpler one could yield an absorbing value.

Early stopping a computation is not useful only to avoid extra work. As we mentioned, many incremental computation systems assume that all sub-computations terminate successfully. Due to the presence of cycles in FluidDB graphs we need to take seriously the case where a sub-expression will fail to terminate. In this case the parent expression needs to exhaust all possibilities that could lead to a value.

$$A = \min(B, C, D)$$

$$B = b_1 + b_2 \cdot D$$

$$C = c_1 + c_2 \cdot A$$

$$D = d_1 + d_2 \cdot B$$

$$b_1 = b_2 = d_1 = d_2 = 1$$

$$c_1 = 3$$

$$c_2 = 0$$

Here B and D are clearly not computable but in the context where these values represent cost of evaluations for nodes A should evaluate to 3 and A should evaluate to

$C = 3$. Furthermore, A should depend *all* parameters for the validity of its value while C should depend only on c_1 and c_2 .

5.2. Background

5.2.1. Arrows and profunctors

The basic element of computation in the purely functional world is so far related to the troika of *functors*, *applicatives*, and *monads* [25], in the interfaces of which we find some recurring sub-structures that relate their modification.

- $a \rightarrow b$ the normal pure function type. We take this for granted as it is provided and managed by the language, in our case Textttell.
- $f\ a \rightarrow f\ b$ A functor morphism. The standard functor stack we saw so far focused on producing such morphisms from simpler arrows.
- $f\ (a \rightarrow b)$ is known as the Cayley morphism
- $a \rightarrow f\ b$ which for an f being a monad is the now familiar Kleisli arrow.

We would like to generalize this notion to an parametric object that we will denote as \rightarrow (occasionally denoted as \mathbf{p}). We follow the arrow definition from Hughes [17] which dictates that an arrow must be a **category** (i.e. arrows must compose in a commuting way and there must be an identity arrow) and must also satisfy the following:

- it must commute with tuple (i.e. must implement `first` and `second`).
- The Hask category must be embeddable in the category defined by the arrows (i.e. via the `arr` function)

In more concrete terms we depend on the haskell typeclasses presented in listing 5.1. An **Arrow** can act on the first or second element of a tuple, and an **ArrowChoice** can act on the the left or the right case of an **Either**.

A related concept to arrows is the profunctor. Profunctors are not necessarily directly composable like arrows nor to they commute with tuple or `either`. They only provide `dimap`. `dimap` is of course implementable in terms of arrow combinators but opting for `dimap` can produce much more efficient code.

```

1  class Category p where
2      id :: p a a -- Maps a value to itself with no effects
3      (.) :: p b c → p a b → p a c -- composition
4
5  class Category c ⇒ Arrow p where
6      arr :: (a → b) → p a b -- arrows do at least what functions can do.
7      (>>>) :: p a b → p b c → p a c -- also composition
8      first :: p a b → p (a,c) (b,c)
9      second :: p a b → p (c,a) (c,b)
10     (***) :: p b c → p b' c' → p (b,b') (c,c')
11     (&&&) :: p b c → p b c' → p b (c,c')
12
13  class Arrow p ⇒ ArrowChoice p where
14      left :: p a b → p (Either a c) (Either b c)
15      right :: p a b → p (Either c a) (Either c b)
16      (+++) :: p a b → p b' c' → p (Either a b') (Either b c')
17      (|||) :: p a c → p b c → p (Either a b) c
18
19  class Profunctor p where
20      dimap :: (a' → a) → (b → b') → p a b → p a' b'

```

Listing 5.1.: Haskell typeclasses related to the notion of **Arrow**.

All implementations of the interfaces in listing 5.1 are expected to conform to some laws which formalize the intuitions described about the meaning of arrows. We describe the laws required for each typeclass to be valid in listings 5.2

```

1  -- | Laws that need to hold for any valid category implementation.
2  -- Right identity
3  f . id = f
4  -- Left identity
5  id . f = f
6  -- Associativity
7  f . (g . h) = (f . g) . h
8
9

```



```

10 -- | Laws that need to hold for any valid arrow implementation.
11 arr id = id
12 arr (f >>> g) = arr f >>> arr g
13 first (arr f) = arr (first f)
14 first (f >>> g) = first f >>> first g
15 first f >>> arr fst = arr fst >>> f
16 first f >>> arr (id *** g) = arr (id *** g) >>> first f
17 first (first f) >>> arr assoc = arr assoc >>> first f
18
19 where
20
21 assoc ((a,b),c) = (a,(b,c))
22
23 -- | Laws that need to hold for any valid arrow choice implementation.
24 left (arr f) = arr (left f)
25 left (f >>> g) = left f >>> left g
26 f >>> arr Left = arr Left >>> left f
27 left f >>> arr (id +++ g) = arr (id +++ g) >>> left f
28 left (left f) >>> arr assocsum = arr assocsum >>> left f
29
30 where
31
32 assocsum (Left (Left x)) = Left x
33 assocsum (Left (Right y)) = Right (Left y)
34 assocsum (Right z) = Right (Right z)
35
36 -- | Laws that need to hold for any valid profunctor implementation.
37 dimap id id ≡ id
38 dimap (f . g) (h . i) ≡ dimap g h . dimap f i

```

Listing 5.2.: Laws for the typcalsses related to arrows.

Examples

The definition and highlighting of important properties of arrows that we have done so far is probably not extremely helpful to understand how computation is defined in such terms. To make the concepts of arrows and profunctors more tangible consider the following examples:

The Kleisli arrow The **Kleisli** arrow [20] as we mentioned earlier, is equivalent to the arrow **Monad** $m \Rightarrow a \rightarrow m b$. In haskell it is defined as:

```
1 newtype Kleisli m a b = Kleisli (a → m b)
```

So the arrow **Kleisli** (**State** s) $a \rightarrow b$ is an arrow that when mapping values from a to b also operates on a state value of type s . Cpmposing two Kleisli arrows is equivalent to binding two $a \rightarrow m b$ objects using \Rightarrow . The identity of the Kleisli category is equivalent to the monadic **return**.

The state arrow We already saw the **State** monad in the previous section. If we expand the **Kleisli** (**State** s) $a \rightarrow b$ arrow:

```
1 Kleisli (State s) a b
2   ~ = a → State s b
3   ~ = a → (s → (b,s))
4   ~ = (a,s) → (b,s)
```

The final way of putting it skips a few intermediate steps, we define a **StateArrow** then as

```
1 newtype StateArrow s a b = StateArrow ((s,a) → (s,b))
```

Which is equivalent to the Kleisli arrow we created before: an arrow that when mapping values from a to b also operates on a state value of type s .

The **StateArrow** is not unique in its equivalence to a Kleisli arrow. This is a property that a couple of useful arrows have and that we will be exploiting in the following sections. Below is the interface we provide to describe the Kleislifiable arrow (named **ArrowFunctor** below) and a few examples

```

1  -- | Some arrows correspond to Kleisli arrows. We should be requiring
2  -- a monad functor but
3  class Functor (ArrFunctor c) => ArrowFunctor c where
4      type ArrFunctor c :: * -> *
5
6      toKleisli :: c a b -> a -> ArrFunctor c b
7      fromKleisli :: (a -> ArrFunctor c b) -> c a b
8
9  -- | A kleisli arrow is trivially a kleisli arrow
10 instance Functor m => ArrowFunctor (Kleisli m) where
11     type ArrFunctor (Kleisli m) = m
12     toKleisli (Kleisli c) = c
13     fromKleisli = Kleisli
14
15
16 -- | A state arrow is a state arrow with a bit of rearrangement.
17 instance ArrowFunctor c => ArrowFunctor (StateArrow s c) where
18     type ArrFunctor (StateArrow s c) =
19         (StateT s (ArrFunctor c))
20     toKleisli (StateArrow c) a = StateT $ \s -> swap <$> toKleisli c (s,a)
21     fromKleisli f = StateArrow $ fromKleisli $ \s,a -> swap
22         <$> runStateT (f a) s
23
24 -- | A writer arrow is also a kleisli arrow on top of a writer monad
25 -- with a bit of rearrangement.
26 instance ArrowFunctor c => ArrowFunctor (WriterArrow w c) where
27     type ArrFunctor (WriterArrow w c) =
28         WriterT w (ArrFunctor c)
29     toKleisli (WriterArrow c) = WriterT . fmap swap . toKleisli c
30     fromKleisli c = WriterArrow $ fromKleisli $ fmap (fmap swap . runWriterT) c

```

5.2.2. Mealy arrows: processes that remember

The concepts that we will dub a *mealy arrow* or *mealy process*, is kin to many different ideas: transducer, automaton, and a coroutine to name a few. A mealy arrow (listing 5.3)

is a function that along with the result gives a new version of itself. Like an *automaton* after every iteration it moves to a new state ready to continue the process. Like a *transducer* every element of a stream consumed changes a hidden state to be taken into account when consuming the next one. Like a *coroutine* it can be conceptualized as a process that can yield computation to be resumed by the caller.

```
1 newtype MealyArrow a b =  
2   MealyArrow (a → (MealyArrow a b,b))
```

Listing 5.3.: Haskell definition of a mealy arrow.

The mealy arrow is at the heart of Antisthenis as it represents the basic building block of computation. We call the arrow returned by a Mealy arrow along with the value an *iteration*. Assuming that a mealy arrow represents an incremental computation, every iteration is produced in such a way as to exploit information about the structure of the problem based the previous computation.

For example, a mealy arrow calculating a multiplication of several sub-arrows may want to try the ones that evaluated to zero in the previous iteration in order to exploit possible domain knowledge that nodes that are zero are likely to remain zero in future iterations.

Mealy arrow transformer

Like monads [8] , arrows can be composed to form composit arrow types when they are represented as *arrow transformers* For brevity we will only talk about mealy arrow transformers but the concept is generalizable to many different kinds of arrows [58].

The motivation behind "transformizing" the mealy arrow is that arrow evolution as in the plain mealy arrow we presented must be pure (due to it being facilitated by `→`). As we will see in detail when describing the internals of Antisthenis, as subprocesses evolve they need to interact with external data structures, may throw irrecoverable errors, need to run monadic computations from external libraries etc. None of that can be done based on a pure haskell function. For that reason we change the mealy arrow definition slightly (listing 5.4)

```
1 newtype MealyArrow c a b =  
2   MealyArrow (c a (MealyArrow a b,b))
```

Listing 5.4.: A MealyArrow can take on the properties of other arrows by swapping out the function \rightarrow type for a parametric one.

Here the mealy arrow is parameterized by an arbitrary arrow `c` which may be a Kleisli arrow or any other arrow which will define what side effects are supported during the evolution of the process.

Building mealy arrows

While arrows are powerful and general, they can be slightly awkward to work with even with the haskell `proc` syntactic sugar [30]. Transferring state between mealy iterations via their closures even more so. Capitalizing on the parallels between an arrow and a coroutine, and taking advantage of the rich haskell ecosystem around monads we define a monad `MB a b m` to help us build mealy arrows more easily. The interface to that monad is presented in listing 5.5.

```

1 yieldMB :: Monad m => b -> MB a b m a
2 mkMealy :: (ArrowFuncor c, Monad (ArrFuncor c))
3         => (a -> MB a b (ArrFuncor c) Void)
4         -> MealyArrow c a b

```

Listing 5.5.: The `MB a b m` monad can be used as a convenience to implement mealy arrows using a conroutine-like interfaces.

The important part is that we want to do something before each iteration, something after, and we want to feed some state to the next iteration like presented in listing 5.6, where it is how how we can transfer variables between iterations via the closure of the `do`-block.

```

1 foo :: ArrowFuncor c => MealyArrow c a b
2 foo = wrapMealy arr $ \a it -> do
3   doSomeStuff
4   (it',r) <- lift $ toKleisli (runMealyArrow it) a
5   a' <- yieldMB r
6   return it'

```

Listing 5.6.: An example the usage of the MB functor to generate mealy arrows.

This involves `toKleisly` in reference to the "kleislifiable" arrows we saw in the previous section and lifting to `MB a b m` arrow which is precisely defined in 5.7. The basis here is the functor `MealyF a b` which structurally looks like a non-recursive version of a mealy arrow we saw earlier. That is if `x` is replaced with `MealyF a b _` we get a mealy arrow with an extra value. Unfortunately there is no direct way to make a law abiding monad out of the `MealyF` type but there is a trivial functor for it, which allows us to use the free monad [22] to get the monad we want. Note how the Mealy arrow is *not* kleislifiable itself. The `MB a b m` monad can't produce any mealy arrow. We can produce only arrows that are from `a` to `b` from `MB a b m Void`.

```
1 -- | Constructing monadic xmealy arrows
2 newtype MealyF a b x = MealyF (a → x,b)
3 type MB a b = Free (MealyF a b)
```

Listing 5.7.: Definition of the MB monad transformer.

5.2.3. Zippers

Now that we have a fairly coherent view of computations we need one more piece of background to start putting together Antisthenis: the zipper data structure. A zipper [10] is an auxiliary data structure that is similar to (but not exactly the same as) a cursor.

For our purposes a zipper is a datastructure that represents an alternative configuration of a container such that it focuses on a specific location in that container, allowing an interface for reading and modifying the element at that location. A zipper should also be able to shift focus to adjacent elements.

Zippers can be very complex data structures with interesting categorical properties arising from the fact that they can implement a comonad interface [18] but for our purposes we only need a rudimentary understanding of it.

As an example, consider the structure `ListZipper`:

```
1 data ListZipper =
2   ListZipper
3   { lzLeft  :: [a]
4     , lzCur :: a
5     , lzRight :: [a]
6   }
```

A list zipper breaks the list in three parts: a left part `lzLeft`, a focused element `lzCur` and a right part `lzRight`.

```

1  -- Make a zipper from a non-empty list
2  mkZipper :: [a] → Maybe (ListZipper a)
3  -- get a list from the zipper
4  getList :: ListZipper a → [a]
5  -- Focus on an element on the left if there is one
6  moveLeft :: ListZipper a → Maybe (ListZipper a)
7  -- Focus on an element on the right if there is one
8  moveRight :: ListZipper a → Maybe (ListZipper a)
9  -- Modify the current element
10 modCur :: (a → a) → ListZipper a
11 -- Read the element in focus.
12 getCur :: ListZipper a → a

```

These are fairly self-explanatory, a zipper can shift focus from the element in question to the neighboring elements.

5.3. Antisthenis core

In this section we discuss the core components of Antisthenis. The parts that comprise Antisthenis are roughly split in two categories.

- The *core components* that are related to building the Antisthenis processes, the cells of computation, and the interfaces that need to be implemented to specialize them to specific operations.
- and the *external components* that generally relate cells to each other and the systems that specialize processes to implement specific operators.

While in this section we focus on the core components we will frequently allude to external components which we will see in details in sections 5.4 on operators and ?? on named processes.

5.3.1. Processes

An Antisthenis process is a node in the evaluation graph and it is built using the concepts described so far in order to implement the following properties:

- *Parameterize* the computation by a cascaded configuration (as we will see **Conf** _w)
- *Adapt* the computation between runs using Mealy arrows.
- Return metadata about the result *deliniating its validity* (as we will see **CoEpoch** _w)
- Disambiguate between an *uncomputable*, a *partial*, and a *concrete* result (as we will see **BndR** _w)

As we described in section 5.2.2 mealy arrows define the basis of an Antisthenis process. Indeed our definition of an Antisthenis incremental cell is a composition of **MealyArrow**, **WriterArrow** and **Kleisli** presented in 5.8

```

1 type ArrProc w m =
2   MealyArrow (WriterArrow (ZCoEpoch w) (Kleisli m))
3   (Conf w) (BndR w)

```

Listing 5.8.: The type of an Antisthenis process is a mealy arrow paired with a kleisli arrow.

The **w** parameter is simply a token related to the operation that the cell is evaluating that we use helps us parameterize the types involved in a process by leveraging the type family feature of haskell's type system [69]. We call it a **ZipperParams** token beacuse, as we will be seing throughout the chapter, it mainly parameterize the evolution of the internal state of the Antisthenis process, the **Zipper**. We will be looking in detail into the methods that **ZipperParams** tag needs to implement in the following sections.

We will describe every type involved in the **ArrProc** type separately, but to do that we first provide a general intuition of what these types mean. The basis is a mealy arrow that emits monoidal value via **WriterArrow**, along with the values, we dub **ZCoEpoch**. It is a datastructure that summarizes the provenance of the value returned (section 5.3.3 on Antisthenis caps sand bounds). It is used to determine when the value is no longer valid due, for example, to changed parameters.

The **MealyArrow** and **WriterArrow** combined transform an arbitrary Kleisli arrow which is used to embed the computation of Antisthenis processes into parent monadic computations, which in the case of FluidDB is the **PlanT** $t \rightarrow n \rightarrow m$ monad as we saw in section 4.2. In the FluidDB case that would be the monad computation associated with planning. This way we give the Antisthenis processes the means to lookup values, throw errors, or generally have effects of any computation we require.

The output side of the **ArrProc** is a **BndR** w type, the returned value, which is parametric to the operation but has some common structure (listing 5.9). This indicates that every Antisthenis process has three options for the value they can return all of which are parameterized on a per-operation basis (ie **ZErr**, **ZBnd** and **ZRes** are type families):

- An error indicating that the value is not computable.
- A final and precise result.
- A partial result, a convex bound indicating that calling the iteration of the arrow will get us closer to ax' final result. We will have a chance to see the bounds in detail in the section 5.3.3 oncaps and bounds.

```

1 data BndR w
2   = BndErr (ZErr w)
3   | BndRes (ZRes w)
4   | BndBnd (ZBnd w)

```

Listing 5.9.: The definition of the return value of an Antisthenis process. It may be a final result, an error indicating that a final result is non-computable, or a bound for the final value.

Indeed, the iteration of the process has a very different meaning depending on whether the previous returned value is a bound or an error/result. In the former case the iteration continues a previous computation, in the latter there are two options.

- the previously returned value is still valid and is therefore returned,
- if it is not and the computation is reset.

To reflect this fundamental difference between a paused and a finished computation we will call the former simply an *iteration* and the latter *coiteration*.

Finally, we turn our attention to **Conf** w which represents the configuration for running the next step of the computation. Much like **BndR** the configuration is parameterized w.r.t. w but has a basic structure presented in listing 5.10

```

1 data Conf w =
2   Conf { confCap :: Cap (ZCap w)
3         , confEpoch :: ZEpoch w
4         }

```

Listing 5.10.: The type definition of a configuration. It is a tuple containing information that can be used to derive whether a value is valid.

The configuration for running a computation includes an epoch (discussed in section 5.3.2) and a cap (discussed in 5.3.3). It is passed as input to each computation and propagated to the children computation. If the process is iterating (as opposed to co-iterating, i.e. it has not returned a concrete result or error) the cap is compared against the previously yielded bound, if the cap is less than the bound returned then the process must keep processing until it found a bound that exceeds that cap. This way we can control how much work a process is allowed to do before switching to one that might help Antisthenis to cut the overall computation short. This way Antisthenis avoids falling into computational rabbit holes, where it evaluates a long computation when switching to a shorter one would help arrive sooner to a final result.

The *epoch* represents state external to Antisthenis like parameter versions. When the epoch indicates that a parameter on which the process depends has been updated, the process needs to be reset and start computation from scratch (see section 5.3.6 on process resetting).

5.3.2. Epochs and coepochs

As mentioned in the previous section the interplay between epochs and coepochs determine the validity of the value of each Antisthenis process. When the returned value is no longer valid the process needs to be reset. In particular, the epoch mainly indicates the "version" of the computation parameters. This might represent the full external state, but the purpose of the epoch is to be used as an increasing value that can be used by the processes to infer whether their progress is based on out-of-date assumptions. For example the epoch may be a map of a natural number per computation parameter indicating the number of times that parameter has been updated. Or it might simply be

the actual value of each parameter if it is cheap the compare against. Each process then can compare the version or value of each parameter of the epoch to the ones it used to build its internal state. When a discrepancy is found the process knows its state is out of date and must reset.

It is clear that not all parameters are relevant to all processes. Most processes only depend on a subset of parameters, or in Antisthenis' terms, most processes only ever depend on parts of the epoch. This reference to a subset of any epoch is reified by the *coepoch*. Coepochs are the monoidal type parameter of the writer arrow transformer of **ArrProc**. In practice this means that due to the writer semantics, the coepochs of evaluated subprocesses are concatenated to produce the coepoch of the parent process. A special case for this are the commutative operations that involve absorbing elements (\wedge , \vee , and \times), where only the coepoch of the subprocess that yields an absorbing element is returned. In short, the coepoch of a subprocess is included in all parent processes that are not constant with respect to the value of that subprocess.

In principle the function that the **ZipperParams** tag needs to implement with respect to the epoch/coepoch interplay is a function that has a type equivalent to listing 5.11.

```
1 combEpochCoepoch0 :: ZEpoch w → ZCoEpoch w → Bool
```

Listing 5.11.: The type of a naive function checking the validity of a value based on epoch and coepoch.

That will indicate whether the value needs to be pushed down. In practice we can do slightly better at that by enforcing the following law on the `combEpochCoepoch0` function:

$$\begin{aligned} \text{combEpochCoepoch}_0(e, c_0 \diamond c_1) &\Rightarrow \\ \text{combEpochCoepoch}_0(e, c_0) \wedge \text{combEpochCoepoch}_0(e, c_1) \end{aligned}$$

Where $e \in ZEpochw$, $c_0, c_1 \in ZCoEpochw$ and \diamond is the monoidal merging operation on coepochs.

We therefore have the option of filtering the epoch into a smaller subset of itself that is only relevant to the subprocesses that contributed to the generation of the coepoch being checked. Thereby the combination function can optionally return a new epoch rendering the type of the function to be something similar to listing 5.12. It is worth noting that if by this process we determine that a reset is required, a new coepoch is required and

therefore there is no way to constrain the coepoch for the reset. In the example above, if $\text{coepoch0} \triangleright \text{coepoch1}$ do not match the epoch, it is possible that re-evaluation of P_1 will not yield **BndRes** \emptyset and thereby we would need to evaluate P_2 including its coepoch to the aggregate.

```
1 combEpochCoepoch :: ZEpoch w → ZCoEpoch w → Maybe (ZEpoch w)
```

Listing 5.12.: The final function for epoch and coepoch combination.

5.3.3. Antisthenis caps and bounds

In section 5.3.2 regarding epochs and coepochs we described how Antisthenis deals with the validity of values returned by a process. In this section we will focus on how the children processes avoid falling into computational rabbit holes using the interplay of *caps* and *bounds*. In the example discussed in the introduction (5.1), we saw that there are case where, as soon as a subprocess can prove that its value is going to exceed a certain threshold (*cap*), any further work on its part is likely futile. It then returns a *bound* which is propagated up the chain of parent processes and handled at the point where the cap was imposed. The reader is encouraged to maintain a mental model where the cap is an threshold and the bound is a lower bound.

We expect that the bound of a function satisfies the *mototonicity criterion*:

$$C \dot{<} B[f(a_0, \dots)] \Rightarrow C \dot{<} f_B(B[a_0], \dots)$$

Where C is the cap, $B[E]$ is the bound of a process expressed by the expression E by only taking into account the top level terms. f is the function if the operator and f_B is a function that combines the bounds of the arguments to come up with a bound for f . $C \dot{<} b$ means that the bound b exceeds the cap C .

In simpler terms this means that if a bound exceeds the cap, evaluating the arguments and coming up with a tighter bounds will still exceed the cap.

There may be more requirements of the properties of a cap w.r.t. to a particular operator. For example, the minimum operator requires that bounds be translatable to caps and that if a bound b is translated to a cap c and $b' \dot{<} c$ then $b' < b$. This requirement is only related to the way the particular operator is implemented and is not necessarily a requirement for any other operator.

The **ZipperParams** tag needs to provide a way for comparing caps to bounds. Unlike the case of coepoch and epoch combination function this one is much more straightforward:

```
1 exceedsCap :: ZCap w → ZBnd w → Bool
```

Simply from a cap and a bound check if the bound exceeds the cap. A special case for the cap is **ForceResult** for which `exceedsCap ForceResult _ = False`. As the name suggests **ForceResult** is used when we want a process to necessarily return a concrete result or an error, as opposed to a bound.

Now that we have seen how the cap is handled on the input side we should discuss how the parent process passes a cap to its child processes. The cap is created in an operator-specific way (ie via the overloading of a function based on the tag `w`). In particular an operator defines a function that "localizes" the entire configuration.

```
1 data MayReset a = DontReset a | ShouldReset
2 zLocalizeConf :: ZCoEpoch w → Conf w → Zipper w p → MayReset (Conf w)
```

The `zLocalizeConf` function accepts a coepoch, the configuration that the parent process receives and returns the configuration passed. The `zLocalizeConf` function also decides whether the process needs to be reset by calling the `combEpochCoepoch` function that we discussed in subsection 5.3.2 on epochs and coepochs.

To make all this more tangible, a process that adds epoch implements the `zLocalizeConf` roughly like shown in listing 5.13.

```
1 zLocalizeConf coepoch conf z =
2   combEpochCoepoch coepoch (confEpoch conf)
3   $ conf { confCap = newCap
4     }
5   where
6     -- zRes is the partial result so far
7     newCap = case zRes z of
8       SumPartErr _ → error
9         $ "The partial sum is an error this thunk "
10        ++ "should never be reacahble beacuse "
11        ++ "no subprocesses should be called."
12       SumPartInit → confCap conf -- No subprocesses have been evaluated.
```

```

13      SumPart partRes → case confCap conf of
14          -- partRes is the total sum so far. We offset the global cap
15          -- by the sum so far so that the local cap ensures that the
16          -- cursor process, when ran, does not cause the overall
17          -- bound to exceed the global cap. This may generate weird
18          -- caps like negative values. That is ok as it should be
19          -- handled by the evolutionControl function. Note again that
20          -- zRes does not include the value under cursor.
21      CapVal cap → CapVal $ subCap cap partRes
22      ForceResult → ForceResult

```

Listing 5.13.: A sample implementatation of the function that transforms the configuration received by a parent process into one suitable for the child process. Checks if the parent process needs to be reset and uses the partial result to constrain the cap.

Since the process knows it will be evaluating the cursor, it subtracts the partial sum so far from the cap to come up with the cap to pass to the subprocess.

5.3.4. Antisthenis zipper

We saw in the introduction on Antisthenis (5.1) what a zipper is in general. Here we will specialize the notion for the specific case of the internal state of the computation. A single process is defined by an operator, the subprocesses and a partial result for the computation. The process evolves by evaluating a subprocess at a time. We define a zipper structure that

- focuses on the particular process to be evaluated next.
- arranges a recet process that restarts the evaluation from a state that possibly takes into account information about the previous evaluation, for example by starting evaluation from a subprocess that previously yileded an absorbing element.

With this in mind we define the zipper data structure keeps track of the state of the node (see figure 5.2).

- A set of `initial` processes that have not been evaluated yet or that have yielded deprecate values (see section 5.3.2 on epochs and coepochs).

- A set of *iteration* processes, processes whose latest evaluation has yielded a result bound is represented by a data structure that associates the iteration processes with *initial* processes the bounds returned. Since these processes that have yielded a value bound and need to be rerun with a different cap (see section 5.3.3). It should be stressed that this still acts like a heap of subprocesses where the internal properties of the data structure decide the top element that is to be popped. Since we have some information about the final result of these processes, namely the bound, the parent process can be smart about the order in which they are evaluated. As the particular strategy is dependent on the operation implemented by the parent process, the particular data structure used is also dependent on said operation. For example, a process implementing a sum between the subprocesses stores them in a list as there is no beneficial order in which to evaluate the subprocesses while minimum operation benefits from evaluating the processes that have lower bounds first (for more details on the minimum operation see section 5.4.1).

- *Coiteration processes* are processes that have been evaluated and yielded a concrete *e* (error or result), ie a value that is predicated on the epoch (see section 5.3.2) but can not be refined by re-evaluating the process with a more permissive cap. Upon reset the coiteration stack is ordered in an operator specific way and moved to the init stack.

- A *cursor* that is the next process to be triggered along with a way to reset it and the previous value

- A *partial* value to which values can be inserted or removed (for example the sum accumulated up to that point). The partial value does *not* include the value stored in the cursor.

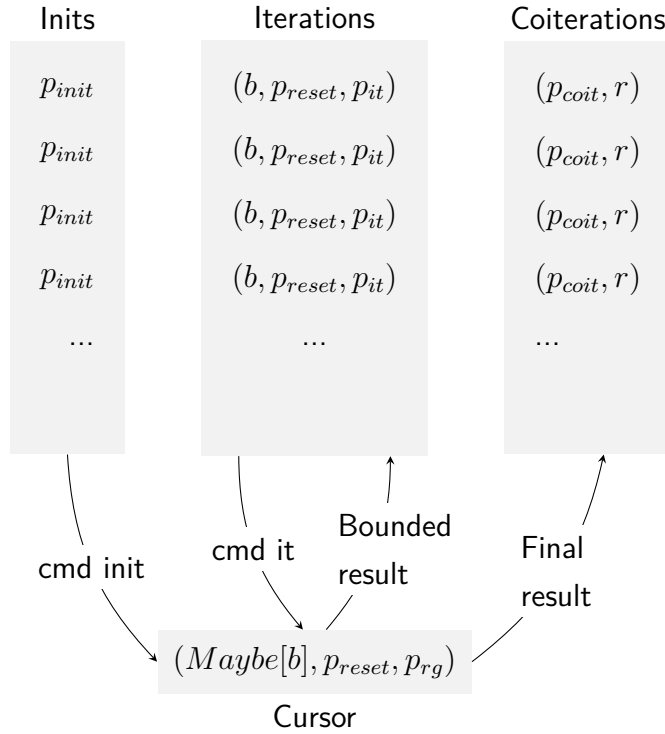


Figure 5.2.: The zipper splits the sub-processes that cooperate to compute values for the owning process in four categories a) initials that have not yielded a currently valid value b) iterations that have not yielded a valid partial value c) coiterations that have yielded a full valid value d) the cursor that is the next sub-process to be evaluated. Depending on the evaluation strategy of the operator once the cursor process evaluates it is replaced with an init process or an iteration process, and depending on the value the cursor process yields it is pushed to the iterations or the coiterations of the zipper.

During computation, the subprocesses owned by the parent process are evaluated, and their iteration are then moved from the **initials** set to the **coiterations** set, possibly with an intermediate stop in the **iterations** set according to the values and strategies. These are further elaborated in the section on section 5.4 on Antisthenis operations. The structure of **iterator** set is highly operator dependent and primarily geared towards efficiently evaluating the parent process, but fundamentally it must not compromise on the correctness of the operator. For example, the if a process is combining elements of an abelian group (commutative, associative, invertible) and the bound has the same type as the result – as it is the case with addition – the structure of the zipper is fairly

straightforward:

- the partial result maintains the aggregation of all values and bounds encountered
- When evaluating the cursor we add the result to the group and remove the bound of the next cursor if we are popping from the iteration set.
- The the partial result is precisely the bound of the final value.

At the other end of the spectrum, a process evaluating a magma (a set with a closed binary operation), where the binary operator has no properties besides closure, must never push anything to the iterations set. The set is equivalent to unit (\emptyset), and the cursor is forced to yield a concrete result or an error by being run with a **ForceResult** cap.

As promised we will expand on these in section 5.4 on Antisthenis operators but to make the concept tangible here are some examples of what the iterator set is defined as for different operations (all numerical operations are assumed to operate on non-negative numbers, in particular query costs):

- Addition is the simplest case: the iteration set is an always-empty container since we always evaluate the cursor until it is sent to the coiteration list. In the context of addition there is no heuristic to predict which subprocess is more beneficial to evaluate.
- For multiplication over positives we want to evaluate the ones that have a zero bound first in the hopes that they will turn out to evaluate to absorbing elements. Therefore, the iteration set is represented as two lists: one where all the zero-bounded iterations are stored and one for all the rest.
- In the case of minimum the iteration set is a heap. We always want to evaluate the one with the minimum bound. When the minimum bound exceeds the minimum concrete result encountered, the min result is the final result.

For completeness we provide the definition of the **zipper** in listing 5.14

```

1  data ZipState w a =
2    ZipState
3    { bgsInits :: [InitProc a]
4      ,bgsIts  :: ZItAssoc w (InitProc a,ItProc a)

```

```

5     ,bgsCoits :: [(Either (ZErr w) (ZRes w),CoitProc a)]
6   }
7
8   data Zipper' w cursf (p :: *) partialRes =
9     Zipper
10    { zBgState :: ZipState w p
11      ,zCursor  :: cursf (Maybe (ZBnd w),InitProc p,p)
12      ,zRes     :: partialRes -- The partial result without the cursor.
13    }
14   type Zipper w p = Zipper' w Identity p (ZPartialRes w)

```

Listing 5.14.: The definition of the zipper.

Finally, each operator defines the partial result type of the zipper `ZPartialRes` and functions for putting and replacing values in it. When the subprocess on the cursor is evaluated the result needs to be either "appended" into the partial result if the new cursor is drawn from the initials set, or to replace the corresponding bound value if it is replaced with a value from the iterations set.

5.3.5. The Cmds functor

As indicated in the section on zippers the process evolution is guided by shifting the zipper focus to different subprocesses and then evaluating them.

A process is internally represented as a mealy arrow (introduced in section 5.2.2) describing the evolution of a `ZProc` which is transformed into the Antistehnis process described in section 5.3.1 by following the strategy in which the `Zipper` is meant to evolve. An operator is allowed to follow a strategy defined specifically its `ZipperParams` tag, but the function for evaluating said strategy applies to the `ZProc` object. There are two important aspects to be noted about `ZProc` (listing 5.15).

- The return value is a zipper full of processes `Zipper w (ArrProc w m)`
- The Kleisli arrow is a free functor of commands `FreeT (Cmds w) m`.

```

1   type ZProc e m =
2     MealyArrow

```

```

3      (WriterArrow (ZCoEpoch w) (Kleisli (FreeT (Cmds w) m)))
4      (LConf w)
5      (Zipper w (ArrProc w m))

```

Listing 5.15.: An internal representation of the process evolving the internal representation of a process: the zipper.

Cmds is a union type of different directions toward which the zipper can be evolved. Its definition is provided in listing 5.16. In this general an operator agnostic form the evolution can take one of a few state transitions at a time:

- It can always be reset
- It can pop a process from the initial to the cursor to be executed
- It can pop a process from the iterator set if it is not empty. Which one is to be popped is decided by the parameters
- The may have reached a final value where it can't be evolved anymore.

```

1  data Cmds' r f a =
2    Cmds { cmdReset :: ResetCmd a
3           ,cmdItCoit :: ItInit r f a
4           }
5
6  data ItInit r f a
7    = CmdItInit (ItProcF f a) (InitProc a)
8    | CmdIt (ItProcF f a)
9    | CmdInit (InitProc a)
10   | CmdFinished r -- when a process finishes it should stick to a
11                     -- value until the epoch/coepoch pair requests a
12                     -- reset.

```

Listing 5.16.: Definition of the commands functor that provides different branches of evolution for zipper.

So the **Cmds** facilitates a tree of actions in conjunction with a free monad . A free monad of functor **f** is actually many **f** nested like **f** (f (f ... f (f a))). This that **FreeT**

Cmds _m is an arbitrary depth tree of **Cmds**. This tree is navigated in an operator specific way guided by the generic function `evolutionStrategy` the type of which is presented in listing 5.17.

```
1 evolutionStrategy
2   :: forall x .
3   FreeT (Cmds w) m x
4   → m (Maybe (RstCmd w m x), Either (ZCoEpoch w, BndR w) x)
```

Listing 5.17.: A function that each antisthenis operator needs to implement that decides the traversal of the tree created by the different possible evolutions of zipper. The implementation may also deem that a good reset point has been discovered.

This function navigates the tree in the way that is most beneficial to the particular operator and comes up with two values:

- A reset command to be used if reset is required in the future (this will be discussed in the section about 5.3.6 resetting)
- When encountering a **CmdFinished** value return the result of the finished operation, otherwise traverse the tree to the preferable leaf and return the leaf.

An abridged example is presented in figure 5.3.

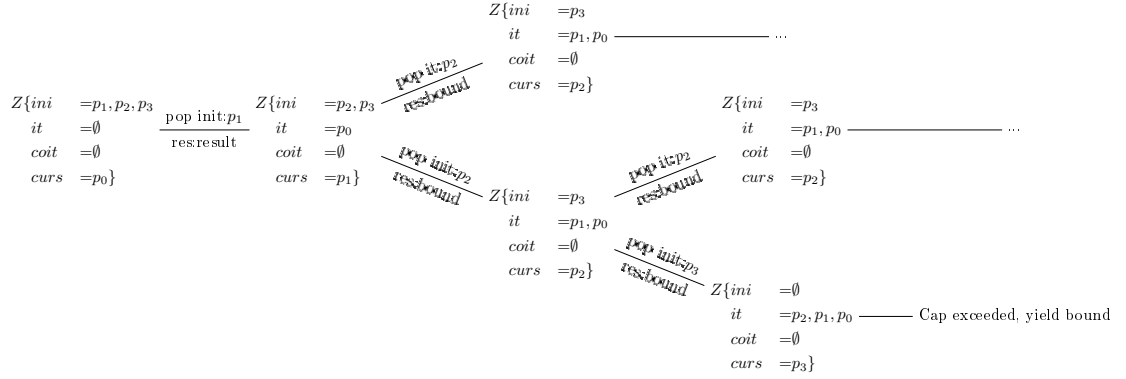


Figure 5.3.: A tree of different evolution paths a zipper may undergo during between the calling of a process and the accumulated bound exceeding the cap. p_0 , p_1 , p_2 , p_3 and p_4 are subprocesses. The values returned, the specific iteration subset data structure are not presented for brevity and to shift focus to the movement of the subprocesses within the zipper. Nodes in this tree that have two children correspond to the **ItInit** constructor, nodes with one child correspond to the **CmdIt** or **CmdInit** constructor depending on which process set the process in the next cursor is drawn from. When a concrete value is reached or the cap is exceeded the process stops and the free monad that wraps **Cmnds** takes the value **Pure**.

5.3.6. Resetting processes

In subsection 5.3.5 discussing the command functor we discussed that it may include a reset process. Why create a reset process when reusing the initial process would reset the computation just as well? In every case a reset moves all the init parts of the elements of iteration set as well as the coiterations into the init set of the zipper. But in what order? From previous runs it is often apparent that some execution strategies are more likely to lead to a quick result than others. For example, under the assumption that processes change their values less often than they maintain them, while computing a logical AND operation (\wedge) we likely want to start the computation from processes that previously yielded **False**.

To accommodate this, the **Cmnds** functor (introduced in section 5.3.5) can yield a reset point. While running a process we maintain in the lexical scope of the mealy arrows computation the latest reset process to switch to when the epoch/coepoch comparison requires as much. When generating a **Cmnds** functor, the operator has the option to provide a reset process that resets a running zipper that is at an arbitrary state.

5.4. General operators

While describing the core concepts of Antisthenis we often dealt with functionality that is operator specific. Indeed, a large part of Antisthenis depends on operator specific logic. We reify operators via uninhabited types that instantiate the `w` parameter of many types we mentioned so far (`BndR w`, `Zipper w`, `ArrProc w m`, etc). We refer to these types as “tags”. We have implemented a few tags in Antisthenis to cover the needs of FluidB (see listing 5.19).

```

1  -- Result types
2  class BndRParams w where
3    -- The result type when the value is non-computable
4    type ZErr w :: *
5    -- The partial result.
6    type ZBnd w :: *
7    -- The full result
8    type ZRes w :: *
9
10 -- Zipper internal types
11 class BndRParams w => ZipperParams w where
12   -- The type of the process cap
13   type ZCap w :: *
14   -- The type of the process epoch
15   type ZEpoch w :: *
16   -- The type of the coepoch
17   type ZCoEpoch w :: *
18   -- The type of the iteration set
19   type ZItAssoc w :: * -> *
20   -- The type of the partial result
21   type ZPartialRes w :: *
```

Listing 5.18.: Operator specific types that need to be implemented by every operator.

```

1  -- The witness type for the sum operator
2  data SumTag p
```

```

3  -- The witness type for the minimum operator
4  data MinTag p
5  -- The witness type for multiplication
6  data MulTag p
7  -- The witness type for logical And/Or
8  data BoolTag op

```

Listing 5.19.: Tags are phantom types that define the Antisthenis operations. The tags themselves may be parameterized using type parameter `p`. For the needs of FluidB we define operations for addition, subtraction, multiplication and boolean operations.

The parametric part of the tags, especially for `SumTag` and `MinTag`, is due to the fact that we want to be able to further parameterize them for concrete cost calculation (using sum and min operators normally) and for stochastic cost calculation (see section on 5.5.4 historical cost). These two both utilize minimums and sums, but they use different types for return values as well as caps/bounds and epochs/coepochs.

Since all types and implementations can be derived from the `ZipperParams` type `w`, the Antisthenis process creation function can be as simple as in listing 5.20, where `mkProc` combines a list of subprocesses into a parent process that evaluates the operator described by the tag `w` without requiring any further information.

```

1  mkProc
2  :: (Monad m, ZipperParams w, Eq (ZCoEpoch w)) => [ArrProc w m] -> ArrProc w m

```

Listing 5.20.: The type `w` fully defines the operator so combining subprocesses into a process is unambiguous.

Compatibility between heterogeneous processes

The astute reader has likely noticed that `w` is the parameter of both the subprocesses and the process even though they may be different operators for example the operands of a min are themselves sum expressions in FluidB cost estimation). The reason is that a process expects to see a specific type of coepoch and bound returned by its subprocesses, while the subprocess must be able to deal with the epoch and cap that the parent process is defined for. To do that we avoid defining an $N \times N$ correspondence between all the operators, but rather we define ad hoc a record that facilitates the conversion, the type

of which is presented in listing 5.21. This way we can compartmentalize the translation among operators from the internal logic of the operators.

```

1  data Conv w w' =
2    Conv
3    { convEpoch :: ZEpoch w' → ZEpoch w
4      , convCoEpoch :: ZCoEpoch w → ZCoEpoch w'
5      , convCap :: ZCap w' → ZCap w
6      , convRes :: ZRes w → ZRes w'
7      , convBnd :: ZBnd w → ZBnd w'
8      , convErr :: ZErr w → ZErr w'
9    }
10 convArrProc :: Monad m ⇒ Conv w w' → ArrProc w m → ArrProc w' m

```

Listing 5.21.: An object of type `Conv w w'` can act as an interface between parent processes of type `w` to subprocesses of type `w'`.

5.4.1. Basic cost estimation operators

General cost estimation involves minimum and sum operators with some special semantics:

- Non-computable values are larger than any value from the perspective of minimum.
- All values are positive (therefore all values are implicitly bounded by $[0 - \epsilon]$)

Sum

We do not implement a general sum operator but rather a sum operator that assumes the properties of query costs as the values it deals with. The sum is probably the simplest operator implemented. The adding zipper evaluates each initial subprocess under a cap $c_{subproc}$ that is

$$c_{subproc} := c_{parent} - r_{partial}$$

where $r_{partial}$ is the partial result omitting the cursor, ie the current subprocess and c_{parent} is the cap that the parent process is operating under. If the parent cap is

ForceResult all arguments are also evaluated under **ForceResult**. Of course if $c_{parent} > r_{partial} + b_{cursor}$ then $r_{partial} + b_{curopr}$ is returned as the result bound.

Minimum

Much like the sum operator, the minimum operator implementation also assumes that it is manipulating cost values, and it is used in conjunction with the sum operator to define our cost evaluations. The operation of the minimum Antisthenis operator depends heavily upon the iteration set which has the interface of a heap. The zipper evolution strategy is described by the pseudo-python code in listing 5.22

```

1  def exceeds_cap(cap, zipper):
2      return zipper.iter_procs.min().bound > conf.cap \
3          and zipper.coiter_procs.min().result > conf.cap \
4
5
6  def minimum_zipper_evolution(conf, zipper):
7      # Run each initial process with a zero bound
8      while len(zipper.initial_procs) > 0:
9          zipper.iter_procs.push(cursor)
10         zipper.cursor = zipper.initial_procs.pop()
11         res, nxt_proc = zipper.cursor.run(conf{cap = 0})
12         # We cant hope to get a cost lower than zero
13         if res == BndRes(0):
14             # Until the epoch is updated
15             return BndRes(0)
16
17         # Push it to the correct iter set.
18         if res.is_bound():
19             zipper.iter_procs.push(res, nxt_proc)
20         else:
21             zipper.coiter_procs.push(res, nxt_proc)
22
23     # Run the iterating processes
24     while True:

```

```

25     # If the result is better than the best iteration, it's the final result
26     if zipper.coiter_procs.min().result < zipper.iter_procs.min().bound:
27         return BndRes(zipper.coiter_procs.min().result)
28
29     # If there are two iter processes
30     if not exceeds_cap(conf.cap, zipper):
31         # Run the best iterating process until it exceeds the
32         # second best iterating process or the final result
33         cap = min(zipper.iter_procs.second_min().bound,
34                 zipper.coiter_procs.min().result)
35         res, nxt = zipper.iter_procs.min().run(conf{cap = cap})
36         # Push it to the correct iter set.
37         if res.is_bound():
38             zipper.iter_procs.push(res, nxt_proc)
39         else:
40             zipper.coiter_procs.push(res, nxt_proc)
41     else:
42         conf = yield zipper.iter_procs.second_min().bound

```

Listing 5.22.: Pseudocode for the algorithm of evaluating a process up to a threshold defined by the cap. For brevity and to avoid too much unnecessary detail we omit sanity checks and the reset handling code.

Boolean conjunctions and disjunctions

Antisthenis implements operators for boolean conjunction and disjunction to answer the question of whether a node is materializable in the presence of a specific inventory of materialized nodes. Booleans are different to the other operations in how they handle caps. Boolean \wedge and \vee do not use the same type for caps and bounds as for results. Instead the bound is a two dimensional value that represents the minimum number of steps required to arrive to each possible value (**True** or **False**).

To efficiently evaluate boolean algebra expressions we make use of absorbing elements,. Consider the evaluation of a boolean expression

$$Q := A \wedge B \wedge C$$

where A , B and C are subexpressions. Obviously if any of the expressions A, B or C evaluates to **False** the entire expression would consequently evaluate to **False** as **False** is the absorbing element for the group of booleans over \wedge . Furthermore without taking into account the structure of each of the subexpressions A , B and C we can tell that the amount of work for evaluating it depends on the result. If the result of this expression is **False** then the best case scenario in terms of amount of work we would need to do is for A to be the expression **False**. Or more succinctly

$$w(Q|Q = False) \geq \min(w(A|A = False), w(B|B = False), w(C|C = False))$$

Where $w(E|E = k)$ is the amount of work for calculating E if E turns out to be k . In this case we would be able to evaluate the expression with just one dereference. If the result is **True** then the best case scenario is for A , B and C to all be **True** and we would be able to come to a result with three dereference steps. Or equivalently

$$w(Q|Q = True) = w(A|A = True) + w(B|B = True) + w(C|C = True)$$

Now let's consider the compound expression

$$X := (A_1 \vee A_2 \vee A_3) \wedge (B_1 \vee B_2) \wedge (C_1 \vee C_2 \vee C_3)$$

Where A_i , B_i and C_i are expressions. How would we navigate this expression to minimize the number of operations? We could completely expand each term and hope that we are lucky, and we encounter absorbing elements early on. As explained in the introduction, especially for terms automatically generated this strategy is unlikely to be efficient.

Instead, at each expansion we track the minimum number of steps required for each value. The best case scenario if $X = False$ is $B_1 = False$ and $B_2 = False$, in which case the expression is evaluated in a minimum of 4 steps (dereferences):

- Evaluate $B_1 = False$ in a single step
- Evaluate $B_2 = False$ in a single step
- Evaluate the subexpression $False \vee False$

- Evaluate $X = (...) \wedge False \wedge (...) = False$

For a $X = (True)$ we need 7 steps:

- Evaluate $A_i = True$ in a single step for any i
- Evaluate the subexpression $True \vee \dots$
- Evaluate $B_i = True$ in a single step for any i
- Evaluate the subexpression $True \vee \dots$
- Evaluate $C_i = True$ in a single step for any i
- Evaluate the subexpression $True \vee \dots$
- Evaluate $X = True \wedge True \wedge True$

A boolean expression in Antisthenis, then, takes the bound value of the minimum cost for trues and falses. If this expression were evaluated with a cap of $\{True : 6, False : \infty\}$ it would return a bound of value $\{True : 7, False : 4\}$. In plain english this means "try to evaluate this if your other options are more expensive than 7 steps for **True** and 4 for **False**". It could be seen as a metric for the size of the tree.

To demonstrate the derivation of caps we present the following expression:

$$X = X_1 \wedge X_2 \wedge X_3$$

Any of the X_i terms has a minimum bound of $\{True : 1, False : 1\}$. Since the absorbing element for \wedge is **False** expanding any of the X_i terms is bounded by $\{True : \infty, False : 1\}$ so an expansion of

$$X_1 = A_1 \vee A_2$$

would cause the subprocess to immediately yield $\{True : 2, False : 3\}$ passing control to the parent process. Thus, a symmetrically growing tree would be traversed in a breadth first fashion and any asymmetries would cause Antisthenis to try to greedily exploit the smaller parts of the tree.

5.5. Node machines

??

A referable machine is a process with a name that can be referenced by more than one other processes. In FluiDB a machine name corresponds to a graph node since the Antisthenis infrastructure is utilized to provide incremental computation of properties on nodes (materializability, cost, stochastic cost, etc). In this section we look into the the layer of Antisthenis that takes into account the dynamics of referable machines.

5.5.1. Antisthenis machine tape

The *machine tape* is a dynamic directory that maps nodes to machines. The meaning of a node is dependent on the implementation, in the context of FluiDB nodes in this context correspond to the FluiDB internal graph nodes. The machine tape is part of the computation context and is updated by new and temporary values of machines (see section [5.5.3](#) on cyclical machines).

Not all processes involved in a computation are involved in the machine tape. Each machine in the tape is composed by a hierarchy of processes that are inaccessible from outside that machine. In fact a machine is discriminated from a process solely by virtue of the fact that it can be referenced via the tape.

Since machines are referable by different other processes it means that, we need to be careful that more than one process do not own their own diverging copy of the same machine, essentially treating it as a process. To avoid this, processes referring to machines do so indirectly via a *machine wrapper process*. A machine wrapper is a process that always does the following:

- look up a name in the tape
- temporarily replace it with a cyclical machine (see section [5.5.3](#))
- evaluate the machine that was just removed from the tape.

The machine tape state is managed via a state monad in the underlying Kleisli arrow. Thereby a recursive reference is formed between the tape and the machine.

5.5.2. Process stack

Before we go further with the details of how exactly named machines are handled, it is important to make the distinction between a process and a computation. In the context of Antisthenis a process is the mealy machine that evolves during evaluation and survives between evaluations. On the other hand a computation is the part of the evaluation itself that refers to the particular machine. In a more abstract sense, a computation is the *relation* between the process and the value produced, regardless of whether that value is an actual result, a bound for a potential result or a witness to the non-computability of the result.

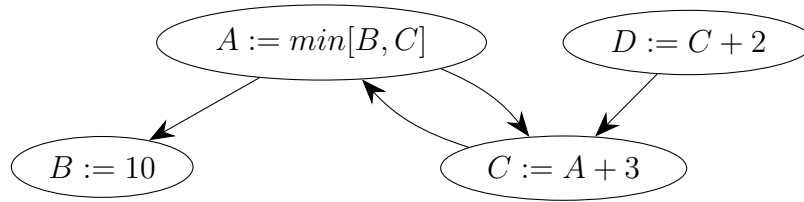
With this distinction in mind we introduce the notion of a process stack which is similar to the call stack in most programming languages. The process stack contains all processes that participate in the computation at any given moment. When a process is called, it is pushed to the process stack and when a result is returned it is popped from the process stack. The process stack is completely ephemeral, acyclic and should be completely empty between computations. In this sense it is part of the *computation context*, the set of effects that affect each computation separately. The computation context may be implicit (like in the case of the process stack as we will see when discussing cyclical machines in section 5.5.3) or explicit in the process' underlying Kleisli arrow.

5.5.3. Cyclical machines

By virtue of machines being referrable by arbitrary other machines, it is possible, and certain in the case of FluiDB, that the Antisthenis machines will form referential cycles. Such cycles in most computational frameworks, especially ones that deploy a non-lazy order of computation, will cause the computation to either fail or to not-terminate. Antisthenis is explicitly designed around allowing the computation to handle such self-referential cases.

We already discussed in a previous section that machines are stored in data structure called a tape that allows them to be referenced by other machines, and that in the period from the time a value is requested of a machine and the time a value is returned (including an error or a value bound) the machine is said to be in the computational stack (introduced in section 5.5.2). During its residence in the stack a machine is in a state where it is unable to produce a value. A cycle occurs when a machine is referenced while in the computational stack.

We make the problem concrete with an example based on the graph in figure 5.4a. A naive approach of dealing with recursive calls is demonstrated in figure 5.4c where naive-Antisthenis looks for cycles by throwing errors (\perp) when encountering a node that is already in the computational stack. Evaluating node A Antisthenis first node B to 10. It then tries to evaluate node C which redirects back to A . C can therefore not yield a value and returns an error \perp . Then, without changing any of the parameters the user tries to evaluate D . Looking at the graph it is clear that D should in sequence evaluate to 15, as demonstrated by the evaluation process in figure 5.4d (the reader is reminded that all values refer to query evaluation costs, and therefore we have asserted that $\min[a, \perp] \equiv a$).



(a) A recursive cost graph.

$$\left. \begin{array}{l} \text{comp}[A]: \\ \text{comp}[B] \Rightarrow 10 \\ \text{comp}[C]: \\ \text{comp}[A] \Rightarrow \perp \end{array} \right] \Rightarrow \perp$$

(b) Computation for A

$$\left. \begin{array}{l} \text{comp}[D]: \\ \text{comp}[C] \Rightarrow \perp \end{array} \right] \Rightarrow \perp$$

$$\left. \begin{array}{l} \text{comp}[D]: \\ \text{comp}[C]: \\ \text{comp}[A]: \\ \text{comp}[B] \Rightarrow 10 \\ \text{comp}[C] \Rightarrow \perp \end{array} \right] \Rightarrow 10$$

(c) Naive computation for D

(d) Computation for D

Figure 5.4.: Node C is recomputed because the \perp value is predicated on A being in the computation stack.

The root of the problem in the naive case is that the process stack is an ephemeral

property and yet it affects the return values that are non-ephemeral. The only solution for this is to lift information about the cycles to the non-ephemeral plane that survives the computation. Indeed, what is needed is for values that are deemed non-computable to *reset* when accessed with a different stack. In other words the value of \perp for C is *predicated* on the process stack containing A , or A being otherwise non-computable. Coepochs fill the bill precisely. We extend the coepoch to include, along with the parameter versions used to derive a particular result, also *a set of machines that need to be non-computable in order for the value to be valid*.

To accomplish this when a machine is looked up, before being evaluated it is removed from the tape and in its place a *cyclical machine* is placed. This machine always returns the same \perp result that indicates that the value is non-computable and sets to the coepoch that it is itself not computable. When the original machine comes up with a value it replaces the cyclical machine with the iteration or coiteration of the process and the coepoch is censored to remove the predicate that the current machine must be \perp . This way whenever any machine tries to evaluate another non-computable machine it finds a cyclical machine which handles the situation automatically. Upon evaluating the cyclical machine, the coepoch of the caller is infected with the predicate that the cyclical machine must be \perp . Pseudo-python algorithm in listing 5.23.

```

1  def cycle_machine(node):
2      # Return a process that returns non-comp and a predicate notifying
3      # that it is itself not computable.
4      class CycleProc:
5          def run(conf):
6              return (lookup_node(node), non_comp, Coepoch(predicates=[node]))
7
8      return CycleProc()
9
10 def run_submachine(node, conf):
11     # Replace the real process with a dummy cycle process.
12     p = lookup_node(node)
13     insert_node(node, cycle_machine(node))
14     nxt, res, coepoch = p.run(conf)
15     insert_node(node, nxt)
16     # Censor the current node from the predicates. If the result is

```



```

17     # computable then the predicate is already falsified, if not it is
18     # confirmed.
19     return (nxt,res,coepoch{predicates=coepoch.predicates.delete(node)})

```

Listing 5.23.: The algorithm for creating predicates.

At the end of the computation Antisthenis goes through all the non-computable predicates and asserts that they are indeed non-computable. The ones that turn out not to be non-computable are collected and set as *copredicates* in the epoch. The node is then reevaluated with the new epoch. Due to the incremental nature of Antisthenis only the parts of the graph that are dependent on the predicates are reevaluated. The process is described in listing 5.24

```

1  def safe_run(n,conf):
2      # To check if new copredicates were added
3      initial_copred_num = len(conf.copredicates)
4      # Run the process
5      coepoch,res,coit_p = lookup_node(n).run(conf)
6      # Remember: machines censor their own reference from the predicate
7      # and no predicates.
8      assert n not in coepoch.predicates
9      assert len(intersection(conf.epoch.predicates,coepoch.predicates)) == 0
10     # For efficiency any result we find we register all results we
11     # find that may be encountered later as predicates.
12     conf.copredicates.push(n)
13     # For each non computable predicate
14     for noncomp_node in coepoch.predicates:
15         # Run the predicate to check if indeed it is non-computable
16         res = safe_run(lookup_node(noncomp_node), conf)
17         if is_computable(res):
18             # If it is actually computable then register it a copredicate.
19             conf.copredicates.push(noncomp_node)
20
21     if len(conf.copredicates) > initial_copred_num:
22         return lookup_node(n).run(conf)

```

```

23     else:
24         return res

```

Listing 5.24.: Algorithm for handling the predicates.

Another way to look at predicates is as deferred parts of the coepoch, to get the full coepoch we also need the coepochs of the predicates. The copredicates then are a forcing mechanism where Antisthenis assumes that some of deferred parts are actually invalid. Antisthenis is lazy with these coepochs first and foremost because it depends on it to handle cycles but also because coepochs may be inhibited by absorbing elements, and it may not need to evaluate them at all. But if a predicate survives both the process stack based censoring that we just described and being inhibited by the operators, Antisthenis has another trick up his sleeve: because it has a completely relativistic view of values, where they are only valid for particular epochs, the only requirement at every moment is that the final result of the computation be correct. Furthermore, the value depends on the machines in the predicate being \perp . Therefore, it is enough to check that they are indeed \perp certifying the predicates. The predicates that are falsified in this way are the ones for which there is no other option but to force their value be taken into account via the epoch's copredicates.

5.5.4. Calculating historical query costs

As we saw in section 4.2 on the query planner, the planner depends on the cost of past nodes to make decisions about the overall usefulness of a plan. A naive approach to this would be to just use the normal min/sum operation set to calculate historical costs. As explained in more detail in that chapter this approach is problematic. We need a way of convincing Antisthenis to "look behind" those materialized nodes a) without completely ignoring the materialized nodes and b) by handling referential cycles more gracefully than declaring nodes non-computable. For this reason we develop the *stochastic costs model*.

The correct way of calculating the expected cost of a node in the future, which we very informally and heuristically attempt to approximate, would be to instead of considering the cost of materialized nodes to be zero, to calculate the likelihood that a materialized node will still be materialized when we encounter it again. This is related not only to an estimation of how many page-writes separate the moment of cost estimation and the actual plan the cost of which is being estimated, but also all the decisions that the garbage collector will make in that time. After FluidB's budget is exhausted for the first

time, for every page write there a page needs to be garbage collected. We considered a couple of options for stochastically modelling the behavior of the GC like assuming that it chooses random pages or random tables, but we could find none that was useful and computationally viable.

For this reason we decided to follow a pragmatic approach to the problem and simply assume that for every materialized node there is a constant probability that it will not still be materialized when we need it, scaling the cost by that factor to get the stochastic cost. Furthermore, under this regime, Antisthenis would encounter many more cycles, rendering virtually all machines non-computable. For that reason we make all values semi computable by introducing a $[0, 1]$ scale where 0 would mean that the value is fully computable and 1 would mean that the value is completely non-computable.

Fortunately this change does not require a complete overhaul of the operator implementations, only the definition of more complex types for cap, bound and result values. Particularly, we qualify each value with two numbers:

- the materialized nodes encountered
- a quantification of the non-computability of the value

Some care needs to be taken so that the values remain independent of the root of the computation.

The cap (listing 5.25) is only qualified by the number of machines corresponding to materialized nodes that are currently in the computation stack, i.e. machines that would have otherwise yielded zero cost.

```

1 data HistCap a =
2   HistCap
3   { hcMatsEncountered :: Int
4     , hcValCap :: Min a
5     , hcNonCompTolerance :: Double
6   }
```

Listing 5.25.: Definition of the type used for capping the cost of historical queries. The bounds and results are augmented by two new dimensions (listing 5.26):

- The $[0, 1]$ computability metric

- The maximum chain of materialized nodes encountered.

And the comparison between bounds and caps is presented in listing 5.27.

```

1 data HistRes a =
2   HistRes
3   { hrComp :: Double -- [0,1] computability metric.
4     ,hrMaxMatTrail :: Int -- Maximum trail of materialized nodes
5                           -- encountered.
6     ,hrRes :: a
7   }

```

Listing 5.26.: Definition of the type used for both bounds and results of historical queries.

```

1 -- | We prefer "more computable" results than "less computable" so we
2 -- scale the actual result by the computability. We also define a
3 -- computability threshold under which we consider the value always
4 -- >.
5 compare :: HistRes a → HistRes a → Comparison
6 compare (HistRes c1 t1 r1) (HistRes c2 t2 r2) =
7   if
8     | c1 > thr && c2 > thr → compare (scale c1 r1) (scale c2 r2)
9     | c1 > thr           → GT
10    | c2 > thr           → LT
11    | otherwise          → compare a b
12   where
13     thr = 0.7
14
15 -- There are no circumstances under which we are going to choose a
16 -- value that has come through more materialized nodes. The scaling
17 -- will have damped it enough.
18 instance Ord a ⇒ Ord (Cert a) where
19   compare (Cert _ a') (Cert _ b') = compare a' b'
20
21 extExceedsCap _ HistCap {..} HistRes {..} =

```

```

22     maybe False (hrRes >) (unMin hcValCap)
23     || hrMaxMatTrail bnd > (maxMatTrail - hcMatsEncountered)
24     || hrComp bnd > hcNonCompTolerance
25 where
26     maxMatTrail = 5 -- actually this is part of a global config.

```

Listing 5.27.: Comparison between bounds and between bound and cap are different. Between bounds we need to account for the computability metric. The cap on the other hand defines a three-dimensional bound () that the bound must fall within in order to not exceed it.

5.6. Conclusion

In this chapter we described in detail the framework for building incremental evaluation systems we developed for FluiDB that we call Antisthenis as well as a brief summary of the way we use it in FluiDB. The main goals of Antisthenis are to compute incrementally, to gracefully handle computational cycles and to order the operations in each operator in a specific way such that the overall computation finishes as soon as possible.

There are two major shortcomings of Antisthenis:

- Extending it to new operators and even using the existing operators involves a lot of boilerplate. This is an immediate effect of the current version of Antisthenis being highly specific to the FluiDB planer.
- Our current implementation relies heavily on the higher order functional capabilities of haskell that are efficient for what they are, but they stress the garbage collector and are hard for GHC to optimize.

Both of these could be mitigated in a future version of Antisthenis that would be detached from FluiDB and that would replace the Antisthenis process with a more "traditional" data structure.

Further work could also be done to parallelize Antisthenis. Some operations like `min` rely on the operands being evaluated serially to take advantage of early stopping but addition is much more flexible with the order of evaluation. It is clear that parallelism would also be implemented in an operator specific way.

Execution engine

Don't hit hard the beads, work
makes a man.

(Vamvakaris)

Chapter summary

- Tables are stored in a memory file systems and records are stored as POD binary objects organized in 4K pages.
- The extent of each relation involved in a plan is converted to C++ struct and each predicate is converted to a C++ callable class.
- These parameterize the templated operators to allow the compiler to generate highly specialized code.
- In this chapter we discuss the implementation important operations and their reverse.

The logical plan is translated into a physical plan that has the form of a c++ program. The motivation for using code generation relates to expression level optimizations so it's less work on some fronts (although more work in others). Also gdb for debugging. Remember that compilation time is not pure overhead, we would have to duplicate at least some low level optimizations in the haskell.

6.1. Codegen introduction

Code generation is becoming a more and more common in RDBMSs. Used mainly in in-memory databases, where disk IO does not dominate the runtime, it aims to minimize the overhead of data access function calls, to optimize the operdicates and numerical expressions and to avoid indirection in tight loops. Approaches to code generation fall generally on a spectrum between two extremes:

- Transpilation of the physical plan to a low level programming language like C or C++ for every query
- JIT compilation of small parts of the plan as the query executes.

FluidDB follows the approach of Krikellas et. al. [29] falling far to the former end of the spectrum. We generate very specific, template-heavy C++ code for every query and we call to an off-the-shelf compiler to generate highly optimized machine code.

6.2. Storage memory management

We opted for delegating the task to the OS and use use use the tmpfs as a storage layer to our database. The tmpfs filesystem depends on the the *shmem* module (as of linux v5.13) for handling file operations. *shmem* is a resizable virtual memory filesystem for linux. Where a typical persistent filesystem stores files in a block device and caches pages in memory for efficiency, *shmem* keeps files exclusively as pages in the page cache. The OS tries to keep all pages in memory but when resources start running out it writes pages in swap.

Assuming that the pages are not in swap, normal reads for *shmemfs* using `read()` are equivalent to copying pages from the page cache to the user space. *shmem* writes on the

other hand operate directly to the pages. We can mitigate the copying overhead of reads using `mmap` which will remap the page to the address space of the application.

The problem with this approach is that, while it allows us to minimize copying, we still need to run system calls in tight loops, which can be very computationally expensive. This can't be completely mitigated unless we move the "storage" layer to the userspace, re-implementing the memory management that we get for "free", in terms of engineering effort, from the `shmem` module so we leave that for a future version of FluidB.

On the other hand this approach allows FluidB to easily be adapted to operate over any filesystem backed by different storage technologies like non-volatile memories.

6.3. Data layout

Before we get into the details of the actual physical plan (in the form of C++ code) we need to state some assumptions about the layout of the code that is assumed by the primary data.

FluidB is a *row store system* and it depends on the filesystem for intermediate result lookup and page management. For code generation to make sense the file system is, in particular, a `tmpfs` filesystem that resides entirely in ram. This is not an ideal solution performance-wise and in the future we plan on using a less OS-reliant way of managing storage resources but it is good enough for now. In particular we use one file per table or intermediate result, and each file is simply a sequence of records organized into pages. For now FluidB does not support any kind of indexing or compression. FluidB can, of course, be run over any filesystem but non-memory based file systems diminish benefit of code generation, making performance IO bound.

With this in mind there are three parts to understanding the principles of FluidB storage:

- The format in which primary data is inserted into the database.
- The layout of the data within the database
- The transformation from the former to the latter.

6.3.1. Initial data conversion

Starting with the initial data, as it stands, for FluidDB to be adapted to a particular dataset it requires the primary data in CSV format and some haskell code describing the shape of the data and the database configuration. Our experiments so far have been revolving around the SSB TPC-H benchmark so the format expected is the plaintext format that dbgen [77] generates. This is comprized by two steps: first a haskell program parses the CSV records into standard-layout binary objects that can be directly cast to C/C++ standard-layout structs. These binary objects are stored one after the other in a flat binary file. The end result of what we want to do to be able to execute code similar to the one presented in listing 6.1.

```

1  template<typename R, size_t batch_size=2000>
2  void bama_to_dat(const std::string& bama_file, const std::string& dat_file) {
3      int fd;
4      size_t read_bytes;
5      std::array<R, batch_size> batch;
6      // Open the file in binary mode
7      fd = ::open(bama_file.c_str(), O_RDONLY);
8      // The writer is the object used by all the libraries and it will
9      // write each record in pages.
10     Writer<R> w(dat_file);
11     do {
12         // Read a batch of data,
13         read_bytes = ::read(fd, batch.data(), sizeof(batch));
14         // Use the writer object to write each record.
15         for (size_t i = 0; i < read_bytes / sizeof(R); i++) {
16             w.write(batch[i]);
17         }
18         // Keep reading until a batch is cut short.
19     } while (read_bytes == sizeof(batch));
20     // Wrap up.
21     w.close();
22     close(fd);
23 }
```

Listing 6.1.: For standard FFI communication C++ structs that do not contain fancy constructors

For this to work we refer to the C++ standard [68]. The types that represent table rows must be standard layout. According to the standard:

A class S is a standard-layout class if it:

- has no non-static data members of type non-standard-layout class (or array of such types) or reference,
- has no virtual functions and no virtual base classes
- has the same access control for all non-static data members,
- has no non-standard-layout base classes,
- has at most one base class subobject of any given type,
- has all non-static data members and bit-fields in the class and its base classes first declared in the same class, and
- has no element of the set $M(S)$ of types as a base class, where for any type X , $M(X)$ is defined as follows. [Note: $M(X)$ is the set of the types of all non-base-class subobjects that may be at a zero offset in X .]
 - If X is a non-union class type with no (possibly inherited) non-static data members, the set $M(X)$ is empty.
 - If X is a non-union class type with a non-static data member of type X_0 that is either of zero size or is the first non-static data member of X (where said member may be an anonymous union), the set $M(X)$ consists of X_0 and the elements of $M(X_0)$.
 - If X is a union type, the set $M(X)$ is the union of all $M(U_i)$ and the set containing all U_i , where each U_i is the type of the i th non-static data member of X .
 - If X is an array type with element type X_e , the set $M(X)$ consists of X_e and the elements of $M(X_e)$.
 - If X is a non-class, non-array type, the set $M(X)$ is empty.

The records we generate certainly conform to these requirements. An example is the supplier row in 6.2. Therefore the object is trivially copyable and occupies contiguous

bytes of storage. This means that we can safely write each record `R` as `sizeof(R)` contiguous binary data to a file and expect to find the same value of `R` when we read it. Based on this we can safely copy binary record objects from memory to disc and visa versa.

```

1  class Record {
2  public:
3      Record(unsigned __s__supkey, fluidb_string<25> __s__name,
4              fluidb_string<40> __s__address, fluidb_string<16> __s__city,
5              fluidb_string<16> __s__nation, fluidb_string<13> __s__region,
6              fluidb_string<15> __s__phone)
7          : s__supkey(__s__supkey),
8            s__name(__s__name),
9            s__address(__s__address),
10           s__city(__s__city),
11           s__nation(__s__nation),
12           s__region(__s__region),
13           s__phone(__s__phone) {}
14
15     Record() {}
16
17     std::string show() const {
18         std::stringstream o;
19         o << s__supkey << " | " << arrToString(s__name) << " | "
20           << arrToString(s__address) << " | " << arrToString(s__city) << " | "
21           << arrToString(s__nation) << " | " << arrToString(s__region) << " | "
22           << arrToString(s__phone);
23         return o.str();
24     }
25
26     bool operator==(const Record& otherRec) const {
27         // compare each field ...
28     }
29
30     bool operator!=(const Record& otherRec) const {
31         // compare each field ...
32     }
33
34     unsigned s__supkey;
35     fluidb_string<25> s__name;
36     fluidb_string<40> s__address;

```

```

32   fluidb_string<16> s__city;
33   fluidb_string<16> s__nation;
34   fluidb_string<13> s__region;
35   fluidb_string<15> s__phone;
36 };

```

Listing 6.2.: The supplier row representation in the generated C++ code. The `fluidb_string` type is a constant size array of characters.

However, in the case of parsing we need to take great care with byte alignment which is compiler dependent. Fortunately Clang and GCC informally agree on the algorithm for `alignof(<cls>, <member>)` for standard layout objects. The algorithm is presented in 6.3

```

1  schemaPostPaddings :: [CppType] → Maybe [Int]
2  schemaPostPaddings [] = Just []
3  schemaPostPaddings [_] = Just [0]
4  schemaPostPaddings schema = do
5    elemSizes ← sequenceA [cppTypeSize t | t ← schema]
6    spaceAligns' ← sequenceA [cppTypeAlignment t | t ← schema]
7    let (_, spaceAligns) = spaceAligns' ++ [maximum spaceAligns']
8    let offsets = 0 : zipWith3 getOffset spaceAligns offsets elemSizes
9    return $ zipWith (-) (zipWith (-) (tail offsets) offsets) elemSizes
10   where
11     getOffset nextAlign off size =
12       (size + off)
13       + ((nextAlign - ((size + off) `mod` nextAlign)) `mod` nextAlign)

```

Listing 6.3.: Algorithm to infer the padding of members according to the Itanium ABI.

Once the bama files are generated C++ code is generated for parsing the file calls into the C++ function `bama_to_dat` that is parametric to the type of the object being and uses the BAMA library to write objects, thus making sure that the data is readable by the operators. `bama_to_dat` simply reads the input `.bama` file as a stream of records of size `sizeof(Record)` casting the bytes with `reinterpret_cast<Record*>` into the record. It then use the bame record writing facility `Writer<R>::write` that takes care of organizing the record into pages. Thus the final *data file* is created that is ready for use by the generated code.

```

1  #include <bamify.hh>
2  class Record {
3      // ...
4  };
5  int main(int argc, char* argv[]) {
6      bama_to_dat<Record>("supplier.bama","supplier.dat");
7  }

```

Listing 6.4.: Convert a bama file to a paged data file.

6.3.2. Pages

We translate the flat binary files into .dat files that contain the table data in its final format. The basic block of the file is the *Page*, and the file is simply a raw sequence of pages. A page is a constant-size data structure that contains up to $\left\lfloor \frac{S_{rec}}{S_{pg}} \right\rfloor$ whole records where S_{rec} is `sizeof(Record)` and S_{pg} is the size of the page which typically is 4KB. All pages must contain as many whole records as can fit except the last one.

All transactions with the storage are made at the page level: we either read or write only entire pages. These operations are abstracted by the `Reader` and `Writer` classes. We use one more level of abstraction for convenience, the higher order `eachRecord` function. To demonstrate what the interface looks like we present the implementation of `eachRecord` 6.5.

```

1  // Fn could be instantiated to std::function<void(const R&)> but that
2  // will *always* forbid f from being inlined.
3  template <typename R, typename Fn>
4  inline void eachRecord(const std::string& inFile, Fn f) {
5      Reader<R> reader;
6      size_t i = 0;
7      reader.open(inFile);
8      while (reader.hasNext()) {
9          i++;
10         f(reader.nextRecord());
11     }
12     reader.close();
13 }

```

Listing 6.5.

As alluded to in the previous section both the `Writer` and `Reader` use `reinterpret_cast` to "serialize" and "deserialize" the data respectively.

6.3.3. C++ row iterators

It is important for the planner to have full control of the storage budget and assume that no significant memory is required for evaluating each operator. For that reason we require that all operators' algorithms have constant space complexity. This may mean major compromises for some algorithms with respect to the time complexity like join and aggregation. Fortunately we can mitigate that by taking advantage of the set semantics of FluidDB relational algebra and sorting input tables in-place before running joins or aggregations. Our particular implementation of in-place sorting hinges on the `RecordMap` type that provides C++ random access iterators to the records of a file, abstracting the page reads and writes. We pass these iterators to `std::sort` (see listing 6.6) which runs insertion sort for small ranges to take advantage of the processors reorder buffer and quicksort for longer ones.

```
1 RecordMap<size_t> fs("/tmp/removeme.dat");
2 std::sort(fs.begin(), fs.end());
```

Listing 6.6.: Using a `RecordMap` to sort the records of a file.

The operation of `RecordMap` is very simple. Each iterator is paired with the page that contains the record it points to, when an iterator is the only iterator pointing inside a page and is incremented or decremented and no longer points to the same page the old page is written back to the file and the new page is read. When the only iterator pointing to a page is deleted the page is written to the file. More than one iterators pointing to the same page do not maintain different copies of that page and the last one to be destroyed or to leave the page triggers the page to be written back to the file.

6.4. Physical planning

The fundamental logic of the code generator is fairly simple: The input of the code generator is a a list of transitions generated by the FluidDB planner. The reader is

reminded that each transition has one of three kinds:

- Trigger of a t-node. The input n-nodes at the time of the trigger are materialized and the trigger itself materializes a subset of the output n-nodes
- Reverse trigger of a t-node, which represents the materialization of a subset of the input n-nodes from materialized output nodes
- The deletion of an n-node.

Transitions are meant to be executed in the order they appear in the received sequence to preserve correctness. However there is no one-to-one correspondence between operators and transitions. Instead, as we saw in section 3.5, each operation corresponds to a cluster of connected t-nodes. The first step of the code generator, therefore, is to group the low level transitions received by the planner into higher level batches that correspond to exactly one operation. This process is driven by intermediate n-nodes, ie helper nodes that do not correspond to materializable relations, but are rather part of the graph to reify the valid combinations that the planner is allowed to materialize. The constraint we are trying to preserve while batching the low level transitions such that the cluster level transitions do not materialize any intermediate nodes.

Since intermediate nodes are always internal to clusters and no cluster shares a t-node with another cluster, each batch of transitions corresponds to exactly one cluster, except for deletion transitions which are standalone. Furthermore, each cluster corresponds to exactly one relational operator, which we also include in the higher order transition.

6.5. Operator implementations

The C++ AST is expressed as a tree of haskell algebraic data types. They do not capture the entire C++ language, only the following concepts concepts:

- Function declaration
- Function application and arguments
- Expressions
- Code symbols

- Assignment
- Includes
- Literals
- Classes
- Member declarations (for record printing)
- Algorithm selection
- Global declarations

The code generated version of each operator is parameterized by the record types of its inputs and outputs and highly specific code is generated by the C++ templating system.

All operator implementations live in the BAMA library which branched out from [29]. The BAMA library makes heavy use of templates and constexpr to generate very query-specific machine code for each operator.

6.5.1. Generated code structure

The FluidB code generator generates a couple of different kinds of C++ components to construct the final file. In this section we will discuss the kinds of C++ structs that FluidB generates in order to parameterize the BAMA operators at compile time, leveraging the C++ template system.

In the following subsections when we say *compile time* we refer to the compile time of the generated code, not compile time of FluidB.

Maybe types

Maybe types are used to indicate optional outputs of operators. For example an operator σ_p may yield $\sigma_p A$ or $\sigma_{\neg p} A$ or both. On one hand the precise outputs that are required are important for performance and memory budget management, and on the other any combination of these outputs can and should be generated through one pass over the input. Furthermore information about which outputs are required is known at compile time so the generated code should be specific to the combination of outputs that are required.

To address this, instead of passing in a simple `std::string` as a path, we pass one of two types `Nothing` or `Just` (listing 6.7) that wrap the file path. Both these types contain an `isNothing` static `constexpr` type that can be used by `if constexpr (...)` expressions in the selection operator to help the compiler generate highly specialized machine code.

```

1  template <typename T>
2  struct Just {
3      Just(T t) : value(t) {}
4      T value;
5      constexpr bool operator==(Just<T> const& j) const {
6          return this->value == j.value;
7      }
8      static constexpr bool isNothing = false;
9  };
10 template <typename T=std::string>
11 struct Nothing {
12     Nothing(T s) {}
13     Nothing() {}
14     static constexpr bool isNothing = true;
15 };

```

Listing 6.7.: The type level maybe

Record types

FluidB is a row store engine and handles table rows by generating record types specific to each table, where objects of that type are rows of the corresponding table. We went over some aspects of the the generated record types when discussing the conversion of primary data to FluidB specific binary data. Records types directly or indirectly parameterize all operators. A selection operator σ_p instance must be parameterized by the record type of its input which will be the same as the output. A projection π is parameterized by the records of both the input and the output type. A generated record type must have some structural properties to function properly with BAMA and the rest of the generated types:

- First and foremost it must be standard-layout type provide so it can be trivially copied from and to binary blob files. This means that there should be no virtual

members, the destructor must be trivial and no static members. We further impose the constraint that they must not have a base class to avoid some more intricate constraints imposed by the C++ standard on standard-layout classes.

- it must provide non-uniquified names for each field of the row it represents so other generated types (like predicates and subtuple extraction types) have direct access to the underlying data.
- They should be comparable by equality ($=$), inequality (\neq) and they should be orderable $<$. The actual semantics of the ordering are not important as long as there is a deterministic way of ordering records of the same type.
- They should be hashable. We get that for free since they are standard layout objects but it is an important property that is useful for equijoins as we will see in the section about equi-joins.
- Finally every record type must implement an `std::string show()` function that serialized the contents into a human readable string.

Predicate types

Predicate types have the function of n-ary boolean functions $A_1 \times \dots \times A_k \rightarrow \{T, F\}$ useful for selections and θ -joins so in practice they are unary or binary. They are passed as template arguments to the operators so the compiler has the chance to inline them to avoid a function call. To avoid too much code clutter they are expected to provide the types of the function domain (demonstrated in listing 6.8) in order to minimize the number of template parameters and keep the generated code relatively human readable.

```

1  class Predicate3421 {
2      typedef Record123 Domain0;
3      typedef Record32 Domain1;
4      bool operator() (const Domain0& rl, const Domain1& rr) {
5          return rl.__field1 < rr.__field2;
6      }
7  }
```

Listing 6.8.: The shape of a generated predicate type.

Record transformation types

Record transformations, much like predicates are callable representing pure functions, only the codomain is a record type, instead of a boolean. They are used by joins to combine the matching records, by projections and aggregations to produce new records from the inputs, by equi-joins to extract the subtuples to be checked for equality, etc. Much like predicates they can be queried for their domain and codomain to reduce the number of template arguments in operators (see listing 6.9)

```

1  class Transform {
2      typedef Record123 Domain0;
3      typedef Record32 Domain1;
4      typedef Record10 Codomain;
5      Codomain operator() (const Domain0& l, const Domain1& r) {
6          return Record10(l.__key, l.__field1, r.__key, r.__field2);
7      }
8  };

```

Listing 6.9.: A record transformation type defines objects with no internal state that are callable.

Operators

Operators are not generated but they are parameterized by all the kinds of generated code we mentioned. They are classes that are constructed using maybe-filenames, record transformers, record types, and predicates and implement the `run()` method which actually runs the internal code. For demonstration purposes an abbreviated version of the σ_p operator, which highlights the shape of an operator class is presented in listing 6.10.

```

1  template <typename Predicate,
2           typename PrimaryOutType,    // Maybe(std::string)
3           typename SecondaryOutType  // Maybe(std::string)
4           >
5  class Select {
6      typedef typename Predicate::Domain0 Record;
7

```

```

8  public:
9      Select(PrimaryOutType prim, SecondaryOutType sec, std::string in)
10         : primary_file(prim), secondary_file(sec), infile(in) {
11         static_assert(!PrimaryOutType::isNothing || !SecondaryOutType::isNothing,
12             "Both primary and secondary output files are Nothing.");
13     }
14
15     ~Select() {}
16
17     void run() {
18         // ...
19     }
20
21     void print_output(size_t x) {
22         // ...
23     }
24
25 private:
26     PrimaryOutType primary_file;
27     SecondaryOutType secondary_file;
28     std::string infile;
29     static Predicate predicate;
30 };

```

Listing 6.10.: The selection operator. It is parameterized by the predicate and the primary and secondary output types. Enough information about these values is known at compile time such that the compiler can generate highly specialized code.

Unfortunately, even recent iterations of the C++ standard do not include template type inference for classes and structs, therefore, to make the generated code simpler we wrap the constructor into a function like shown in listing 6.11. This allows the generated code for selection to look like the one presented in listing 6.12.

```

1  // C++17 can only infer typenamees (primaryout secondaryout) in
2  // function templates.

```

```
3  template<typename Predicate,
4          typename PrimaryOutType,    // Maybe(std::string)
5          typename SecondaryOutType // Maybe(std::string)
6          >
7  auto mkSelect (const PrimaryOutType prim,
8                const SecondaryOutType sec,
9                const std::string& in) {
10     return Select<Predicate,
11                 PrimaryOutType,    // Maybe(std::string)
12                 SecondaryOutType> // Maybe(std::string)
13     (prim, sec, in);
14 }
```

Listing 6.11.: The C++ declaration of the select.

```
1  int main() {
2      // ...
3      {
4          auto op = mkSelect<Predicate32>(Just("data123.dat"), Just("data53.dat"),
5                                          "lineitem.dat");
6          op.run();
7          op.print_output();
8      }
9      // ...
10     return 0;
11 }
```

Listing 6.12.: A block representing a particular operator.

6.5.2. Select algorithm

The selection algorithm is likely the simplest of the implemented ones because FluidB does not support indexes and makes no assumptions about the ordering of the data. In the forward variety it is implemented as either a partition or a selection depending on which of the outputs it is materializing. In its backward variety it is essentially a union

$A \equiv \sigma_p A \cup \sigma_{\neg p} A$. As we discussed already (listing 6.11) the BAMA library expects only a predicate class as a template parameter and three filenames as its arguments.

It is important to remind the reader that no order is expected to be preserved by union, which is to say that the transition $A \Rightarrow \{\sigma_p A, \sigma_{\neg p} A\} \Rightarrow A$ does not preserve the order of the nodes. The reader is also reminded that FluidDB does not support **NULL** values to preserve the correctness of the reverse select operation.

6.5.3. Projection algorithm

The implementation of the projection algorithm is parameterized by two template arguments that extract complementary subtuples from an input relation, affording each one with a common unique subtuple to facilitate the inverse operation. The inverse of a projection then is simply an equijoin between the two produced slices based in that shared unique subtuple.

While the BAMA side of the calculation is fairly straightforward, on the Haskell side (code generation) it is slightly more complex. The first piece of the puzzle is the representation of the projection. When the query is initially processed so that projections are augmented to expose unique subtuple we actually change the representation of projections from **QProj** [(e, **Expr** e)] to **QProjI** [e] [e] [(e, **Expr** e)]. The two extra parameters typed [e] represent the complement of the projection (the columns from the input not exported by the projection) and a unique subtuple that the complement and the projection must have in common, while [(e, **Expr** e)] is the normal projection. We make this transformation because the easiest and cheapest place to calculate the first and second parameters (complement and correspondence) is during the projection augmentation phase that happens immediately after parsing. Therefore we have all projections carry information about the complement at the level of the RA operator.

In particular this process is combined with remapping of unique subtuples into the projection. For example, during the unique column exposure process the query **select** p_color **from** part is found not to be valid because p_color is not unique for each row. For that reason, during preprocessing we remap the p_partkey column since p_partkey is the only unique column of part this way the constraint that every valid relation in FluidDB needs to have at least one subtuple that is unique within that relation is satisfied and the query becomes

```
1 select p_partkey, p_color from part
```

In order to be able to reconstruct the input we require that at least one unique subtuple is shared between the primary projection relation and the complement. In the example demonstrated this means the `p_partkey` column is duplicated in both output tables of the operator. The particular valid projection operator that would be generated the **QProj** is presented in listing 6.13.

```

1  op =
2    QProjI
3    -- All keys except p_color
4    [p_partkey,p_name,p_mfgr,p_category,p_brand1,p_type,p_size,p_container]
5    -- The unique subtuple (just one column)
6    [p_partkey]
7    -- The the actual projection
8    [(p_partkey,E0 p_partkey),(p_color,E0 p_color)]

```

Listing 6.13.: The projection operator produced from the SQL query `select p_partkey,p_color from part`.

The algorithm for reverse projection is almost a standard join. The only difference is that in our semantics of join, the relation $A \bowtie_{a=b} B$ contains both the column a and the column b , while when joining a projection with it's complement we only create a single copy of the shared subtuple. We will see more details about join in the next section.

6.5.4. The join operator

We distinguish between two kinds of join operators: the equijoin and the general θ -join. The former is specific to cases where subtuples from each table are compared and the implementation is simply a merge join, while the latter is used for general predicates and is essentially a pipelined $\sigma(A \times B)$. The important detail worth noting about join algorithm implementations in FluiDB is that they must be able to produce antijoin relations $A \ltimes B$, $A \ltimes B$ simultaneously with the join relation $A \bowtie B$.

It is important for join to be space efficient, and particularly constant-space, due to the assumption made by the planner that it has control over memory. For general θ -joins we use nested loop joins. For equijoins we take advantage of the assumption that the tables are unordered and use merge join after sorting the input tables in-place. A sketch of the algorithm is presented in listing 6.14


```

1  # This is a pythonic way of detecting which input stream finished is.
2  class LeftFinished:
3      pass
4  class RightFinished:
5      pass
6
7
8  def lnext(it):
9      """If iteration finishes throw LeftFinished"""
10     try:
11         return next(x)
12     except:
13         raise LeftFinished
14
15  def rnext(it):
16      """If iteration finishes throw RightFinished"""
17     try:
18         return next(x)
19     except:
20         raise RightFinished
21
22  @template(extractl,extractr,combine)
23  def merge_equijoin(in_l,in_r,anti_l,out,anti_r):
24      # In place sort each of the inputs
25      sort(in_l,key=extractl)
26      sort(out_l,key=extractr)
27      # Define iterators and current values for each of the inputs.
28      it_l = iter(in_l)
29      val_l = next(it_l)
30      it_r = iter(in_r)
31      val_r = next(it_r)
32
33     try:
34         # The actual merge join algorithm

```

```

35     while True:
36         # Gather blocks of equal records from left and right and write
37         # out their product.
38         if extractl(val_l) == extractr(val_r):
39             # Gather a sequence of equal-key values from the left
40             val_ls = []
41             tmp = val_l
42
43             while extractl(val_l) == extractl(tmp):
44                 val_ls.append(val_l)
45                 val_l = lnext(it_l, None)
46
47             # Gather a sequence of equal-key values from the right
48             tmp = val_r
49             val_rs = []
50             while extractl(val_r) == extractl(tmp):
51                 val_rs.append(val_r)
52                 val_r = rnext(it_r)
53
54             # Write out their product.
55             for l,r in product(val_ls, val_rs):
56                 out.write(combine(l,r))
57
58             # Push the non-equal records to the antijoins
59             while extractl(val_l) < extractr(val_r):
60                 anti_l.write(val_r)
61                 val_r = rnext(it_r)
62             while extractl(val_r) < extractr(val_l):
63                 anti_r.write(val_l)
64                 val_l = lnext(it_l)
65         except LeftFinished:
66             for l in it_l:
67                 left_anti.write(r)
68         except RightFinished:

```

```

69         for r in it_r:
70             right_anti.write(r)

```

Listing 6.14.: The equi-join algorithm first sorts in place the inputs w.r.t. the equal subtuples and then merges them.

The type of the join operator is defined in terms of a template such that it can be specialized at compile time (see listing 6.15). The template arguments are the following:

- Three file paths are provided, `OutFile`, `LeftAnti` and `RightAnti` are each either of type `Nothing` indicating to the compiler not to generate any code that relates to the particular output or of type `Just<std::string>` indicating to the compiler the opposite.
- The subtuple extraction functions `LeftExtract` and `RightExtract`. These are record transformation functions that extract the subtuple that needs to

```

1  template <typename LeftExtract, typename RightExtract, typename Combine,
2          typename OutFile,      // Maybe(std::string)
3          typename LeftAnti,     // Maybe(std::string)
4          typename RightAnti>    // Maybe(std::string)
5  class Join;

```

Listing 6.15.: Class declaration of the join operator

The reverse join operator is composed of two pipelined steps:

- A projection and deduplication based on the unique subtuple of the input on the join output to get a left or right semijoin
- A union of the semijoin with the antijoin to get the input table.

As with the selection operator, the correctness of this reverse operation is predicated on the fact that FluidDB has no notion of `NULL` values.

6.5.5. Aggregation and sort algorithms

The main challenge with implementing aggregations in the context of FluiDB is maintaining constant space during its evaluation, as the planner tries to use as much of the memory budget as possible to maintain intermediate results. Therefore we opted against using an auxiliary hash map used to group records. Like we did to implement joins we perform aggregation in two steps: first we sort the records in place based on the grouping columns and then aggregate in a single pass.

6.6. Conclusion

In this chapter we discussed some physical attributes of how FluiDB evaluates the final plans it constructs via code generation. We discussed the storage layout and the fundamental ideas the code generator was built around in order to generate highly specialized code. Finally, we discussed some specific algorithms used by the BAMA library that is called by the generated code to implement the actual operators.

We further discussed the some caveats that relate to the FluiDB's approach to code generation, namely that highly the specific code generated can be expensive to compile from scratch and that the shmем in-memory filesystem provided by linux is a cheap solution in terms of engineering effort but would be better replaced with a solution that does not involve system calls.

Evaluation

Men know the material past,
future deeds are known by gods

(C. P. Kavafy)

Chapter summary

- FluiDB is well suited for join-heavy star schemata so we evaluated using the SSB-TPC-H benchmark.
- Our evaluation shows that FluiDB is able to plan around the space constraints and come up with plans that materialize intermediate results that are useful for future queries.
- FluiDB is generally faster than the baseline due to caching but at times may be slower in individual queries when it receives “unexpected” queries.
- FluiDB generally performs better when allowed larger memory budgets but this speedup is based on heuristic assumptions that sometimes break in interesting ways.

We based our evaluation on the Star Schema Benchmarkx (SSB) [47] which is a variation of TPC-H. Like TPC-H it models the data warehouse of a wholesale supplier.

As the name suggests these queries query a star schema. The star schema is derived from the TPC-H schema (figure 7.2) by merging the `linenumber`, `orders` and `partsupp` and completely dropped tables `nation` and `region` tables as shown in figure 7.1.

The particular queries contained in the workload are presented in listing A.1 in the appendix (A).

We generated a workload where queries were compiled and planned one by one in a loop for and run them over a database of data generated with a modified version of the TPC-H dbgen [78] with scaling of size 1, meaning that the total primary data is around 1GB in CSV format. The exact size of the data generated after formatting them to the format required by FluiDB (described in chapter 6) are

```
1 $ du -sh *.dat
2 4.1M    customer.dat
3 256K    date.dat
4 757M    lineorder.dat
5 31M     part.dat
6 268K    supplier.dat
```

Which makes a total of approximately 200k pages.

The lowest budget within which FluiDB was able to plan all 13 queries of SSB-TPC-H was 2300k pages.

In our experiment we use page IO as a proxy for performance, despite the fact that FluiDB is an in-memory database. We think that this is a reasonable experimental approach because, as FluiDB leans heavily on code generation, it is unlikely that the actual instruction retiring will have a major impact on the performance. Instead the performance cost is dominated by page IO that will certainly cause cache misses at the LLC. Another thing to note is that FluiDB generates code that focuses on performance and not on compilation time, making heavy use of metaprogramming like `constexpr` and templates. This makes for fairly slow compilation times. There are ways to speed up compilation time like using pre-compiling header files [79] and fine-tuning the compiler optimization passes. These techniques are beyond the scope of this work, FluiDB is focused on analytics workloads that do not include sub-second queries.

Figure 7.3 demonstrates that FluiDB running a workload versus running individual

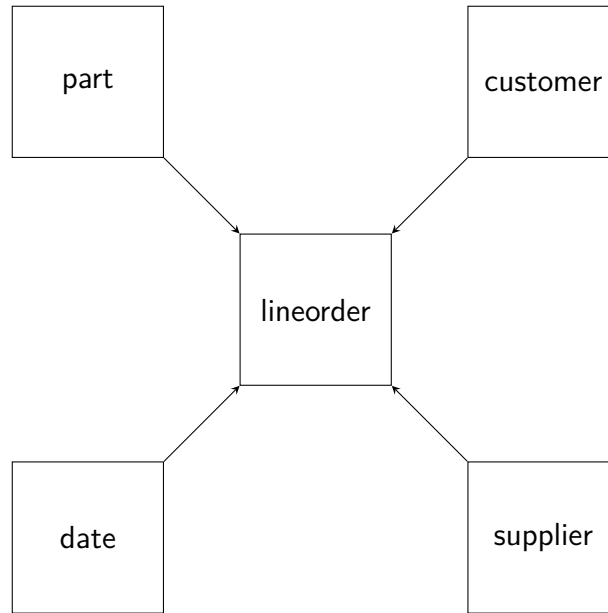


Figure 7.1.: The foreign key links in a SSB-TPC-H

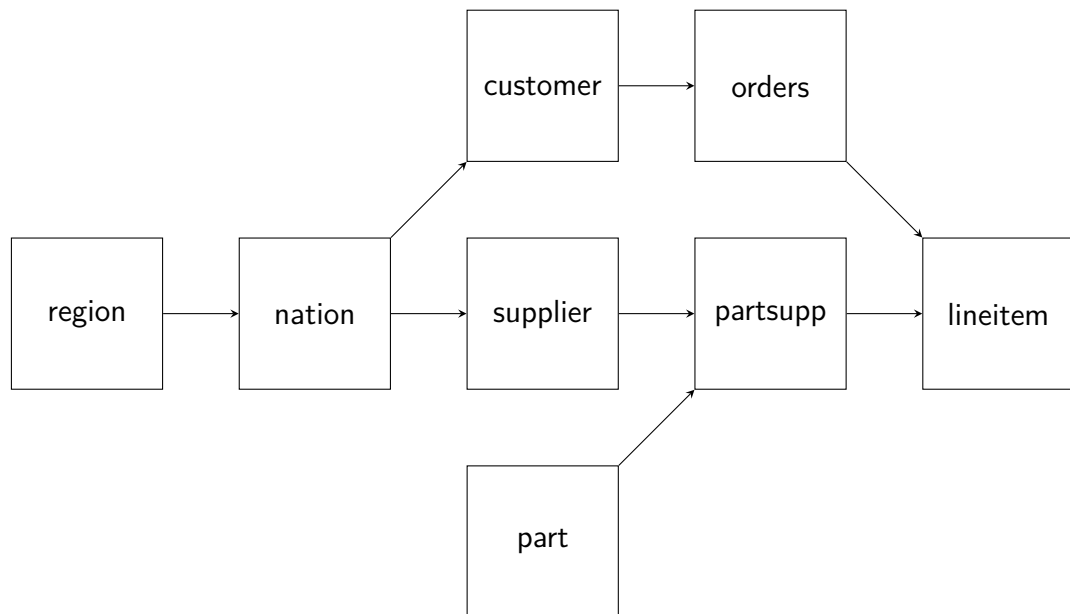


Figure 7.2.: The foreign key links in a traditional TPC-H schema

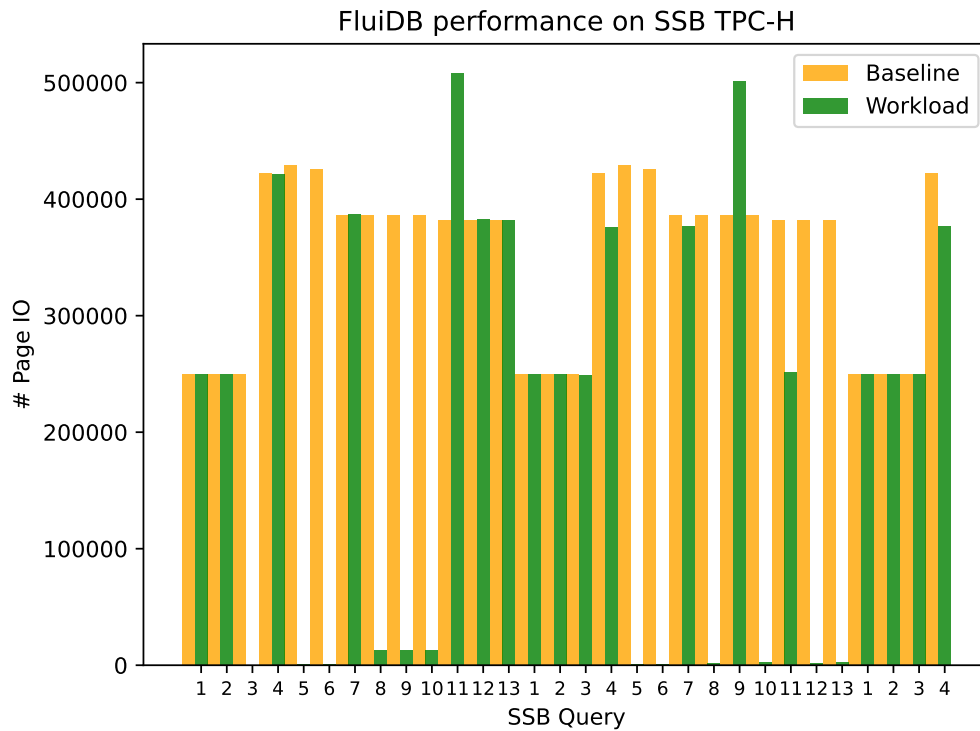


Figure 7.3.: The total page read/writes for each query a sequence of queries. The baseline is the query being run by FluiDB directly without any materialized nodes. The workload bars represent the cost of each query in a workload being accumulated into the same QDAG. The budget allowed for this is the minimum budget within which FluiDB can run each individual query (2300K pages).

queries causes significant speedups even in constrained budgets. In some cases, however the garbage collector is forced to delete tables that need to be recreated later in the workload causing FluiDB to be sporadically less performant than the base case.

For a larger budget the FluiDB is able to store more useful intermediate results as demonstrated in figure 7.4

An interesting point here is the plan for evaluating query 7 is more expensive in the workload run under laxer budgetary constraints. This may seem strange but it is an example demonstrative of the fundamental operation of FluiDB.

The top level explanation is that the lax-budget plan needs to evaluate *lineorder* via the reverse trigger of a join as it was deleted during evaluation of query 6. Mysteriously during the strict-budget planning the *lineorder* relation is readily available during planning of query 7! The key lies slightly earlier in the workload (see listing 7.1).

FluiDB materializes in query 4 the join

$$Q_{36} := \text{supplier} \bowtie_{lo_suppkey=s_suppkey} \text{lineorder}$$

and corresponding antijoins Q_{35} and Q_{37} making the node *lineorder* deletable. However, when running the workload under strict budgetary constraints, it is forced to garbage collect shortly after materializing said join at a moment while both Q and *lineorder* are protected (see section 4.2.1 on garbage collection). Therefore, FluiDB is forced to delete both Q_{35} and Q_{37} . Making *lineorder* non-deletable when planning for query 4 finishes.

```

1  # Query 4
2  Query sγπσ(supplier ⋈ lineorder ⋈ date ⋈ part) {
3    # There is enough space to keep both and the complements
4    Q36, Q35, Q37 :=
      ↪ Materialize[supplier ⋈ lineorder, supplier⋈lineorder, supplier⋈lineorder]
5    Q41, Q40, Q42 := Materialize[Q36 ⋈ date, Q36⋈date, Q36⋈date]
6    Q46, Q45, Q47 := Materialize[Q41 ⋈ part, Q41⋈part, Q41⋈part]
7    Q50 := Materialize[σQ46]
8    Q90 := Materialize[γπQ50]
9    Q91 := Materialize[sQ90]
10 }
11
12 # Query 5

```

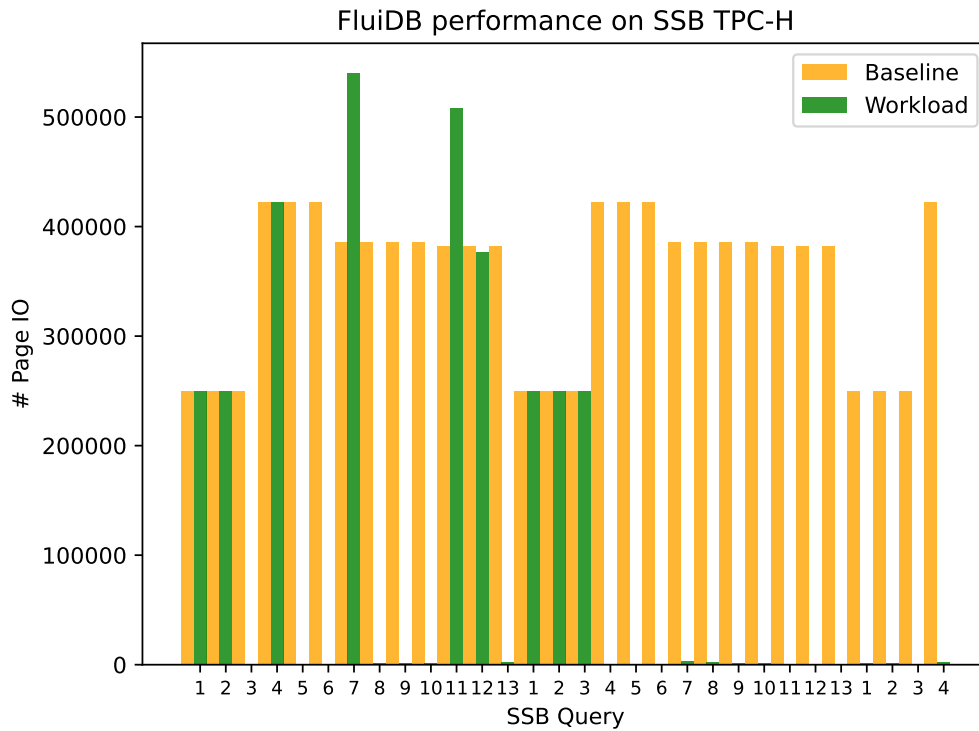


Figure 7.4.: The total page read/writes for each query a sequence of queries. The baseline is the query being run by FluiDB directly without any materialized nodes. The workload bars represent the cost of each query in a workload being accumulated into the same QDAG. The budget allowed for this is about triple the minimum budget within which FluiDB can run each individual query (6100K pages). A badly timed GC run causes some individual queries to be slower than their counterparts from a workload being run in a tight budget (figure 7.3).

```

13 Query  $s\gamma\pi\sigma(\text{supplier} \bowtie \text{lineorder} \bowtie \text{date} \bowtie \text{part})$  {
14    $Q_{92} := \text{Materialize}[\sigma Q_{46}]$ 
15    $Q_{118} := \text{Materialize}[\gamma\pi Q_{92}]$ 
16    $Q_{119} := \text{Materialize}[sQ_{118}]$ 
17 }
18
19 # Query 6
20 Query  $s\gamma\pi\sigma(\text{supplier} \bowtie \text{lineorder} \bowtie \text{date} \bowtie \text{part})$  {
21   # FluidB decides delete lineorder since it has the complements
22   GC { Delete[..., lineorder, ...] }
23    $Q_{120} := \text{Materialize}[\sigma Q_{46}]$ 
24    $Q_{146} := \text{Materialize}[\gamma\pi Q_{120}]$ 
25    $Q_{147} := \text{Materialize}[sQ_{146}]$ 
26 }
27
28 # Query 7
29 Query  $s\gamma\pi\sigma(\text{customer} \bowtie \text{date} \bowtie \text{lineorder} \bowtie \text{supplier})$  {
30   # Oops ... this would be avoided if we hadn't deleted lineorder.
31    $\text{lineorder} := \text{Materialize}[\bar{\pi}Q_{36} \cup Q_{37}]$ 
32    $\text{date} := \text{Materialize}[\bar{\pi}Q_{41} \cup Q_{42}]$ 
33   GC { Delete[...] }
34    $Q_{149}, Q_{148}, Q_{150} := \text{Materialize}[\text{date} \bowtie \text{lineorder}, \text{date} \Join \text{lineorder}, \text{date} \Join \text{lineorder}]$ 
35   GC { Delete[...] }
36    $Q_{154}, Q_{153}, Q_{155} :=$ 
37      $\hookrightarrow \text{Materialize}[\text{customer} \bowtie Q_{149}, \text{customer} \Join Q_{149}, \text{customer} \Join Q_{149}]$ 
38   GC { Delete[...] }
39    $Q_{165} := \text{Materialize}[\sigma Q_{154}]$ 
40    $Q_{182}, Q_{181}, Q_{183} := \text{Materialize}[Q_{165} \bowtie \text{supplier}, Q_{165} \Join \text{supplier}, Q_{165} \Join \text{supplier}]$ 
41   GC { Delete[...] }
42    $Q_{163} := \text{Materialize}[\sigma Q_{182}]$ 
43    $Q_{203} := \text{Materialize}[\gamma\pi Q_{163}]$ 
44    $Q_{204} := \text{Materialize}[sQ_{203}]$ 
45 }

```

Listing 7.1.: Abbreviated version of the plans of queries 4 to 7 of SSB TPC-H. This demonstrates how an unfortunately timed GC can cause FluidDB to make some bad decisions

On the other hand, with laxer budgetary constraints, no garbage collection is triggered during query 4. In fact the next garbage collection is triggered during query 6, at a moment when *lineorder* is deletable, unprotected, and a prime candidate for deletion based on the GC heuristics.

Alas, when query 7 requires *lineorder* for its plan, FluidDB needs to re-construct it in the case of lax budgetary constraints but, not in the case of strict constraints.

This example of FluidDB being forced to locally produce more expensive plans is an effect of FluidDB being more opportunistic the lower the available budget is, and more adventurous when operating with high budgets. When FluidDB is frugal it is generally prone to miss opportunities to share computation between queries. There are times however that this fugality saves it from bad heuristic-based decisions that it is allowed to make otherwise.

FluidDB aspires to deal with the entire workload as if it were planning a single query. While any decision during the planning of a single query can be scrapped in the backtracking process, FluidDB is tragically forced to commit to whatever adventurous or conservative decisions it makes at the end of every planning iteration, doomed to pay dearly for every misstep but to reap the rewards of every insightful choice.

Interestingly this problem goes away when we run with a budget of 3500k pages (figure 7.5) as the GC run is delayed to a time when other relations are better candidates for deletion than *lineorder*. A carefully designed set of heuristics for the GC should avoid this problem in most cases.

7.1. A note on size estimation

The size of the budget may seem fairly excessive for the size required for the primary tables. To understand why that is we need to look at the largest nodes in the QDAG:

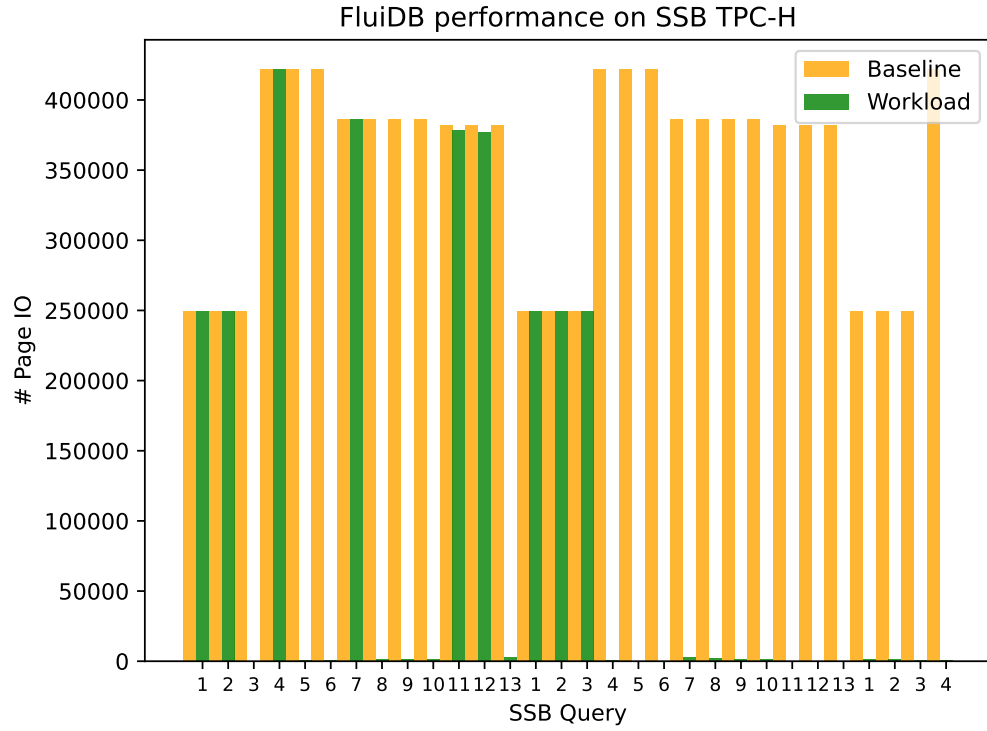


Figure 7.5.: The total page read/writes for each query a sequence of queries. The baseline is the query being run by FluiDB directly without any materialized nodes. The workload bars represent the cost of each query in a workload being accumulated into the same QDAG.. The budget allowed for this is about triple the minimum budget within which FluiDB can run each individual query (6500K pages).

Pages	Expr
188000	<i>lineorder</i>
547100	<i>customer</i> ⋈ (<i>date</i> ⋈ <i>lineorder</i>)
376200	$\sigma((\textit{customer} \bowtie (\textit{date} \bowtie \textit{lineorder})) \bowtie \textit{supplier})$
376200	$\sigma((\textit{date} \bowtie \textit{lineorder})) \bowtie \textit{supplier}$
273600	$\sigma(\textit{customer} \bowtie (\textit{date} \bowtie \textit{lineorder}))$
188100	$\sigma((\sigma(\textit{customer} \bowtie (\textit{date} \bowtie \textit{lineorder}))) \bowtie \textit{supplier})$

From looking at those nodes it seems that a sufficiently advanced garbage collector should be able to support plans that materialize the nodes in about half the budget as FluidDB should have been able to plan the query needing around double the pages of the largest equijoin which for us is *customer date lineorder*.

The main issue however, as it is with many query processing systems [48], is the cardinality estimator. With a more sophisticated cardinality estimator the required pages. Which assumes that in a natural join there are no foreign key lookup failures, that is, that all natural joins are extension joins. Therefore

$$|\sigma_{p(\textit{customer})}(\textit{customer} \bowtie (\textit{date} \bowtie \textit{lineorder}))| = |(\sigma_{p(\textit{customer})}\textit{customer} \bowtie (\textit{date} \bowtie \textit{lineorder}))|$$

This causes FluidDB to vastly overestimate the size of plans that have selections pushed down making these plans disproportionately unappealing

7.2. Conclusion

FluidDB is a very complex database system built from scratch focused on a specific idea: making optimal use of the storage budget by fully adapting the data layout to the workload. While FluidDB itself it is far from stable enough system to be used in a production setting, the results presented in this chapter demonstrate that it is an idea worth considering to incorporate in the design of a database system.

Conclusions and future perspectives

We had our fun

(Someone Else)

Chapter summary

- Summary of FluidDB concepts.
- Future work.

In this thesis we covered the important and interesting aspects of FluiDB. Implementing the system itself involved much work that is not covered in this thesis in order to not flood the reader with details that might remove focus from the important aspects of the system. Some of these are:

- Implementation of a backtracking-based, caching pretty printing library.
- A query decorrelation algorithm.
- A set of tools for debugging and tracing backtracking algorithms.

Aside from those, in the introduction we discussed some examples where data layout can radically adapt to the workload by making use of reversible operations and integrating a garbage collector into the planner. We further introduced the relational database model and provided an overview of the various stages of query processing. We focused two subfields of database research that are of particular relevance to FluiDB: code generation for in-memory databases, and intermediate result recycling.

The first non-introductory chapter (chapter 3) goes in depth about the semantics of reversible operations and how they, along with possible intermediate results for each query form the QDAG. It covers topics that relate to the construction of the QDAG, the use of the QDAG for enumerating plans, higher levels of organization of the QDAG nodes into clusters and how those clusters are used by FluiDB to infer properties of intermediate query results like cardinality and relation extent using propagators.

Chapter 4 on physical planning starts off by introducing a novel framework for implementing search algorithms, the **HCntT** monad transformer. Computation expressed in terms of **HCntT** can dynamically switch between branches of the backtracking algorithm essentially implementing a weighted search. The main novelty lies in the implementation of **once** and **fallback** operations that, to our knowledge, are not supported by any similar weighted search framework.

Based on **HCntT** we describe the implementation of the physical planner that lies at the core of FluiDB. The physical planner uses an A^* -like algorithm to traverse the QDAG in order to enumerate query plans. Query plan fragments are interleaved with GC runs that produce sequences of table deletions when the planner detects that the storage budget would be exceeded.

The decisions made by the planner are based on computations that depend on the inventory of materialized relations, namely variations of intermediate result cost estimation

and intermediate result materializability. It makes sense to implement computations in terms of an incremental computation framework as there are major opportunities for sharing between evaluations. Furthermore, the performance of these computations is sensitive to the order of evaluation due to absorbing elements and `min` operations. To solve these problems we developed *Antisthenis* which was described in detail in chapter 5.

We wrap up the technical description of FluiDB with chapter 6 that goes in depth about how FluiDB transpiles physical plans into C++ and what algorithms are used to implement each operation and their reverse.

Finally, we present some experimental results to demonstrate the performance characteristics of FluiDB and any system that might implement similar concepts to it. Particularly, we run the SSB TPC-H benchmark and show that FluiDB can successfully outperform a similar system that only does per-query optimization in amortized cost and even in most per-query cases. We also demonstrate the expected result that FluiDB's repertoire of heuristics may cause individual queries planned within the context of a workload to be less performant than the same query planned in isolation. We further demonstrate how FluiDB usually succeeds at making better use of larger storage budgets but occasionally, like every heuristic based system, will make decisions that hurt performance.

8.1. Future perspectives

FluiDB steps into a new path towards what we believe to be a *truly* adaptive storage, but it does only a few steps. There are many interesting tangents to follow and many places where FluiDB does not take full advantage of the state of the art in database research.

Starting with the latter, the most important place where FluiDB can do better is cardinality estimation. As mentioned in section (7.1) FluiDB is very naive w.r.t. to how it estimates sizes. That said the propagator network structure used to estimate cardinalities can be augmented to propagate information about selectivity and statistics about each node. These statistics could be collected for every materialized intermediate relation to enrich the much needed information available during planning.

Another low hanging fruit that can be exploited is parallelism. All plans that FluiDB produces are single threaded but there is no fundamental reason for that to be so. A simple data dependency analysis could reveal opportunities for parallelism and operators themselves could be easily replaced with parallel versions. A more complex but very

interesting problem would be how to teach the planner to take this into account in order to produce plans that aren't just accidentally parallel, but are designed taking such speedups into account.

Another important shortcoming of FluiDB as it was implemented is its complete lack of support for updates. There are heaps of literature to draw from and many excellent approaches to the problem. It would be interesting to find out which of these and how they could be applied to the peculiar case of FluiDB's data model where primary tables are likely not materialized in memory.

Furthermore, a system based on FluiDB could take advantage of non-relational operations like table compression, materialization of indexes, movement of data in a memory hierarchy in order to produce plans that can better take advantage of the various possibilities provided by the state of the art in database research.

Finally, a major bottleneck especially for cheap queries, is the latency of the C++ compiler that can take up to 10 sec to build a query plan into an executable. The FluiDB code generation make absolutely no attempt to make life easier for the C++ compiler, making heavy use of metaprogramming facilities like templates, traits, and `constexpr` expressions. There has been a lot of work to mitigate this problem, in the literature, some good examples of which are outlined in the background section [2.3.1](#)

Appendix

```
1  -- Flight 1
2  -- Query 1
3  select sum(lo_extendedprice*lo_discount) as revenue
4  from lineorder, date
5  where lo_orderdate = d_datekey
6  and d_yearmonthnum = 199401
7  and lo_discount between 4 and 6
8  and lo_quantity between 26 and 35
9
10 -- Query 2
11 select sum(lo_extendedprice*lo_discount) as revenue
12 from lineorder, date
13 where lo_orderdate = d_datekey
14 and d_year = 1993
15 and lo_discount between 1 and 3
16 and lo_quantity < 25
17
18 -- Query 3
19 select sum(lo_extendedprice*lo_discount) as revenue
20 from lineorder, date
21 where lo_orderdate = d_datekey
22 and d_yearmonthnum = 199401
```

Appendix A. Appendix

```
23 and lo_discount between 4 and 6
24 and lo_quantity between 26 and 35
25
26 -- Flight 2
27 -- Query 4
28 select sum(lo_revenue), d_year, p_brand1
29 from lineorder, date, part, supplier
30 where lo_orderdate = d_datekey
31 and lo_partkey = p_partkey
32 and lo_suppkey = s_suppkey
33 and p_category = 'MFGR#12'
34 and s_region = 'AMERICA'
35 group by d_year, p_brand1
36 order by d_year, p_brand1
37
38 -- Query 5
39 select sum(lo_revenue), d_year, p_brand1
40 from lineorder, date, part, supplier
41 where lo_orderdate = d_datekey
42 and lo_partkey = p_partkey
43 and lo_suppkey = s_suppkey
44 and p_brand1 between 'MFGR#2221' and 'MFGR#2228'
45 and s_region = 'ASIA'
46 group by d_year, p_brand1
47 order by d_year, p_brand1
48
49 -- Query 6
50 select sum(lo_revenue), d_year, p_brand1
51 from lineorder, date, part, supplier
52 where lo_orderdate = d_datekey
53 and lo_partkey = p_partkey
54 and lo_suppkey = s_suppkey
55 and p_brand1 = 'MFGR#2221'
56 and s_region = 'EUROPE'
```

```

57  group by d_year, p_brand1
58  order by d_year, p_brand1
59
60  -- Flight 3
61  -- Query 7
62  select c_nation, s_nation, d_year, sum(lo_revenue) as revenue
63  from customer, lineorder, supplier, date
64  where lo_custkey = c_custkey
65  and lo_suppkey = s_suppkey
66  and lo_orderdate = d_datekey
67  and c_region = 'ASIA' and s_region = 'ASIA'
68  and d_year ≥ 1992 and d_year ≤ 1997
69  group by c_nation, s_nation, d_year
70  order by d_year, revenue desc
71
72  -- Query 8
73  select c_city, s_city, d_year, sum(lo_revenue) as revenue
74  from customer, lineorder, supplier, date
75  where lo_custkey = c_custkey
76  and lo_suppkey = s_suppkey
77  and lo_orderdate = d_datekey
78  and c_nation = 'UNITED STATES'
79  and s_nation = 'UNITED STATES'
80  and d_year ≥ 1992 and d_year ≤ 1997
81  group by c_city, s_city, d_year
82  order by d_year, revenue desc
83
84  -- Query 9
85  select c_city, s_city, d_year, sum(lo_revenue) as revenue
86  from customer, lineorder, supplier, date
87  where lo_custkey = c_custkey
88  and lo_suppkey = s_suppkey
89  and lo_orderdate = d_datekey
90  and (c_city='UNITED KI1' or c_city='UNITED KI5')

```

Appendix A. Appendix

```
91 and (s_city='UNITED KI1' or s_city='UNITED KI5')
92 and d_year ≥ 1992 and d_year ≤ 1997
93 group by c_city, s_city, d_year
94 order by d_year, revenue desc
95
96 -- Query 10
97 select c_city, s_city, d_year, sum(lo_revenue) as revenue
98 from customer, lineorder, supplier, date
99 where lo_custkey = c_custkey
100 and lo_suppkey = s_suppkey
101 and lo_orderdate = d_datekey
102 and (c_city='UNITED KI1' or c_city='UNITED KI5')
103 and (s_city='UNITED KI1' or s_city='UNITED KI5')
104 and d_yearmonth = 'Dec1997'
105 group by c_city, s_city, d_year
106 order by d_year, revenue desc
107
108 -- Flight 4
109 -- Query 11
110 select d_year, c_nation, sum(lo_revenue - lo_supplycost) as profit
111 from date, customer, supplier, part, lineorder
112 where lo_custkey = c_custkey
113 and lo_suppkey = s_suppkey
114 and lo_partkey = p_partkey
115 and lo_orderdate = d_datekey
116 and c_region = 'AMERICA'
117 and s_region = 'AMERICA'
118 and (p_mfgr = 'MFGR#1' or p_mfgr = 'MFGR#2')
119 group by d_year, c_nation
120 order by d_year, c_nation
121
122 -- Query 12
123 select d_year, s_nation, p_category, sum(lo_revenue - lo_supplycost) as profit
124 from date, customer, supplier, part, lineorder
```

```

125 where lo_custkey = c_custkey
126 and lo_suppkey = s_suppkey
127 and lo_partkey = p_partkey
128 and lo_orderdate = d_datekey
129 and c_region = 'AMERICA'
130 and s_region = 'AMERICA'
131 and (d_year = 1997 or d_year = 1998)
132 and (p_mfgr = 'MFGR#1'
133 or p_mfgr = 'MFGR#2')
134 group by d_year, s_nation, p_category order by d_year, s_nation, p_category
135
136 -- Query 13
137 select d_year, s_city, p_brand1, sum(lo_revenue - lo_supplycost) as profit
138 from date, customer, supplier, part, lineorder
139 where lo_custkey = c_custkey
140 and lo_suppkey = s_suppkey
141 and lo_partkey = p_partkey
142 and lo_orderdate = d_datekey
143 and c_region = 'AMERICA'
144 and s_nation = 'UNITED STATES'
145 and (d_year = 1997 or d_year = 1998)
146 and p_category = 'MFGR#14'
147 group by d_year, s_city, p_brand1 order by d_year, s_city, p_brand1

```

Listing A.1.: The SQL code for the SSB queries

Bibliography

- [1] Yehoshua Sagiv and Mihalis Yannakakis. “Equivalences Among Relational Expressions with the Union and Difference Operators”. In: *Journal of the ACM* 27.4 (Oct. 1980), pages 633–655.
- [2] Donald D. Chamberlin et al. “A History and Evaluation of System R”. In: *Communications of the ACM* 24.10 (Oct. 1, 1981), pages 632–646.
- [3] Christopher T. Haynes. “Logic Continuations”. In: *The Journal of Logic Programming* 4.2 (June 1, 1987), pages 157–176.
- [4] G. Graefe and K. Ward. “Dynamic Query Evaluation Plans”. In: *ACM SIGMOD Record* 18.2 (June 1, 1989), pages 358–366.
- [5] W. Pugh and T. Teitelbaum. “Incremental Computation via Function Caching”. In: *Proceedings of the 16th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages - POPL '89*. The 16th ACM SIGPLAN-SIGACT Symposium. Austin, Texas, United States: ACM Press, 1989, pages 315–328.
- [6] G. Graefe and W.J. McKenna. “The Volcano Optimizer Generator: Extensibility and Efficient Search”. In: *Proceedings of IEEE 9th International Conference on Data Engineering*. Proceedings of IEEE 9th International Conference on Data Engineering. Apr. 1993, pages 209–218.
- [7] Goetz Graefe. “The Cascades Framework for Query Optimization”. In: *IEEE Data Eng. Bull.* 18.3 (1995), pages 19–29.
- [8] Sheng Liang, Paul Hudak, and Mark Jones. “Monad Transformers and Modular Interpreters”. In: *Proceedings of the 22nd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*. POPL '95. New York, NY, USA: Association for Computing Machinery, Jan. 25, 1995, pages 333–343.
- [9] Himanshu Gupta. “Selection of Views to Materialize in a Data Warehouse”. In: *Database Theory — ICDT '97*. Edited by Foto Afrati and Phokion Kolaitis. Lecture Notes in Computer Science. Berlin, Heidelberg: Springer, 1997, pages 98–112.
- [10] Gérard Huet. “The Zipper”. In: *Journal of Functional Programming* 7.5 (Sept. 1997), pages 549–554.

Bibliography

- [11] Dimitrios Theodoratos and Timos Sellis. “Data warehouse configuration”. In: *Proceedings of the 23rd International Conference on Very Large Databases, VLDB 1997*. 23rd International Conference on Very Large Databases, VLDB 1997. Morgan Kaufmann, Jan. 1, 1997, pages 126–135.
- [12] Jeffrey D. Ullman. “Information Integration Using Logical Views”. In: *International Conference on Database Theory*. Springer, 1997, pages 19–40.
- [13] Ralf Hinze. “Deriving Backtracking Monad Transformers”. In: *ACM SIGPLAN Notices* 35 (Aug. 15, 2000).
- [14] Prasan Roy et al. “Efficient and Extensible Algorithms for Multi Query Optimization”. In: *Proceedings of the 2000 ACM SIGMOD International Conference on Management of Data*. 2000, pages 249–260.
- [15] Hoshi Mistry et al. “Materialized View Selection and Maintenance Using Multi-Query Optimization”. In: *Proceedings of the 2001 ACM SIGMOD International Conference on Management of Data*. SIGMOD '01. New York, NY, USA: Association for Computing Machinery, May 1, 2001, pages 307–318.
- [16] Michael Stillger et al. “LEO-DB2's Learning Optimizer”. In: *VLDB*. Volume 1. 2001, pages 19–28.
- [17] John Hughes. “Programming with Arrows”. In: *Advanced Functional Programming*. Edited by Varmo Vene and Tarmo Uustalu. Lecture Notes in Computer Science. Berlin, Heidelberg: Springer, 2005, pages 73–129.
- [18] Tarmo Uustalu and Varmo Vene. *Comonadic Functional Attribute Evaluation*. Volume 6. Jan. 1, 2005, page 162. 145 pages.
- [19] Gang Gou, J. X. Yu, and Hongjun Lu. “A/Sup */ Search: An Efficient and Flexible Approach to Materialized View Selection”. In: *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)* 3.36 (2006), pages 411–425.
- [20] Jeremy E Dawson. “Compound Monads and the Kleisli Category”. In: *Unpublished note*. Available online at <http://users.cecs.anu.edu.au/~jeremy/pubs/cmkc> (2007).
- [21] Surajit Chaudhuri, Vivek Narasayya, and Ravi Ramamurthy. “A Pay-as-You-Go Framework for Query Execution Feedback”. In: *Proceedings of the VLDB Endowment* 1.1 (2008), pages 1141–1152.
- [22] Janis Voigtländer. “Asymptotic Improvement of Computations over Free Monads”. In: *Mathematics of Program Construction*. Edited by Philippe Audebaud and Christine Paulin-Mohring. Lecture Notes in Computer Science. Berlin, Heidelberg: Springer, 2008, pages 388–403.
- [23] Alexey Andreyevich Radul. “Propagation Networks: A Flexible and Expressive Substrate for Computation”. In: (2009).
- [24] Gerald Jay Sussman and Alexey Radul. “The Art of the Propagator”. In: (Jan. 26, 2009).

- [25] Brent Yorgey. “The Typeclassopedia”. In: (2009).
- [26] Pradeep Kumar Gunda et al. “Nectar: Automatic Management of Data and Computation in Datacenters.” In: *OSDI*. Volume 10. 2010, pages 1–8.
- [27] Herodotos Herodotou and Shivnath Babu. “Xplus: A Sql-Tuning-Aware Query Optimizer”. In: *Proceedings of the VLDB Endowment* 3.1-2 (2010), pages 1149–1160.
- [28] Milena G. Ivanova et al. “An Architecture for Recycling Intermediates in a Column-Store”. In: *ACM Transactions on Database Systems (TODS)* 35.4 (2010), pages 1–43.
- [29] Konstantinos Krikellas, Stratis D. Viglas, and Marcelo Cintra. “Generating Code for Holistic Query Evaluation”. In: *2010 IEEE 26th International Conference on Data Engineering (ICDE 2010)*. 2010 IEEE 26th International Conference on Data Engineering (ICDE 2010). Mar. 2010, pages 613–624.
- [30] Simon Marlow. “Haskell 2010 Language Report”. In: *Available on: <https://www.haskell.org/onlinereport/haskell2010>* (2010).
- [31] Tomasz Nykiel et al. “MRShare: Sharing across Multiple Queries in MapReduce”. In: *Proceedings of the VLDB Endowment* 3.1-2 (Sept. 1, 2010), pages 494–505.
- [32] Pramod Bhatotia et al. “Incoop: MapReduce for Incremental Computations”. In: *Proceedings of the 2nd ACM Symposium on Cloud Computing. SOCC '11*. New York, NY, USA: Association for Computing Machinery, Oct. 26, 2011, pages 1–14.
- [33] Swati V. Chande and Madhavi Sinha. “Genetic Optimization for the Join Ordering Problem of Database Queries”. In: *2011 Annual IEEE India Conference*. IEEE, 2011, pages 1–5.
- [34] Alfons Kemper and Thomas Neumann. “HyPer: A Hybrid OLTP&OLAP Main Memory Database System Based on Virtual Memory Snapshots”. In: *2011 IEEE 27th International Conference on Data Engineering*. IEEE, 2011, pages 195–206.
- [35] Iman Elghandour and Ashraf Aboulnaga. *ReStore: Reusing Results of MapReduce Jobs*. Comment: VLDB2012. Feb. 29, 2012. arXiv: [1203.0061](https://arxiv.org/abs/1203.0061) [cs]. URL: <http://arxiv.org/abs/1203.0061> (visited on 11/26/2021).
- [36] Imene Mami and Zohra Bellahsene. “A Survey of View Selection Methods”. In: *Acm Sigmod Record* 41.1 (2012), pages 20–29.
- [37] Srinath Shankar et al. “Query Optimization in Microsoft SQL Server PDW”. In: *Proceedings of the 2012 International Conference on Management of Data - SIGMOD '12*. The 2012 International Conference. Scottsdale, Arizona, USA: ACM Press, 2012, page 767.
- [38] Lindsey Kuper and Ryan R. Newton. “LVars: Lattice-Based Data Structures for Deterministic Parallelism”. In: *Proceedings of the 2nd ACM SIGPLAN Workshop on Functional High-Performance Computing. FHPC '13*. New York, NY, USA: Association for Computing Machinery, Sept. 23, 2013, pages 71–84.

- [39] Fabian Nagel, Peter Boncz, and Stratis D. Viglas. “Recycling in Pipelined Query Evaluation”. In: *2013 IEEE 29th International Conference on Data Engineering (ICDE)*. IEEE, 2013, pages 338–349.
- [40] Guoping Wang and Chee-Yong Chan. “Multi-Query Optimization in Mapreduce Framework”. In: *Proceedings of the VLDB Endowment 7.3* (2013), pages 145–156.
- [41] Matthew A. Hammer et al. “Adapton: Composable, Demand-Driven Incremental Computation”. In: *Proceedings of the 35th ACM SIGPLAN Conference on Programming Language Design and Implementation*. PLDI ’14: ACM SIGPLAN Conference on Programming Language Design and Implementation. Edinburgh United Kingdom: ACM, June 9, 2014, pages 156–166.
- [42] Luis L. Perez and Christopher M. Jermaine. “History-Aware Query Optimization with Materialized Intermediate Views”. In: *2014 IEEE 30th International Conference on Data Engineering*. 2014 IEEE 30th International Conference on Data Engineering. Mar. 2014, pages 520–531.
- [43] Luis L. Perez and Christopher M. Jermaine. “History-Aware Query Optimization with Materialized Intermediate Views”. In: *2014 IEEE 30th International Conference on Data Engineering*. IEEE, 2014, pages 520–531.
- [44] Mohamed A. Soliman et al. “Orca: A Modular Query Optimizer Architecture for Big Data”. In: *Proceedings of the 2014 ACM SIGMOD International Conference on Management of Data*. 2014, pages 337–348.
- [45] Mohamed A. Soliman et al. “Orca: A Modular Query Optimizer Architecture for Big Data”. In: *Proceedings of the 2014 ACM SIGMOD International Conference on Management of Data*. SIGMOD ’14. New York, NY, USA: Association for Computing Machinery, June 18, 2014, pages 337–348.
- [46] Michael Armbrust et al. “Spark SQL: Relational Data Processing in Spark”. In: *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data*. SIGMOD ’15. New York, NY, USA: Association for Computing Machinery, May 27, 2015, pages 1383–1394.
- [47] Melyssa Barata, Jorge Bernardino, and Pedro Furtado. “An Overview of Decision Support Benchmarks: TPC-DS, TPC-H and SSB”. In: *New Contributions in Information Systems and Technologies*. Edited by Alvaro Rocha et al. Advances in Intelligent Systems and Computing. Cham: Springer International Publishing, 2015, pages 619–628.
- [48] Viktor Leis et al. “How Good Are Query Optimizers, Really?” In: *Proceedings of the VLDB Endowment 9.3* (Nov. 2015), pages 204–215.
- [49] Jack Chen et al. “The MemSQL Query Optimizer: A Modern Optimizer for Real-Time Analytics in a Distributed Database”. In: *Proceedings of the VLDB Endowment 9.13* (Sept. 1, 2016), pages 1401–1412.

- [50] Wentao Wu, Jeffrey F. Naughton, and Harneet Singh. “Sampling-Based Query Re-Optimization”. In: *Proceedings of the 2016 International Conference on Management of Data*. 2016, pages 1721–1736.
- [51] Tarun Kathuria and S. Sudarshan. “Efficient and Provable Multi-Query Optimization”. In: *Proceedings of the 36th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems*. 2017, pages 53–67.
- [52] Prashanth Menon, Todd C. Mowry, and Andrew Pavlo. “Relaxed Operator Fusion for In-Memory Databases: Making Compilation, Vectorization, and Prefetching Work Together at Last”. In: *Proceedings of the VLDB Endowment* 11.1 (Sept. 1, 2017), pages 1–13.
- [53] William Arthur Rogers et al. “Multi-Query Optimization”. Patent. US Patent App. 15/222,229. Mar. 23, 2017.
- [54] E. Yu. Sharygin et al. “Query Compilation in PostgreSQL by Specialization of the DBMS Source Code”. In: *Programming and Computer Software* 43.6 (Nov. 1, 2017), pages 353–365.
- [55] Edmon Begoli et al. “Apache Calcite: A Foundational Framework for Optimized Query Processing over Heterogeneous Data Sources”. In: *Proceedings of the 2018 International Conference on Management of Data*. 2018, pages 221–230.
- [56] Bailu Ding et al. “Plan Stitch: Harnessing the Best of Many Plans”. In: *Proceedings of the VLDB Endowment* 11.10 (June 2018), pages 1123–1136.
- [57] Stavros Harizopoulos et al. “OLTP through the Looking Glass, and What We Found There”. In: *Making Databases Work: The Pragmatic Wisdom of Michael Stonebraker*. 2018, pages 409–439.
- [58] Sven Keidel and Sebastian Erdweg. “Sound and Reusable Components for Abstract Interpretation”. In: *Proceedings of the ACM on Programming Languages* 3 (OOPSLA Oct. 10, 2019), 176:1–176:28.
- [59] Ryan Marcus et al. “Neo: A Learned Query Optimizer”. In: *Proceedings of the VLDB Endowment* 12.11 (July 2019), pages 1705–1718.
- [60] Jennifer Ortiz et al. *An Empirical Analysis of Deep Learning for Cardinality Estimation*. Sept. 11, 2019. arXiv: [1905.06425 \[cs\]](https://arxiv.org/abs/1905.06425). URL: <http://arxiv.org/abs/1905.06425> (visited on 11/28/2021).
- [61] Mark Raasveldt and Hannes Mühleisen. “DuckDB: An Embeddable Analytical Database”. In: *Proceedings of the 2019 International Conference on Management of Data*. SIGMOD/PODS ’19: International Conference on Management of Data. Amsterdam Netherlands: ACM, June 25, 2019, pages 1981–1984.

Bibliography

- [62] Ana Nunes Alonso et al. “Building a Polyglot Data Access Layer for a Low-Code Application Development Platform: (Experience Report)”. In: *Distributed Applications and Interoperable Systems*. Edited by Anne Remke and Valerio Schiavoni. Volume 12135. Lecture Notes in Computer Science. Cham: Springer International Publishing, 2020, pages 95–103.
- [63] Daniel Selsam, Simon Hudon, and Leonardo deMoura. “Sealing Pointer-Based Optimizations behind Pure Functions”. In: *Proceedings of the ACM on Programming Languages* 4 (ICFP Aug. 2, 2020), 115:1–115:20.
- [64] Chris Hanson and Gerald Jay Sussman. *Software Design for Flexibility: How to Avoid Programming Yourself into a Corner*. MIT Press, Mar. 9, 2021. 449 pages. Google Books: [iM_tDwAAQBAJ](https://books.google.com/books?id=iM_tDwAAQBAJ).
- [65] Donnacha Oisín Kidney and Nicolas Wu. “Algebras for Weighted Search”. In: *Proceedings of the ACM on Programming Languages* 5 (ICFP Aug. 18, 2021), 72:1–72:30.
- [66] Guoliang Li, Xuanhe Zhou, and Lei Cao. “Machine Learning for Databases”. In: *The First International Conference on AI-ML-Systems*. AIMLSystems 2021. New York, NY, USA: Association for Computing Machinery, Oct. 21, 2021, pages 1–2.
- [67] Pietro Michiardi, Damiano Carra, and Sara Migliorini. “Cache-Based Multi-Query Optimization for Data-Intensive Scalable Computing Frameworks”. In: *Information Systems Frontiers* 23.1 (2021), pages 35–51.
- [68] 14:00-17:00. *ISO/IEC 14882:2020*. ISO. URL: <https://www.iso.org/cms/render/live/en/sites/isoorg/contents/data/standard/07/93/79358.html> (visited on 10/28/2021).
- [69] *Associated Type Synonyms | ACM SIGPLAN Notices*. URL: https://dl.acm.org/doi/abs/10.1145/1090189.1086397?casa_token=HLcvyM59laEAAAAA:3yH9-bJAMmzqia0eww_MReCDVnbfydL5nsIokMy4g8KWDC12huZV0_X7kMSVaE257aTivPkkuZL00A (visited on 10/10/2021).
- [70] *Implementing Data Cubes Efficiently | Proceedings of the 1996 ACM SIGMOD International Conference on Management of Data*. URL: <https://dlnext.acm.org/doi/10.1145/233269.233333> (visited on 11/26/2021).
- [71] Oleg Kiselyov et al. “Backtracking, Interleaving, and Terminating Monad Transformers”. In: (), page 12.
- [72] Edward Kmett. *Propagators*.
- [73] Tim Kraska et al. “SageDB: A Learned Database System”. In: (), page 10.
- [74] Daan Leijen. “Parsec: Direct Style Monadic Parser Combinators For The Real World”. In: (), page 22.
- [75] Thomas Neumann. *Evolution of a CompilingQuery Engine*. URL: <https://www.semanticscholar.org/paper/Evolution-of-a-CompilingQuery-Engine-Neumann/2cd16bb2d3680d04948f4536270c36d98fd4e768> (visited on 11/28/2021).

- [76] Luis L. Perez and Christopher M. Jermaine. *History-Aware Query Optimization with Materialized Intermediate Views*.
- [77] Chris Perivolaropoulos. *Fakedrake/Ssb-Dbgen*.
- [78] Chris Perivolaropoulos. *Fakedrake/Ssb-Dbgen*.
- [79] *Precompiled Headers (PCH)*. URL: <https://releases.llvm.org/3.1/tools/clang/docs/PCHInternals.html> (visited on 12/03/2021).

