

Push vs. Pull-Based Loop Fusion in Query Engines

Amir Shaikhha, Mohammad Dashti, and Christoph Koch
{firstname}.{lastname}@epfl.ch
École Polytechnique Fédérale de Lausanne

ABSTRACT

Database query engines use pull-based or push-based approaches to avoid the materialization of data across query operators. In this paper, we study these two types of query engines in depth and present the limitations and advantages of each engine. Similarly, the programming languages community has developed loop fusion techniques to remove intermediate collections in the context of collection programming. We draw parallels between the DB and PL communities by demonstrating the connection between pipelined query engines and loop fusion techniques. Based on this connection, we propose a new type of pull-based engine, inspired by a loop fusion technique, which combines the benefits of both approaches. Then we experimentally evaluate the various engines, in the context of query compilation, for the first time in a fair environment, eliminating the biasing impact of ancillary optimizations that have traditionally only been used with one of the approaches. We show that for realistic analytical workloads, there is no considerable advantage for either form of pipelined query engine, as opposed to what recent research suggests. Also, by using microbenchmarks we show that our proposed engine dominates the existing engines by combining the benefits of both.

1. INTRODUCTION

Database query engines successfully leverage the compositionality of relational algebra-style query plan languages. Query plans are compositions of operators that, at least conceptually, can be executed in sequence, one after the other. However, actually evaluating queries in this way leads to grossly suboptimal performance. Computing (“materialising”) the result of a first operator before passing it to a second operator can be very expensive, particularly if the intermediate result is large and needs to be pushed down the memory hierarchy. The same observation has been made by the programming languages and compilers community and has led to work on loop fusion and deforestation (the elimination of data structure construction and destruction for intermediate results).

Already relatively early on in the history of relational database systems, a solution to this problem has been proposed in the form of the Volcano Iterator model [19]. In this model, tuples are *pulled* up through a chain of operators that are linked by iterators that advance in lock-step. Intermediate results between operators are not accumulated, but tuples are produced on demand, by request by conceptually “later” operators.

More recently, an operator chaining model has been proposed that shares the advantage of avoiding materialisation of intermediate results but which reverses the control flow; tuples are *pushed* forward from the source relations to the operator producing the final result. Recent papers [42, 30] seem to suggest that this push-model

consistently leads to better query processing performance than the pull model, even though no direct, fair comparisons are provided.

One of the main contributions of this paper is to debunk this myth. As we show, if compared fairly, push and pull based engines have very similar performance, with individual strengths and weaknesses, and neither is a clear winner. Push engines have in essence only been considered in the context of query compilation, conflating the potential advantages of the push paradigm with those of code inlining. To compare them fairly, one has to decouple these aspects.

Figure 1 shows a performance comparison of these two engines for several TPC-H queries using 8 GBs of data in a fair scenario. There is no clear winner among these two engines. In the case of two queries (TPC-H queries 12 and 14), the pull engine is performing better than the push engine. However, in some cases, the push-based query engine is performing marginally better. The advantages and limitations of these engines are explained in more detail in Section 2.

In this paper, we present an in-depth study of the tradeoffs of the push versus the pull paradigm. Choosing among push and pull – or any reasonable alternative – is a fundamental decision which drives many decisions throughout the architecture of a query engine. Thus, one must understand the relevant properties and tradeoffs deeply, and should not bet on one’s ability to overcome the disadvantages of a choice by a hack later.

Furthermore, we illustrate how the same challenge and trade-off has been met and addressed by the PL community, and show a number of results that can be carried over from the lessons learned there. Specifically, we study how the PL community’s answer to the problem, *stream fusion* [11], can be adapted to the query processing scenario, and show how it combines the advantages of the pull and push approaches. Furthermore, we demonstrate how we can use ideas from the push approach to solve well-known limitations of stream fusion. As a result, we construct a query engine which combines the benefits of both push and pull approaches. In essence, this engine is a pull-based engine on a coarse level of granularity, however, on a finer level of granularity, it pushes the individual data tuples.

In summary, this paper makes the following contributions:

- We discuss pipelined query engines in Section 2. After presenting loop fusion for collection programming in Section 3, we show the connection between these two concepts in Section 3.3. Furthermore, we demonstrate the limitations associated to each approach.
- Based on this connection with loop fusion, we propose a new pipelined query engine in Section 4 inspired by the stream fusion [11] technique developed for collection programming

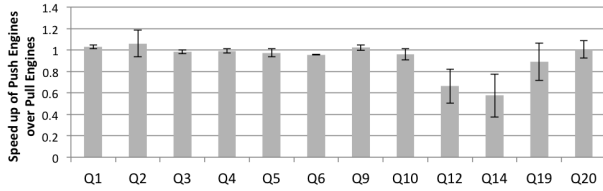


Figure 1: Performance comparison of push-based and pull-based query engines using TPC-H queries running 8 GBs of data.

in the PL community. Also, we discuss implementation concerns and compiler optimizations required for the proposed pipelined query engine in Section 5.

- We experimentally evaluate the various query engine architectures in Section 6. Using microbenchmarks, we discuss the weaknesses of the existing engines and how the proposed engine circumvents these weaknesses by combining the benefits of both worlds. Then we demonstrate using TPC-H queries that good implementations of these engines do not show a considerable advantage for either form of pipelined query engine.

Throughout this paper, we are using the Scala programming language for all code snippets, interfaces and examples. None of the concepts and ideas require specifically this language – other impure functional object-oriented programming languages such as OCaml, F#, C++11, C#, or Java 8 could be used instead.

2. PIPELINED QUERY ENGINES

Database management systems accept a declarative query (e.g., written in SQL). Such a query is passed to a query optimizer to find a fast physical query plan, which then is either interpreted by the query engine or compiled to low-level code (e.g. C code).

Physical query plans perform calculations and data transformations. A sequence of query operators can be *pipelined*, which means that the output of one operator is *streamed* into the next operator without materializing the intermediate data.

There are two approaches for pipelining. The first approach is demand-driven pipelining in which an operator repeatedly *pulls* the next data tuple from its *source* operator. The second approach is data-driven pipelining in which an operator *pushes* each data tuple to its *destination* operator. Next, we give more details on the pull-based and push-based query engines.

2.1 Pull Engine – a.k.a. the Iterator Pattern

The iterator model is the most widely used pipelining technique in query engines. This model was initially proposed in XRM [37]. However, the popularity of this model is due to its adoption in the Volcano system [19], in which this model was enriched with facilities for parallelization.

In a nutshell, in the iterator model, each operator pipelines the data by requesting the next element from its source operator. This way, instead of waiting until the whole intermediate relation is produced, the data is *lazily* generated in each operator. This is achieved by invoking the `next` method of the source operator by the destination operator. The design of pull-based engines directly corresponds to the iterator design pattern in object-oriented programming [54].

Figure 2 shows an example query and the control flow of query processing for this query. Each query operator performs the role of a destination operator and *requests* data from its source operator (the predecessor operator along the flow direction of data). In a

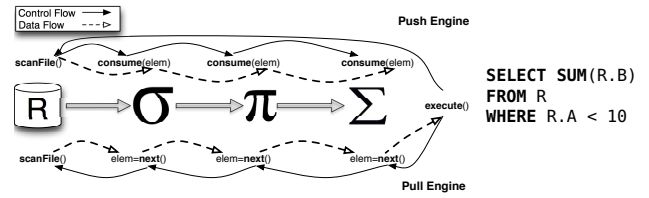


Figure 2: Data flow and control flow for push and pull-based query engine for the provided SQL query.

pull engine, this is achieved by invoking the `next` function of the source operator, and is shown as control flow edges. In addition, each operator serves as source operator and *generates* result data for its destination operator (the successor operator along the flow direction of data). The generated data is the return value of the `next` function, and is represented by the data flow edges in Figure 2. Note the opposing directions of control-flow and data-flow edges for the pull engine in Figure 2.

From a different point of view, each operator can be considered as a `while` loop in which the `next` function of the source operator is invoked per iteration. The loop is terminated when the `next` function returns a special value (e.g., a `null` value). In other words, whenever this special value is observed, a `break` statement is executed to terminate the loop execution.

There are two main issues with a pull-based query engine. First, the `next` function invocations are implemented as virtual functions – operators with different implementations of `next` have to be chained together. There are many invocations of these functions, and each of invocation requires looking up a virtual table, which leads to bad instruction locality. Query compilation solves this issue by inlining these virtual function calls, which is explained in Section 2.3.

Second, although a pull engine pipelines the data through pipelining operators, in practice, selection operators are problematic. When the `next` method of a selection operator is invoked, the destination operator should wait until the selection operator returns the next data tuple satisfying its predicate. This makes the control flow of the query engine more complicated by introducing more loops and branches, which is demonstrated in Figure 3c. This complicated control flow graph degrades branch prediction. Intuitively, this is because there is no construct for skipping the irrelevant results (such as the `continue` construct). This problem is solved in push-based query engines.

2.2 Push Engine – a.k.a. the Visitor Pattern

Push-based engines are widely used in streaming systems [24]. The Volcano system uses data-driven pipelining (which is a push-based approach) for implementing inter-operator parallelism in query engines. In the context of query compilation, stream processing engines such as StreamBase [1] and Spade [16], as well as HyPer [42] and LegoBase [30] use a push-based query engine approach.

In push-based query engines, the control flow is reversed compared to that of pull-based engines. More concretely, instead of destination operators requesting data from their source operators, data is pushed from the source operators towards the destination operators. This is achieved by the source operator passing the data as an argument to the `consume` method of the destination operator. This results in *eagerly* transferring the data tuple-by-tuple instead of requesting it *lazily* in pull-engines.

A push engine can be implemented using the *Visitor* design pattern [54] from object-oriented programming. This design pattern allows separating an algorithm from a particular type of data. In the

case of query engines, the visitor pattern allows us to separate the query operators (data processing algorithms) from a relation of elements. To do so, each operator should be defined as a visitor class, in which the consume method has the functionality of the visit method. The process of the initialization of the chain of operators is performed by using the accept method of the Visitor pattern, which corresponds to the produce method in push engines.

Figure 2 shows the query processing workflow for the given example query. Query processing in each operator consists of two main phases. In the first phase, operators prepare themselves for producing their data. This is performed only once in the initialization. In the second phase, they consume the data provided by the source operator and produce data for the destination operator. This is the main processing phase, which consists of invoking the consume method of the destination operator and passing the produced data through it. This results in the same direction for both control-flow and data-flow edges, as shown in Figure 2.

Push engines solve the problem pull engines have with selection operators. This is achieved by ignoring the produced data if it does not satisfy the given predicate by using a construct which allows to skip the current iteration of the loop (e.g., using `continue`). This simplifies the control flow and improves branch prediction in the case of selection operators. This is in contrast with pull-engines in which the destination operator should have waited for the source operator to serve the request.

However, push engines experience difficulties with *limit* and *merge join* operators. For limit operators, push engines do not allow terminating the iteration by nature. This is because, in push engines, the operators cannot control when the data should no longer be produced by their source operator. This causes the production of elements which will never be used.

The merge join operator suffers from a similar problem. There is no way for the merge join operator to guide which one of its two source operators (which are both sorted and there is a 1-to-n relationship between them) should produce the next data item. Hence, it is not possible to pipeline the data coming from both source operators in merge join. As a result, at least for one of the source operators, the pipeline needs to be broken. Hence, the incoming data coming from one of the source operators can be pipelined (assuming it is correctly sorted, of course), but the input data coming from the other source operator must be materialized.

The mentioned limitation is not specific to operators such as merge joins. A similar situation can happen in the case of more sophisticated analytical tasks where one has to use collection programming APIs (such as Spark RDDs [57]). The `zip` method in collection programming has a similar behavior to the merge join operator and cannot be easily pipelined in push-based engines.

Note that these limitations can be resolved by providing special cases for these two operators in the push engine. In the case of *limit*, one can avoid producing unnecessary elements by manually fusing this operator with its source operator (which in most cases is an ordering operator). Also, one can implement a variant of merge join which uses different threads for its source operators and uses synchronization constructs in order to control the production of data by its two inputs, which can be costly. However, in this paper, by push engine, we mean a *pure* push engine without such augmentations.

2.3 Compiled Engines

In general, the runtime cost of a given query is dependent on two factors. The first factor is the time it takes to transfer the data across storage and computing components. The second factor is the time taken for performing the actual computation (i.e., running the instructions of the query). In disk-based DBMSes, the domi-

```

1 var sum = 0.0
2 var index = 0
3 while(index < RSize) {
4   var rec = null
5   while(index < RSize) {
6     val elem = R(index)
7     index += 1
8     if(elem.A < 10) {
9       rec = elem
10      break
11    }
12  }
13  if(rec == null)
14    break
15  sum += rec.B
16 }
17 return sum

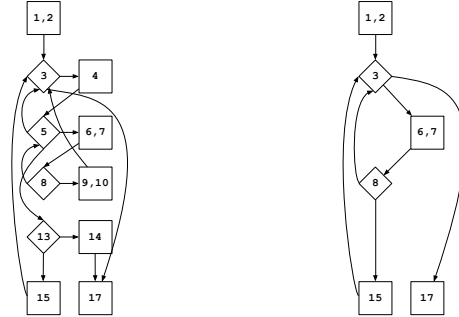
```

```

var sum = 0.0
var index = 0
while(index < RSize) {
  val rec = R(index)
  index += 1
  if(rec.A < 10)
    sum += rec.B
}
return sum

```

(a) Inlined query in pull engine. (b) Inlined query in push engine.



(c) The CFG of the inlined query in pull engine. (d) The CFG of the inlined query in push engine.

Figure 3: Specialized version of the example query in pull and push engines and the corresponding control-flow graphs (CFG).

nating cost is usually the data transfer from/to the secondary storage. Hence, as long as the pipelining algorithm does not break the pipeline, there is no difference between pull engines and push engines. As a result, the practical problem with selections in pull engines (c.f. Section 2.1) is obscured by data transfer costs.

With the advent of in-memory DBMSes, the code layout of the instructions becomes a very important factor. As a result, query compilation uses code generation and compilation techniques in order to inline virtual functions and further specialize the code to improve cache locality [20, 2, 32, 35, 42, 33, 34, 30, 53, 12, 41, 29, 3, 48, 28]. As a result of that, the code pattern used in each pipelining algorithm really matters. Hence, it is important to investigate the performance of each pipelining algorithm for different workloads.

Figure 3a shows the inlined pull-engine code for the example SQL query given in Figure 2. Note that for the selection operator, we need an additional `while` loop. This additional loop creates more branches in the generated code, which makes the control-flow graph (CFG) more complicated. Figure 3c demonstrates the CFG of the inlined pull-engine code. Each rectangle in this figure corresponds to a block of statements, whereas diamonds represent conditionals. The edges between these nodes represent the execution flow. The backward edges represent the jumps inside a loop. This complicated CFG makes the code harder to understand and optimize for the optimizing compiler. As a result, during the runtime execution, performance degrades mainly because of the worse branch prediction.

Figure 3b shows the specialized query for a push engine of the previous example SQL query. The selection operator here is summarized in a single `if` statement. As a result, the CFG for the

inlined push-engine code is simpler in comparison with the one for pull engine, as it is demonstrated in Figure 3d. This makes the reasoning and optimization easier for the underlying optimizing compiler, leading to better branch prediction during runtime execution.

Up to now, there is no separation of the concept of pipelining from the associated specializations. For example, HyPer [42] is in essence a push engine which uses compiler optimizations by default, without identifying the individual contributions to performance by these two factors. As another example, LegoBase [30] assumes that a push engine is followed by operator inlining, whereas the pull engine does not use operator inlining [31]. On the other hand, there is no comparison between an inlined pull engine – we suspect Hekaton [13] to be of that class – with a push-based inlined engine in the same environment. Hence, there is no comparison between pull and push engines which is under completely fair experimental conditions, sharing environment and code base to the maximum degree possible. In Section 6, we attempt such a fair comparison.

Furthermore, naïvely compiling the pull engine does not lead to good performance. This is because a naïve implementation of the iterator model does not take into account the number of `next` function calls. For example, the naïve implementation of the selection operator invokes the `next` method of its source operator twice, as it is demonstrated below:

```

1 class SelectOp[R] (p: R => Boolean) {
2   def next(): R = {
3     var elem: R = source.next()
4     while(elem != null && !p(elem)) {
5       elem = source.next()
6     }
7     elem
8   }
9 }

```

The first invocation is happening before the loop for the initialization (line 3), and the second invocation is inside the loop (line 5). Inlining can cause an explosion of the code size, which can lead to worse instruction cache behavior. Hence, it is important to take into account these concerns while implementing query engines. For example, our implementation of the selection operator in a pull-based query engine invokes the `next` method of its source operator only once by changing the shape of the `while` loop (c.f. Figure 5e). Section 6 shows the impact of this *inline-aware* implementation of pull engines.

3. LOOP FUSION IN COLLECTION PROGRAMMING

Collection programming APIs are getting more and more popular. Ferry [21, 20] and LINQ [39] use such an API to seamlessly integrate applications with database back-ends. Spark *RDDs* [57] use the same operations as collection programming APIs. Also, functional collection programming abstractions exist in main-stream programming languages such as Scala, Haskell, and recently Java 8. The theoretical foundation of such APIs is based on Monad Calculus and Monoid Comprehensions [7, 8, 56, 22, 52, 15].

Similar to query engines, the declarative nature of collection programming comes with a price. Each collection operation performs a computation on a collection and produces a transformed collection. A chain of these invocations results in creating unnecessary intermediate collections.

Loop fusion or Deforestation [55] removes the intermediate collections in collection programs. This transformation is a non-local and brittle transformation which is difficult to apply to impure functional programs (i.e., in languages which include imperative features) and is thus absent from mainstream compilers for such

languages. In order to provide a practical implementation, one can restrict the language to a pure functional DSL for which the fusion rules can be applied locally. These approaches are known as *short-cut* deforestation, which remove intermediate collections using local transformations instead of global transformations. This makes it more realistic for them to be integrated into real compilers; short-cut approaches have been successfully implemented in the context of Haskell [50, 11, 18] and Scala-based DSLs [27, 48].

Next, we present two approaches for short-cut deforestation in the order they were discovered. Both approaches employ two methods of “collection” micro-instructions each, to which a large number of collection operations can be mapped. This allows to implement fusion using very few rewrite rules (in terms of these micro-instructions).

3.1 Fold Fusion

In this approach, every collection operation is implemented using two constructs: 1) the `build` method for *producing* a collection, and 2) the `foldr` method for *consuming* a collection. Some methods such as `map`, which transform a collection, use both of these constructs for consuming the given collection and producing a new collection. However, some methods such as `sum`, which produce an aggregated result from a collection, require only the `foldr` method for consuming the given collection.

We consider an imperative variant of this algorithm, in which the `foldr` method is substituted by `foreach`. The main difference is that the `foldr` method explicitly handles the state, whereas in the case of `foreach`, the state is handled internally and is not exposed to the interface.

Using Scala syntax, the signature of the `foreach` method on lists is as follows:

```

class List[T] {
  def foreach(f: T => Unit): Unit
}

```

The `foreach` method consumes a collection by iterating over the elements of that collection and applying the given function to each element. The `build` function is the corresponding producer for the `foreach` method. This function produces a collection for which the `foreach` method applies the consumer higher-order function to the function `f`. The signature of the `build` function is as follows:

```

def build[T](consumer: (T => Unit) => Unit): List[T]

```

We illustrate the meanings of these two methods by an example. Consider the `map` method of a collection, which transforms a collection by applying a given function to each element. This method is expressed in the following way using the `build` and `foreach` functions:

```

class List[T] {
  def map[S](f: T => S): List[S] = build { consume =>
    this.foreach(e => consume(f(e)))
  }
}

```

The implementation of several other collection operators using these two methods is given in Figure 5b.

After rewriting the collection operations using the `build` and `foreach` constructs, a pipeline of collection operators involves constructing intermediate collections. These intermediate collections can be removed using the following rewrite rule:

Fold-Fusion Rule:

$$\text{build}(f1).\text{foreach}(f2) \rightsquigarrow f1(f2)$$

For example, there is a loop fusion rule for the `map` function, which fuses two consecutive `map` operations into one. More concretely, the expression `list.map(f).map(g)` is converted into

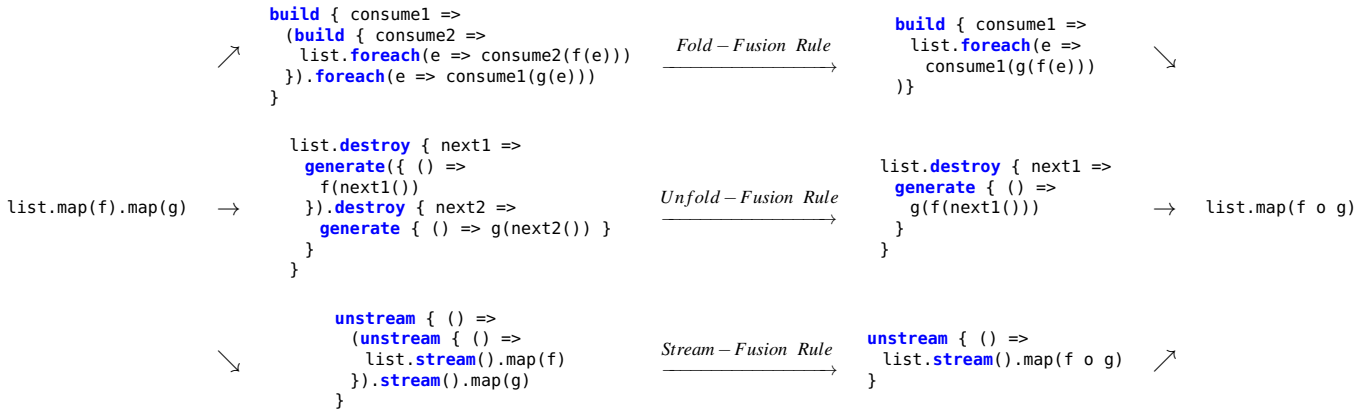


Figure 4: Different fusion techniques on a simple example.

`list.map(f o g)`. Figure 4 demonstrates how the fold-fusion technique can derive this conversion by expressing the map operator in terms of `foreach` and `build`, following by application of the fold-fusion rule.

One of the key advantages of this approach is that instead of writing fusion rewrite rules for every combination of collection operations, it is sufficient to only express these operations in terms of the `build` and `foreach` methods. This way, instead of writing $O(n^2)$ rewrite rules for n collection operations, it is sufficient to express these operations in terms of `build` and `foreach`, which is only $O(n)$ rewrite rules. Hence, this approach greatly simplifies the maintenance of the underlying compiler transformations [48].

This approach successfully deforests most collection operators very well. However, it is not successful in the case of `zip` and `take` operations. The `zip` method involves iterating over two collections, which cannot be expressed using the `foreach` construct which iterates only over one collection. Hence, this approach can deforest only one of the collections and for the other one, an intermediate collection must be created. Also, for the `take` method, there is no way to stop the iteration of the `foreach` method halfway to finish. Hence, the fold fusion technique does not perform well in these two cases. The next fusion technique solves the problem with these two methods.

3.2 Unfold Fusion

This is considered a dual approach to fold fusion. Every collection operation is expressed in terms of the two constructs `generate` and `destroy`, which have the following prototypes:

```

class List[T] {
  def destroy[S](f: (() => T) => S): S
}
def generate[T](next: () => T): List[T]
  
```

The `destroy` method consumes the given list. Each element of this collection is accessible by invoking the `next` function available by the `destroy` method. The `generate` function generates a collection, that its elements are specified by the input function passed to this method. In the case of `map` operator, the elements of the result collection are the elements of the input collection after the `f` function being applied to them.

The `map` method of collections is expressed in the following way using the `generate` and `destroy` methods:

```

class List[T] {
  def map[S](f: T => S): List[S] = this.destroy { next =>
  
```

¹We are presenting an imperative version of unfold fusion here; the purely functional version employs an `unfold` function instead of `generate`, and the approach derives its name from that.

Pipelined Query Engines	Object-Oriented Design Pattern	Collection Loop Fusion
Pull Engine	Iterator	Unfold fusion [50] Stream fusion [11]
Push Engine	Visitor	Fold fusion [18]

Table 2: Correspondence among pipelined query engines, object-oriented design patterns, and collection programming loop fusion.

```

generate { () =>
  val elem = next()
  if(elem == null) null
  else f(elem)
}
}
}
  
```

The implementation of some other collection operators using these two methods is given in Figure 5f.

In order to remove the intermediate collections, the chain of intermediate `generate` and `destroy` can be removed. This fact is shown in the following transformation rule:

Unfold-Fusion Rule:

`generate(f1).destroy(f2) ~> f2(f1)`

Figure 4 demonstrates how this rule fuses the previous example, `list.map(f).map(g)` into `list.map(f o g)`. Note that the null checking statements, which are for checking the end of a list, are removed for brevity.

This approach introduces a recursive iteration for the `filter` operation. In practice, this can cause performance issues, however from a theoretical point of view the deforestation is applied successfully [23]. Also, this approach does not fuse operations on nested collections, which is beyond the scope of this paper.

3.3 Loop Fusion is Operator Pipelining

By chaining query operators, one can express a given (say, SQL) query. Similarly, a given collection program can be expressed using a pipeline of collection operators. The relationship between relational queries and collection programs has been well studied. In particular, one can establish a precise correspondence between relational query plans and a class of collection programs [44].

Operators can be divided into three categories: 1) The operators responsible for *producing* a collection from a given source (e.g., a file or an array), 2) The operators which *transform* the given collection to another collection, and 3) The *consumer* operators which aggregate the given collection into a single result.

The mapping between query operators and collection operators is summarized in Table 1. Most join operators do not have a directly

Operator Category	Producer	Transformer						Consumer
Query Operator	Scan	Selection	Projection	OrderBy	Limit	Join*	Merge Join [†]	Agg [‡]
Collection Operator	List.fromArray	filter	map	sortBy	take	flatMap*	zip [†]	fold [‡]

* Nested loop join can be expressed using two nested flatMaps, but there is no equivalent for hash-based joins. Also, flatMaps can express nested collections, whereas in relational query engines every relation is considered to be flat.

[†] Merge join and zip have similar behavior, but their semantic is different.

[‡] An Agg operator representing a Group By is a transformer, whereas the one which folds into only a single result is a consumer.

Table 1: Mapping between query operators and collection operators

corresponding collection operator with two exceptions: Nested loop joins can be expressed using nested flatMaps and the zip collection operator is very similar to the merge join query operator. Both operators need to traverse two input sequences in parallel. For the rest of join operators, we extend collection programming with join operators (e.g. hashJoin, semiHashJoin, etc.). A similar mapping between the LINQ [39] operators and Haskell lists is shown in Steno [40]. Note that we do not consider nested collections here, although straightforward to support in collection programming, in order to emphasize similarity with relational query engines.

Pipelining in query engines is analogous to loop fusion in collection programming. Both concepts remove the intermediate relations and collections, which break the stream pipeline. Also, pipelining in query engines matches well-known design patterns in object-oriented programming [54]. The correspondence among pipelining in query engines, design patterns in object-oriented languages, and loop fusion in collection programming is summarized in Table 2.

Push Engine = Fold Fusion. There is a similarity between the Visitor pattern and fold fusion. On one hand it has been proven that the Visitor design pattern corresponds to the Church-encoding [6] of data types [9]. On the other hand, the foldr function on a list corresponds to the Church-encoding of lists in λ -calculus [45, 49]. Hence, both approaches eliminate intermediate results by converting the underlying data structure into its Church-encoding. In the former case, specialization consists of inlining, which results in removing (virtual) function calls. In the latter case, the fold-fusion rule and β -reduction are performed to remove the materialization points and inline the λ expressions. The correspondence between these two approaches is shown in Figure 5 (compare (a) vs. (b)). The invocations of the consume method of the destination operators in the push engine corresponds to the invocation of the consume function which is passed to the build operator in fold fusion.

Pull Engine = Unfold Fusion. In a similar sense, the Iterator pattern is similar to unfold fusion. Although the category-theoretic essence of the iterator model was studied before [17], there is no literature on the direct correspondence between the unfold function and the Iterator pattern. However, Figure 5 shows how a pull engine is similar to unfold fusion (compare Figure 5 (e) vs. (f)), to the best of our knowledge for the first time. Note the correspondence between the invocation of the next function of the source operator in pull engines, and the invocation of the next function which is passed to the destroy operator in unfold fusion, which is highlighted in the figure.

4. AN IMPROVED PULL-BASED ENGINE

In this section, we first present yet another loop-fusion technique for collection programs. Then, we suggest a new pull-based query engine inspired by this fusion technique based on the correspondence between queries and collection programming.

4.1 Stream Fusion

In functional languages, loops are expressed as recursive functions. Reasoning about recursive functions is very hard for optimizing compilers. Stream fusion tries to solve this issue by converting all recursive collection operations to non-recursive stream operations. To do so, first all collections are converted to streams using the stream method. Then, the corresponding method on the stream is invoked which results in a transformed stream. Finally, the transformed stream is converted back to a collection by invoking the unstream method.

The signature of the unstream and stream methods is as follows:

```
def unstream[T](next: () => Step[T]): List[T]
class List[T] {
  def stream(): Step[T]
}
```

For example, the map method is expressed in using these two methods as:

```
class List[T] {
  def map[S](f: T => S): List[S] = unstream { () =>
    this.stream().map(f)
  }
}
```

The stream method converts the input collection to an intermediate stream, which is specified by the Step data type. The function f is applied to this intermediate stream using the map function of the Step data type. Afterwards, the result stream is converted back to a collection by the unstream method.

As discussed before, one of the main advantages of the intermediate stream, the Step data structure, is that its operations are mainly non-recursive. This simplifies the task of the optimizing compiler to further specialize the program. The implementation of several methods of the Step data structure is given in Figure 6c.

Such transformations do not result in direct performance gain – they may even degrade performance. This is because of the intermediate conversions between streams and collections. However, these intermediate conversions can be removed using the following rewrite rule:

Stream-Fusion Rule:

$$\text{unstream}(() \Rightarrow e).\text{stream}() \rightsquigarrow e$$

Figure 4 demonstrates how the stream fusion technique transforms list.map(f).map(g) into list.map(f o g). Note that for the Step data type, the step.map(f).map(g) expression is equivalent to step.map(f o g).

The idea behind stream fusion is very similar to unfold fusion. The main difference is the filter operator. Stream fusion uses a specific value, called Skip, to implement the filter operator. This is in contrast with the unfold fusion approach for which the filter operator is implemented using an additional nested while loop for skipping the unnecessary elements. Hence, stream fusion solves the practical problem of unfold fusion associated with the filter operator.

Next, we define a new pipelined query engine based on the ideas of stream fusion.

```

class ProjectOp[R, P](f: R => P) {
  def consume(e: R): Unit =
    dest.consume(f(e))
}
class SelectOp[R](p: R => Boolean) {
  def consume(e: R): Unit =
    if(p(e))
      dest.consume(e)
}
class AggOp[R, S](f: (R, S) => S) {
  var result = zero[S]
  def consume(e: R): Unit = {
    result = f(e, result)
  }
  def getResult: S = result
}
class HashJoinOp[R, R2]
(leftHash: R => Int)
(rightHash: R2 => Int)
(cond: (R, R2) => Boolean) {
  val hm = new MultiMap[Int, R]()
  def consumeLeft(e: R): Unit = {
    hm.addBinding(leftHash(e) -> e)
  }
  def consumeRight(e: R2): Unit = {
    hm.get(rightHash(e)) match {
      case Some(list) =>
        for(l <- list) {
          if(cond(l, e)) {
            dest.consume(l.concat(e))
          }
        }
      case None =>
    }
  }
}
}
}

```

(a) Push-based query engine

```

class QueryMonad[R] {
  def map[S](f: R => S) = build { consume =>
    for(e <- this)
      consume(f(e))
  }
  def filter(p: R => Boolean) = build { consume =>
    for(e <- this)
      if(p(e))
        consume(e)
  }
  def fold[S](zero: S)(f: (R, S) => S): S = {
    var result = zero
    for(e <- this) {
      result = f(e, result)
    }
    result
  }
  def hashJoin[R2](rightList: QueryMonad[R2])
(leftHash: R => Int)
(rightHash: R2 => Int)
(cond: (R, R2) => Boolean) = build { consume =>
  val hm = new MultiMap[Int, R]()
  for(e <- this) {
    hm.addBinding(leftHash(e) -> e)
  }
  for(e <- rightList) {
    hm.get(rightHash(e)) match {
      case Some(list) =>
        for(l <- list) {
          if(cond(l, e)) {
            consume(l.concat(e))
          }
        }
      case None =>
    }
  }
}
}
}

```

(b) Fold fusion of collections.

```

type Cont[T] = (T => Unit) => Unit
class Consumer[T]
(val cont: Cont[T])
extends QueryMonad[T] {
  def foreach(f: T => Unit): Unit =
    cont(f)
}
def build[T](cont: Cont[T])
: QueryMonad[T] =
  new Consumer[T](cont)

```

(c) The constructs for fold fusion.

```

build(f1).foreach(f2)
↓ (inline build definition)
new Consumer(f1).foreach(f2)
↓ (inline foreach definition)
f1(f2)

```

(d) The derivation of the fold-fusion rule.

```

class SelectOp[R](p: R => Boolean) {
  def next(): R = {
    var elem: R = null
    do {
      elem = source.next()
    } while (elem != null && !p(elem))
    elem
  }
}
class ProjectOp[R, P](f: R => P) {
  def next(): P = {
    val elem = source.next()
    if(elem == null) null
    else f(elem)
  }
}
class AggOp[R, S]
(f: (R, S) => S) {
  def next(): S = {
    var result = zero[S]
    var elem: R = source.next()
    while(elem != null){
      result = f(elem, result)
      elem = source.next()
    }
    result
  }
}
class LimitOp[R](n: Int) {
  var count = 0
  def next(): R = {
    if(count < n) {
      count += 1
      source.next()
    } else {
      null
    }
  }
}
}
}

```

(e) Pull-based query engine

```

class QueryMonad[R] {
  def filter(p: R => Boolean) = destroy { next =>
    generate { () =>
      var elem: R = null
      do {
        elem = next()
      } while(elem != null && !p(elem))
      elem
    }
  }
  def map[P](p: R => P) = destroy { next =>
    generate { () =>
      val elem = next()
      if(elem == null) null
      else f(elem)
    }
  }
  def fold[S](zero: S)
(f: (R, S) => S): S =
  destroy { next =>
    var result = zero
    var elem: R = next()
    while(elem != null){
      result = f(elem, result)
      elem = next()
    }
    result
  }
  def take(n: Int) = {
    var count = 0
    destroy { next =>
      if(count < limit) {
        count += 1
        next()
      } else {
        null
      }
    }
  }
}
}
}

```

(f) Unfold fusion of collections.

```

type Gen[T] = () => T
type Dest[T, S] = (Gen[T] => S) => S
class DestroyGen[T]
(val next: Gen[T])
extends QueryMonad[T] {
  def destroy[S](f: Gen[T] => S): S =
    f(next)
}
def generate[T](next: Gen[T])
: QueryMonad[T] =
  new DestroyGen[T](next)

```

(g) The constructs for unfold fusion.

```

generate(f1).destroy(f2)
↓ (inline generate definition)
new DestroyGen(f1).destroy(f2)
↓ (inline destroy definition)
f2(f1)

```

(h) The derivation of the unfold-fusion rule.

Figure 5: Comparison of pull-based and push-based pipelining and loop fusion algorithms; code snippets in Scala.

```

class SelectOp[R](p: R => Boolean) {
  def stream(): Step[R] = {
    source.stream().filter(p)
  }
}
class ProjectOp[R, P](f: R => P) {
  def stream(): Step[P] = {
    source.stream().map(f)
  }
}
class AggOp[R, S](f: (R, S) => S) {
  def stream(): Step[S] = {
    var result = zero[S]
    var done = false
    while(!done){
      source.stream().fold(
        e => { result = f(e, result) },
        () => ,
        () => { done = true }
      )
    }
    return result
  }
}
class LimitOp[R](n: Int) {
  var count = 0
  def stream(): Step[R] = {
    if(count < n) {
      source.stream().map(e => {
        count += 1
        e
      })
    } else {
      Done
    }
  }
}
} } }

class QueryMonad[R] {
  def filter(p: R => Boolean) = {
    unstream { () =>
      stream().filter(p)
    }
  }
  def map[P](f: R => P) = {
    unstream { () =>
      stream().map(f)
    }
  }
  def fold[S](z: S)(f: (R, S) => S): S = {
    unstream { () =>
      var result = zero[S]
      var done = false
      while(!done){
        stream().fold(
          e => { result = f(e, result) },
          () => ,
          () => { done = true }
        )
      }
      return result
    }
  }
  def take(n: Int) = {
    var count = 0
    unstream { () =>
      if(count < n) {
        stream().map(e => {
          count += 1
          e
        })
      } else {
        Done
      }
    }
  }
} } }

trait Step[T] {
  def filter(p: T => Boolean): Step[T]
  def map[S](f: T => S): Step[S]
  def fold[S](yld: T => S,
    skip: () => S, done: () => S): S
}

case class Yield[T](e: T) extends Step[T] {
  def filter(p: T => Boolean) =
    if(p(e)) Yield(e) else Skip
  def map[S](f: T => S) =
    Yield(f(e))
  def fold[S](yld: T => S,
    skip: () => S, done: () => S): S =
    yld(e)
}

case object Skip extends Step[Nothing] {
  def filter(p: Nothing => Boolean) =
    Skip
  def map[S](f: Nothing => S) =
    Skip
  def fold[S](yld: Nothing => S,
    skip: () => S, done: () => S): S =
    skip()
}

case object Done extends Step[Nothing] {
  def filter(p: Nothing => Boolean) =
    Done
  def map[S](f: Nothing => S) =
    Done
  def fold[S](yld: Nothing => S,
    skip: () => S, done: () => S): S =
    done()
}

```

(a) Stream-Fusion Query Engine (b) Stream fusion of collections. (c) The operations of the Step data type.

Figure 6: Stream-based query engine and the stream fusion technique.

4.2 Stream-Fusion Engine

The proposed query engine follows the same design as the iterator model. Hence, this engine is also a pull engine. However, instead of invoking the `next` method, this engine invokes the `stream` method, which returns a wrapper object of type `Step`. We refer to our proposed engine as the *stream-fusion engine*.

As we mentioned in Section 2.1, one of the main practical problems with a pull engine is the case of the selection operator. In this case, an operator waits until the selection operator returns the next satisfying element. The proposed engine solves this issue by using the `Skip` object which specifies that the current element should be ignored. Hence, selection operators are no longer a blocker for their destination operator.

The correspondence between the stream fusion algorithm and the stream-fusion engine is shown in Figure 6. Every query operator provides an appropriate implementation for the `stream` method, which invokes the `stream` method of the source operator to request the next element. Similarly, stream fusion uses the `stream` method to fetch the next element. Then, by invoking the `unstream` method the generated stream is converted back to a collection.

From a different point of view, a push engine can be expressed using a `while` loop and a construct for skipping to the next iteration (e.g. `continue`). By nature, it is impossible for a push-based engine to finish the iteration before the producer’s `while` loop finishes its job. In contrast, a pull engine is generally expressible using a `while` loop and a construct for terminating the execution of the `while` loop (e.g. `break`). This is because of the demand-driven nature of pull engines. However, in a pull-based engine there is no way to skip an iteration. As a result, this should be expressed us-

ing a nested `while` loop which results in performance issues (c.f. Section 2).

The stream-fusion engine solves the mentioned problem by adding a `Skip` construct which results in skipping to the next iteration. This has an equivalent effect to the `continue` construct. Table 3 summarizes the differences among the aforementioned query engines.

Consider a relation of two elements for which we select its first element and the second element is filtered out. The first call to the `stream` method of the selection operator in the stream-fusion engine, produces a `Yield` element, which contains the first element of the relation. The second invocation of the same method returns a `Skip` element, specifying that this element, which is the second element of the relation, is filtered out and should be ignored. The next invocation of this method, results in a `Done` element, denoting that there is no more element to be produced by the selection operator. The `Done` value has the same role as the `null` value in the pull engine.

The specialized version of the example query (which was introduced in Figure 2) based on the stream-fusion engine is shown in Figure 7a. The code is as compact as the push engine code. However, it suffers from some performance problems due to the intermediate `Step` objects created. The next section discusses implementation aspects and the optimizations needed for tuning the performance of the stream-fusion engine.

5. IMPLEMENTATION

In this section, we discuss the implementation of the presented query engines. First, we discuss how the fusion rules are implemented for each approach. Then, we show how the problem asso-


```

var index = 0
var sum = 0.0
while(true) {
  val step1 =
    if(index < RSize) {
      val rec = R(index)
      index += 1
      Yield(rec)
    } else
      Done
  step1.filter(x => x.A < 10)
  .map(x => x.B)
  .fold(x => sum += x,
        () => ,
        () => break)
}
return sum

```

```

var index = 0
var sum = 0.0
while(true) {
  if(index < RSize) {
    val rec = R(index)
    index += 1
    if(rec.A < 10)
      sum += rec.B
  }
  else
    break
}
return sum

```

(a) Inlined query in stream-fusion engine without further specializations. (b) Inlined query in stream-fusion engine by inlining the visitor model of Step.

Figure 7: Specialization of the example query in stream-fusion engine.

Pipelined Query Engines	Looping Constructs
Push Engine	while + continue
Pull Engine	while + break
Stream-Fusion Engine	while + break + continue

Table 3: The supported looping constructs by each pipelined query engine.

ciated with intermediate objects is resolved for the stream-fusion engine.

We used the open-source DBLAB framework [48]² to implement different query engines and the associated optimizations. This framework allows us to implement these engines in the high-level programming language Scala. The input programs can either be expressed using physical (relational algebra-style) query plans in the QPlan language or collection programming using the QMonad language. Furthermore, we implement the optimizations using rewrite rules which are provided by the transformation framework of DBLAB.

We implemented the collection programming operations and the corresponding loop fusion techniques. Due to the equivalence which was shown in Section 3.3 between query engines and collection programming, it is clear how they can be implemented for query engines. As a result, the experimental results presented in the next section for different fusion techniques matches the results for different approaches for pipelined query engines. Next, we discuss how the fusion rules for different loop fusion algorithms can be expressed in this framework.

5.1 Fusion By Inlining

As mentioned in Section 3.3, in loop fusion techniques, the fusion rule is expressed as a local transformation rule which is applied as an extension to the host language compiler (which is GHC [26] in the case of the mentioned papers). In this section, we show how these fusion rules are implemented by only using inlining. This was proposed for implementing fold fusion in Scala [27]. Here, we use a similar approach for other fusion techniques.

Figure 5c shows the definition of the `build` operator. By inlining the definition of this operator, an object of type `QueryMonad` is created. The `foreach` method of this object applies the higher-order function passed to the `build` method (`f1`) to the input parameter of the `foreach` method (`f2`). By inlining this `foreach` method, we derive the same rule as the fold-fusion rule which was introduced

²<http://github.com/epfldata/dblab>

```

type GenStream[T] = () => Step[T]

class Streamer[T]
  (val next: GenStream[T])
  extends QueryMonad[T] {
  def stream(): Step[T] =
    next()
}

def unstream[T](next: GenStream[T])
  : QueryMonad[T] =
  new Streamer[T](next)

```

Figure 8: The constructs for stream fusion.

```

unstream(() => e).stream()
  ↓ (inline unstream definition)
new Streamer(() => e).stream()
  ↓ (inline stream definition)
e

```

Figure 9: The derivation of the stream-fusion rule.

in Section 3. This derivation is shown in Figure 5c. The constructs and derivation of unfold fusion are shown in Figure 5g and Figure 5h. Stream fusion follows a similar pattern which is given in Figure 8 and Figure 9.

Next, we discuss the problematic creation of intermediate objects by the stream-fusion engine, as well as our solution.

5.2 Removing Intermediate Results

Although the stream-fusion engine removes intermediate relations, it creates intermediate `Step` objects. There are two problems with these intermediate objects. First, the `Step` data type operations are virtual calls. This causes poor cache locality and degrades the performance. Second, normally these intermediate objects lead to heap allocations. This causes more memory consumption and a worse runtime. This is why the original stream fusion approach is dependent on optimizations provided by its source language compiler (i.e., the GHC [26] compiler). Implementing an effective version of it for other languages requires supporting similar optimizations supported by the GHC compiler.

The first problem with virtual calls can be solved by rewriting the `Step` operations by enumerating all cases for the `Step` object. This is possible because there are only three possible concrete cases (1. `Yield` 2. `Skip` 3. `Done`) for this data type. To do so, one can use `if`-statements. In functional languages, the pattern matching feature can be used. Although this approach solves the first problem, still there are heap allocations which are not removed.

The good news is that these heap allocations can be converted to stack allocations. This is because the created objects are not escaping their usage scope. For example, these objects are not copied into an array and not used as an argument to a function. This fact can be verified by the well-known compilation technique of *escape analysis* [10]. Based on that, the heap allocations can be converted to stack allocations.

The compiler optimizations can go further and remove the stack allocations as well. Instead of the stack allocation for creating a `Step` object, the fields necessary to encode this type are converted to local variables. Hence the `Step` abstraction is completely removed. This optimization is known as *scalar replacement* in compilers.

```

trait StepVisitor[T] {
  def yld(e: T): Unit
  def skip(): Unit
  def done(): Unit
}

trait Step[T] { self =>
  def __match(v: StepVisitor[T]): Unit
  def filter(p: T => Boolean): Step[T] =
    new Step[T] {
      def __match(v: StepVisitor[T]): Unit =
        self.__match(new StepVisitor[T] {
          def yld(e: T): Unit =
            if (p(e)) v.yld(e) else v.skip()
          def skip(): Unit = v.skip()
          def done(): Unit = v.done()
        })
    }
  def map[S](f: T => S): Step[S] =
    new Step[S] {
      def __match(v: StepVisitor[S]): Unit =
        self.__match(new StepVisitor[T] {
          def yld(e: T): Unit = v.yld(f(e))
          def skip(): Unit = v.skip()
          def done(): Unit = v.done()
        })
    }
}

case class Yield[T](e: T) extends Step[T] {
  def __match(v: StepVisitor[T]): Unit = v.yld(e)
}
case object Skip extends Step[Nothing] {
  def __match(v: StepVisitor[Nothing]): Unit = v.skip()
}
case object Done extends Step[Nothing] {
  def __match(v: StepVisitor[Nothing]): Unit = v.done()
}

```

Figure 10: Step data type implemented using the Visitor pattern.

From a different point of view, removing the intermediate `Step` objects is a similar problem to removing the intermediate relations and collections in query engines and collection programming. Hence, one can borrow similar ideas and apply it for the `Step` objects in a fine-grained granularity.

To do so, we implement a variant of the `Step` data type using the Visitor pattern. As we discussed in Section 3.3, this is similar to the Church-encoding of data types. This encoding results in *pushing* `Step` objects down the pipeline. Hence, the stream-fusion engine implements a pull engine on a coarse-grained level and pushes the tuples on a fine-grained level. The Visitor pattern version of the `Step` data type is shown in Figure 10.

The result of applying this enhancement to our working example is shown in Figure 7b. By comparing this code to the code produced by a push engine, we see a clear similarity. First, there are no more additional virtual calls associated with the `Step` operators. Second, there is no more materialization of the intermediate `Step` objects. Finally, similar to push engines, the produced code does not contain any additional nested `while` loop for selection, hence it is easier to understand and optimize by an underlying compiler.

As an alternative implementation, one can implement the `Step` data type as a *sum* type, a type with different distinct cases in which an object can be one and only one of those cases. Hence, the implementation of the `Step` methods can use the pattern matching feature of the Scala programming language. However, it has been proven that the Visitor pattern is a way to encode the sum types in object-oriented languages [9]. On the other hand, pattern matching in Scala is a way to express the Visitor pattern [14]. Hence, from a conceptual point of view there is no difference between these implementations [25].

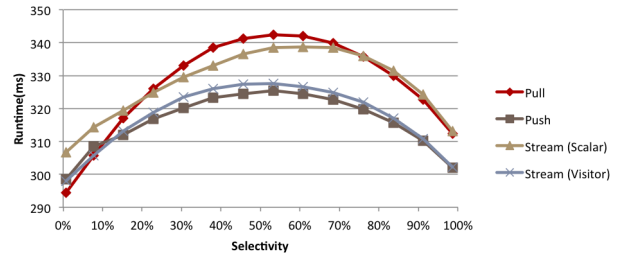


Figure 11: Sensitivity of query engines to selectivity.

6. EXPERIMENTAL RESULTS

For the experimental evaluation, we use a server-type x86 machine equipped with two Intel Xeon E5-2620 v2 CPUs running at 2GHz each, 256GB of DDR3 RAM at 1600Mhz and two commodity HDDs of 2TB. The operating system is Red Hat Enterprise 6.7.

Our query compiler uses the same set of transformations for different pipelining techniques to have a fair comparison. These transformations consist of *dead code elimination* (DCE), *common subexpression elimination* (CSE) or *global value numbering* (GVN), and *partial evaluation* (inlining and constant propagation). These transformations are provided out of the box by DBLAB [48], which we use as our testbed. Also, the scalar replacement transformation is always applied unless it is clearly specified. We do not use any data-structure specialization transformations or inverted indices for these experiments. Finally, all experiments use DBLAB's in-memory row-store representation.

For compiling the generated programs throughout our evaluation we use version 2.9 of the CLang compiler. We use the most aggressive optimization strategy provided by the CLang compiler (the "-O3" optimization flag)³. Finally, for C data structures we use the GLib library (version 2.42.1).

Our evaluation consists of two parts. First, by using micro benchmarks we demonstrate better the differences between different query engines. Then, for more complex queries we use the TPC-H [51] benchmark. We demonstrate how the different query engines behave in more complicated scenarios.

6.1 Micro Benchmarks

The micro benchmarks belong to three categories: First, queries consisting of only selection and aggregation without group by leading to a single result. Second, queries consisting of selection, projection, sort, and limit operations, which will return a list of results. Finally, queries with selection and different join operators, such as hash join, merge join, and semi hash join, which are followed by an aggregation operator resulting to a single result. All these queries use generated TPC-H databases at scaling factor 8, unless specified otherwise. The corresponding SQL queries for all these queries are shown in Table 4.

Sensitivity to Selectivity. The behavior of different engines for a simple query with one selection operator followed by an aggregation for different selectivities are shown in Figure 11. For highly selective queries, the Volcano pull engine is behaving better. This is because the unnecessary elements are skipped faster in the inner tight loop, whereas in the other engines the outer loop is respon-

³Although with this optimization flag it takes more time for the optimizing compiler to compile the queries, we observed similar results with the "-O1" optimization flag. This optimization flag provides all the transformation passes used in HyPer [42] except global value numbering (GVN). This transformation is not needed in our case, as it is already provided by DBLAB [48].

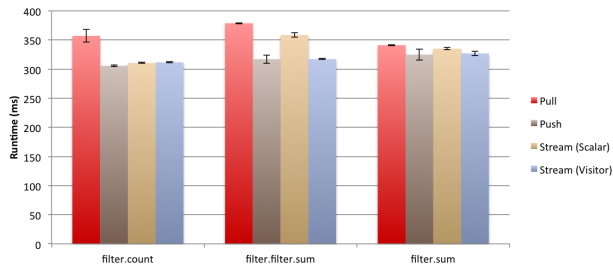


Figure 12: Simple queries with aggregated result.

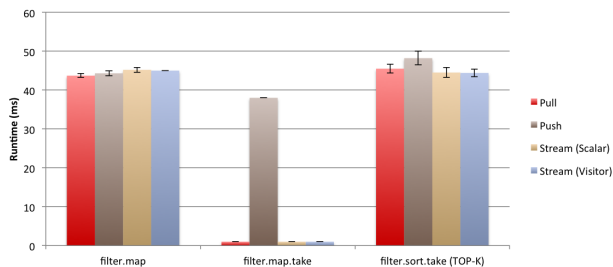


Figure 13: Simple queries with a list of results.

sible for skipping them. A similar effect was shown in [43] in the context of push engines and vectorized engines. For higher selectivities, in most cases, the push engine performs best. The Visitor-based stream-fusion engine offers almost the same performance as push engine, whereas the scalar-based stream-fusion engine (the one which only uses the scalar replacement transformation and not the Visitor pattern) is worse than other engines in most cases.

Aggregated Single Pipeline. Figure 12 demonstrates the performance obtained by each engine for queries with a single pipeline which produce a single result by summing over one column. The push engine is behaving slightly better than the pull engine in the presence of only a single filter operation. However, the stream-fusion engines hide this limitation of pull engines.

The difference is more obvious whenever there are chains of selection operations. A similar effect was shown in HyPer [42] in the case of using up to four consecutive selection operations. Again the Visitor-based stream-fusion engine is resolving this practical limitation of pull engines. From a practical point of view, as the query optimizer is merging all conjunctive predicates into a single selection operator, this case never happens in practice.

Single Pipeline. Figure 13 shows the results for single pipeline queries which do not contain any aggregation, hence producing a list of elements. For this experiment, we use 1GB of generated data. In the first query, in which a selection is followed by a projection operator, all engines behave similarly. However, in the third query, which returns top-k elements after filtering unnecessary elements, the push engine is performing worse than pull engines. This is because the push engine breaks the pipeline when using the limit query operator (c.f. Section 2.2). This situation is more obvious in the second query which consists of a selection, a projection, and a limit operator. In this case the pull engines do not require traversing all the elements and can stop immediately after reaching the limit. However, the push engine should wait until all elements are produced to be able to finish the execution. A similar phenomenon has been observed for pull-based and push-based fusion techniques for Java 8 streaming API in [4].

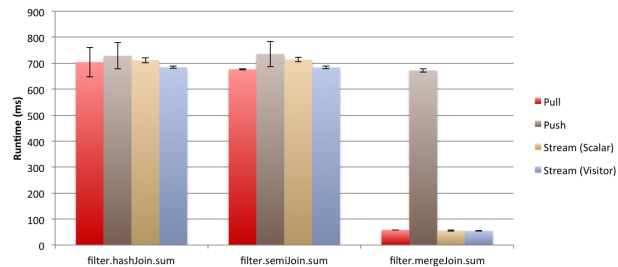


Figure 14: Join queries with aggregated result.

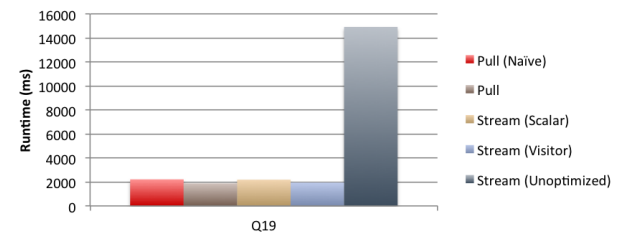


Figure 15: The performance comparison of several variants of pull-based engines on TPC-H query 19.

Aggregated Single Join. Finally, we investigate the performance of different join operations, which is demonstrated in Figure 14. In the case of hash join and semi hash join operators, there is no obvious difference among the engines. However, in the case of merge join, there is a great advantage for pull engines in comparison with the push engine. This is mainly because the push engine cannot pipeline both inputs of a merge join. Hence, it is forced to break the pipeline in one of the inputs (c.f. Section 2.2)⁴.

6.2 Macro Benchmarks

In this section, we investigate scenarios which are happening more often in practice. To do so, we use the larger and more complicated analytical queries defined in the TPC-H benchmark. First, we show the impact of fine-grained optimizations as well as our inline-aware way of implementing pull engines on one of TPC-H queries. Then, we investigate the impact of different engines on 12 TPC-H queries. All these experiments use 8 GBs of TPC-H generated data.

Inline-Aware Pull Engine Implementation. A naïve implementation of the selection operator in a pull-based query engine, invokes the next method of its source operator twice. This can exponentially grow the code size in the case of chain of selection operators. This case is not frequent in practice, since the selection operator is mainly used rightly after the scan operator. However, in the case of TPC-H query 19 the selection operator is used after a join⁵. Figure 15 shows that the inline-aware implementation of the se-

⁴The stream-fusion engine should have a special care for handling merge joins followed by filter operations. By skipping the elements in the main loop of merging, many CPU cycles are wasted for retrieving the next satisfying element. However, accessing them by using a similar approach to the Iterator model (keep iterating until the next satisfying element is found in a tight loop) gives a better performance.

⁵An alternative implementation is to fuse the selections happening after joins in the join operator itself. The experiments performed in [47] are based on this assumption for join operators. This means that the join operator is not a pure join operator, but a super operator containing a join operator followed by a selection operator. For the purposes of this paper we do not consider such cases.

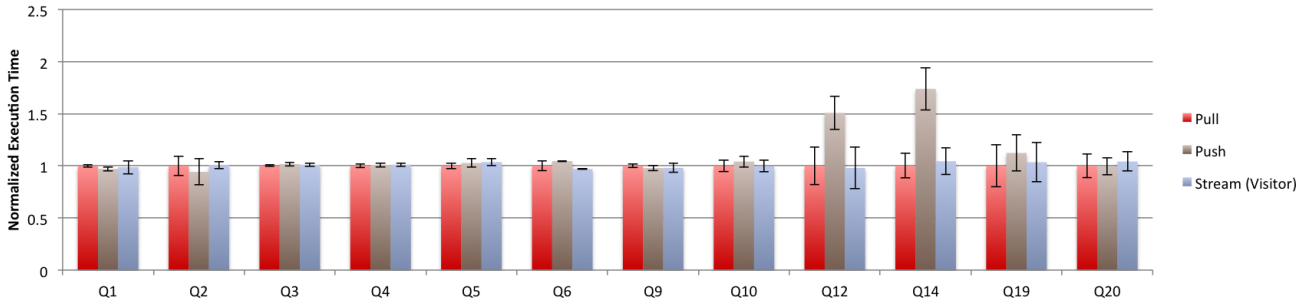


Figure 16: Performance of different query engines for TPC-H queries.

lection operator in pull engines, which is shown using the “Pull” label throughout this section, improves performance by 15%. One of the main reasons is that the inline-aware implementation generates around 40% less query processing code in comparison with the naïve implementation for query processing in these two queries. This improves instruction cache locality as a larger part of the code can fit into the instruction cache.

Removing Intermediate Object Allocations. Figure 15 shows that heap allocating intermediate `Step` objects degrades performance in an order of magnitude. Also, the Visitor pattern for `Step` objects improves performance by 50% in comparison with the case in which heap allocations are converted to stack allocations. Furthermore, our experiments show that for TPC-H queries converting heap allocations to stack allocations (either by Visitor pattern or scalar replacement) decreases the memory consumption from 14 GBs to 11 GBs.

Different Engines on Analytical Queries. Figure 16 shows the performance of several TPC-H queries using different engines. Overall, this figure shows that the difference between engines is not in terms of “orders of magnitude”; in most cases, improvements are minor. This is because the comparison is performed in a fair scenario in which specialization is performed on all engines, in contrast with previous work in which operator inlining was not applied to pull engines [30].

The benchmarked queries can be divided into the following two categories. The first category consists of the queries which are performing almost equally in all engines. In some cases, we see a minor improvement in push engines mainly because of the better control-flow of the generated code which allows the underlying compiler to generate better machine code. However, even in such cases the difference is marginal.

The second category consists of the queries which perform better in pull engines. This is mainly because of using merge join and limit operators. Query 14 falls into this category because of its use of the limit operator. This query has an average 80% speed up for a pull engine in comparison with a push engine. Also, query 12 uses the merge join operator and has an average 70% speed up in comparison with a push-based query engine. It is important to note that in query 12, the query plan that uses a merge join is almost two times faster than the one that uses hash join. This is because both input relations are already sorted on the join key. Hence, the merge join implementation can perform the join on the fly, as opposed to the hash join implementation which needs to construct an intermediate hash table while joining two input relations.

The stream-fusion engine always uses the Visitor pattern throughout this experiment. Interestingly, it is performing as well as push engines, whenever control flow is important. Furthermore, in the cases where push engines require to break the pipeline (the limit

and merge join operators) the stream-fusion engine is performing as well as pull engines. This makes the stream-fusion engine an appropriate choice for query engines.

7. DISCUSSION: PARALLELISM

In this section, we discuss how the results of this paper apply to parallel query engines. *Intra-operator parallelism* is one of the main ways of achieving parallelism in query engines. In this approach, data is *split* among different threads and each thread is responsible for performing the computation on its associated chunk in a sequential manner. At the end, the results computed by different threads are *merged* into a single result.

The split and merge operators [38] are injected between pipeline breaker operators (such as aggregation and hash join operators). Hence, each thread is sequentially computing the result of a chain of pipelining operators. As a result, by using the same split and merge operators, the only difference between pull and push-based engines is how efficiently they compute the result of a chain of pipelining operators. Hence, given a fair environment for pull and push-based engines for intra-operator parallelism, the experimental results we show in this paper for single-threaded scenarios can be expected to match those for multi-threaded scenarios.

One of the key decisions for intra-operator parallelism is the time when the partitioning decision is made. If this decision is made during query compilation time, it is called *plan-driven*. If making this decision is postponed until the runtime, it is called *morsel-driven* [36]. The key advantage of the latter one is using the runtime information and performing better load balancing by using work-stealing [5]. However, the partitioning decision choice is also independent of the type of query engine. All types of query engines can use both approaches, and the impact of work-stealing scheduling is similar on both of them.

Similar efforts have been conducted in the PL community for parallelism in collection programs. As an example, morsel-driven parallelism [36] shares similar ideas, such as work-stealing scheduling, with parallel collection programming libraries [46].

8. CONCLUSION

If one effects a fair comparison of push and pull-based query processing – particularly if one attempts to inline code in both approaches as much as possible – neither approach clearly outperforms the other. We have discussed the reasons for this, and indeed, when considered closely how each approach fundamentally works, it should seem rather surprising if either approach dominated the other performance-wise.

We have also drawn close connections to three fundamental approaches to loop fusion in programming languages – fold, unfold, and stream fusion. As it turns out, there is a close analogy between

pull engines and unfold fusion on one hand and push engines and fold fusion on the other.

Finally, we have applied the lessons learned about the weaknesses of either approach and propose a new approach to building query engines which draws its inspiration from stream fusion and combines the individual advantages of pull and push engines, avoiding their weaknesses.

9. REFERENCES

- [1] StreamBase Systems, <http://www.streambase.com>.
- [2] Y. Ahmad and C. Koch. DBToaster: A SQL compiler for high-performance delta processing in main-memory databases. *PVLDB*, 2(2):1566–1569, 2009.
- [3] M. Armbrust, R. S. Xin, C. Lian, Y. Huai, D. Liu, J. K. Bradley, X. Meng, T. Kaftan, M. J. Franklin, A. Ghodsi, and M. Zaharia. Spark SQL: Relational Data Processing in Spark. SIGMOD '15, pages 1383–1394, New York, NY, USA, 2015. ACM.
- [4] A. Biboudis, N. Palladinos, G. Fourtounis, and Y. Smaragdakis. Streams à la carte: Extensible pipelines with object algebras. In *29th European Conference on Object-Oriented Programming*, page 591, 2015.
- [5] R. D. Blumofe and C. E. Leiserson. Scheduling multithreaded computations by work stealing. *Journal of the ACM (JACM)*, 46(5):720–748, 1999.
- [6] C. Böhm and A. Berarducci. Automatic synthesis of typed λ -programs on term algebras. *Theoretical Computer Science*, 39:135–154, 1985.
- [7] V. Breazu-Tannen, P. Buneman, and L. Wong. *Naturally embedded query languages*. Springer, 1992.
- [8] V. Breazu-Tannen and R. Subrahmanyam. *Logical and computational aspects of programming with sets/bags/lists*. Springer, 1991.
- [9] P. Buchlovsky and H. Thielecke. A type-theoretic reconstruction of the visitor pattern. *Electronic Notes in Theoretical Computer Science*, 155:309 – 329, 2006.
- [10] J.-D. Choi, M. Gupta, M. Serrano, V. C. Sreedhar, and S. Midkiff. Escape analysis for java. *Acm Sigplan Notices*, 34(10):1–19, 1999.
- [11] D. Coutts, R. Leshchinskiy, and D. Stewart. Stream fusion. from lists to streams to nothing at all. In *ICFP '07*, 2007.
- [12] A. Crotty, A. Galakatos, K. Dursun, T. Kraska, U. Çetintemel, and S. B. Zdonik. Tupleware: "big" data, big analytics, small clusters. In *CIDR*, 2015.
- [13] C. Diaconu, C. Freedman, E. Ismert, P.-A. Larson, P. Mittal, R. Stonecipher, N. Verma, and M. Zwillig. Hekaton: Sql server's memory-optimized oltp engine. In *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data*, SIGMOD '13, pages 1243–1254, New York, NY, USA, 2013. ACM.
- [14] B. Emir, M. Odersky, and J. Williams. Matching objects with patterns. ECOOP'07, pages 273–298, Berlin, Heidelberg, 2007. Springer-Verlag.
- [15] L. Fegaras and D. Maier. Optimizing object queries using an effective calculus. *ACM Trans. Database Syst.*, 25(4):457–516, Dec. 2000.
- [16] B. Gedik, H. Andrade, K.-L. Wu, P. Yu, and M. Doo. SPADE: the System S declarative stream processing engine. In *SIGMOD*, 2008.
- [17] J. Gibbons and B. C. d. S. Oliveira. The essence of the iterator pattern. *Journal of Functional Programming*, 19(3-4):377–402, 2009.
- [18] A. Gill, J. Launchbury, and S. L. Peyton Jones. A short cut to deforestation. FPCA, pages 223–232. ACM, 1993.
- [19] G. Graefe. Volcano-an extensible and parallel query evaluation system. *IEEE Transactions on Knowledge and Data Engineering*, 6(1):120–135, 1994.
- [20] T. Grust, M. Mayr, J. Rittinger, and T. Schreiber. FERRY: database-supported program execution. SIGMOD 2009, pages 1063–1066. ACM.
- [21] T. Grust, J. Rittinger, and T. Schreiber. Avalanche-safe LINQ compilation. *PVLDB*, 3(1-2):162–172, Sept. 2010.
- [22] T. Grust and M. Scholl. How to comprehend queries functionally. *Journal of Intelligent Information Systems*, 12(2-3):191–218, 1999.
- [23] R. Hinze, T. Harper, and D. W. H. James. Theory and practice of fusion. In *Proceedings of the 22Nd International Conference on Implementation and Application of Functional Languages*, IFL'10, pages 19–37, Berlin, Heidelberg, 2011. Springer-Verlag.
- [24] M. Hirzel, R. Soulé, S. Schneider, B. Gedik, and R. Grimm. A catalog of stream processing optimizations. *ACM Comput. Surv.*, 46(4):46:1–46:34, Mar. 2014.
- [25] C. Hofer and K. Ostermann. Modular domain-specific language components in scala. In *Proceedings of the Ninth International Conference on Generative Programming and Component Engineering*, GPCE '10, pages 83–92, New York, NY, USA, 2010. ACM.
- [26] S. P. Jones, C. Hall, K. Hammond, W. Partain, and P. Wadler. The glasgow haskell compiler: a technical overview. In *Proc. UK Joint Framework for Information Technology (JFIT) Technical Conference*, volume 93. Citeseer, 1993.
- [27] M. Jonnalagedda and S. Stucki. Fold-based fusion as a library: A generative programming pearl. In *Proceedings of the 6th ACM SIGPLAN Symposium on Scala*, pages 41–50. ACM, 2015.
- [28] M. Karpathiotakis, I. Alagiannis, and A. Ailamaki. Fast queries over heterogeneous data through engine customization. *Proceedings of the VLDB Endowment*, 9(12):972–983, 2016.
- [29] M. Karpathiotakis, I. Alagiannis, T. Heinis, M. Branco, and A. Ailamaki. Just-in-time data virtualization: Lightweight data management with vida. In *CIDR*, 2015.
- [30] Y. Klonatos, C. Koch, T. Rompf, and H. Chafi. Building efficient query engines in a high-level language. *PVLDB*, 7(10):853–864, 2014.
- [31] Y. Klonatos, C. Koch, T. Rompf, and H. Chafi. Errata for "building efficient query engines in a high-level language": PvlDb 7(10):853-864. *Proc. VLDB Endow.*, 7(13):1784–1784, Aug. 2014.
- [32] C. Koch. Incremental query evaluation in a ring of databases. PODS 2010, pages 87–98. ACM, 2010.
- [33] C. Koch. Abstraction without regret in database systems building: a manifesto. *IEEE Data Eng. Bull.*, 37(1):70–79, 2014.
- [34] C. Koch, Y. Ahmad, O. Kennedy, M. Nikolic, A. Nötzli, D. Lupei, and A. Shaikhha. DBToaster: higher-order delta processing for dynamic, frequently fresh views. *VLDBJ*, 23(2):253–278, 2014.
- [35] K. Krikellas, S. Viglas, and M. Cintra. Generating code for holistic query evaluation. In *ICDE*, pages 613–624, 2010.
- [36] V. Leis, P. Boncz, A. Kemper, and T. Neumann. Morsel-driven Parallelism: A NUMA-aware Query Evaluation Framework for the Many-core Age. SIGMOD '14, pages 743–754, New York, NY, USA, 2014. ACM.
- [37] R. A. Lorie. *XRM: An extended (N-ary) relational memory*. IBM, 1974.
- [38] M. Mehta and D. J. DeWitt. Managing intra-operator parallelism in parallel database systems. In *VLDB*, volume 95, pages 382–394, 1995.
- [39] E. Meijer, B. Beckman, and G. Bierman. LINQ: Reconciling Object, Relations and XML in the .NET Framework. SIGMOD '06, pages 706–706. ACM, 2006.
- [40] D. G. Murray, M. Isard, and Y. Yu. Steno: Automatic optimization of declarative queries. PLDI '11, pages 121–131, New York, NY, USA, 2011. ACM.
- [41] F. Nagel, G. Bierman, and S. D. Viglas. Code generation for efficient query processing in managed runtimes. *PVLDB*, 7(12):1095–1106.
- [42] T. Neumann. Efficiently Compiling Efficient Query Plans for Modern Hardware. *PVLDB*, 4(9):539–550, 2011.
- [43] S. Pantela and S. Idreos. One loop does not fit all. In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data*, pages 2073–2074. ACM, 2015.
- [44] J. Paredaens and D. V. Gucht. Possibilities and limitations of using flat operators in nested algebra expressions. In *Proceedings of the Seventh ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, March 21-23, 1988, Austin, Texas, USA*, pages 29–38, 1988.
- [45] B. C. Pierce. *Types and programming languages*. MIT press, 2002.
- [46] A. Prokopec, P. Bagwell, T. Rompf, and M. Odersky. A generic parallel collection framework. In *Euro-Par 2011 Parallel Processing*, pages 136–147. Springer, 2011.
- [47] S. Schuh, X. Chen, and J. Dittrich. An experimental comparison of thirteen relational equi-joins in main memory. 2016.

[48] A. Shaikhha, Y. Klonatos, L. Parreaux, L. Brown, M. Dashti, and C. Koch. How to architect a query compiler. SIGMOD'16, 2016.

[49] O. Shivers and M. Might. Continuations and transducer composition. PLDI '06, pages 295–307. ACM, 2006.

[50] J. Svenningsson. Shortcut fusion for accumulating parameters & zip-like functions. ICFP '02, pages 124–132. ACM, 2002.

[51] Transaction Processing Performance Council. TPC-H, a decision support benchmark. <http://www.tpc.org/tpch>.

[52] P. Trinder. Comprehensions, a Query Notation for DBPLs. In *Proc. of the 3rd DBPL workshop*, DBPL3, pages 55–68, San Francisco, CA, USA, 1992. Morgan Kaufmann Publishers Inc.

[53] S. Viglas, G. M. Bierman, and F. Nagel. Processing Declarative Queries Through Generating Imperative Code in Managed Runtimes. *IEEE Data Eng. Bull.*, 37(1):12–21, 2014.

[54] J. Vlissides, R. Helm, R. Johnson, and E. Gamma. Design patterns: Elements of reusable object-oriented software. *Reading: Addison-Wesley*, 49(120):11, 1995.

[55] P. Wadler. Deforestation: Transforming programs to eliminate trees. In *ESOP'88*, pages 344–358. Springer, 1988.

[56] P. Wadler. Comprehending monads. In *Proceedings of the 1990 ACM Conference on LISP and Functional Programming*, LFP '90, pages 61–78, New York, NY, USA, 1990. ACM.

[57] M. Zaharia, M. Chowdhury, T. Das, A. Dave, J. Ma, M. McCauley, M. J. Franklin, S. Shenker, and I. Stoica. Resilient Distributed Datasets: A Fault-tolerant Abstraction for In-memory Cluster Computing. NSDI'12. USENIX Association.

APPENDIX

A. MICRO BENCHMARK QUERIES

filter.count: <pre>SELECT COUNT(*) FROM LINEITEM WHERE L_SHIPDATE >= DATE '1995-12-01'</pre>	filter.sum: <pre>SELECT SUM(L_DISCOUNT * L_EXTENDEDPRICE) FROM LINEITEM WHERE L_SHIPDATE >= DATE '1995-12-01'</pre>	filter.filter.sum: <pre>SELECT SUM(L_DISCOUNT * L_EXTENDEDPRICE) FROM LINEITEM WHERE (L_SHIPDATE >= DATE '1995-12-01') AND (L_SHIPDATE < DATE '1997-01-01')</pre>
filter.map: <pre>SELECT L_DISCOUNT * L_EXTENDEDPRICE FROM LINEITEM WHERE L_SHIPDATE >= DATE '1995-12-01'</pre>	filter.sort.take: <pre>SELECT L_EXTENDEDPRICE FROM LINEITEM WHERE L_SHIPDATE >= DATE '1995-12-01' ORDER BY L_ORDERKEY LIMIT 1000</pre>	filter.map.take: <pre>SELECT L_DISCOUNT * L_EXTENDEDPRICE FROM LINEITEM WHERE L_SHIPDATE >= DATE '1995-12-01' LIMIT 1000</pre>
filter.XJoin(filter).sum: <pre>SELECT SUM(O_TOTALPRICE) FROM LINEITEM, ORDERS WHERE O_ORDERDATE >= DATE '1998-11-01' AND L_SHIPDATE >= DATE '1998-11-01' AND O_ORDERKEY = L_ORDERKEY</pre>		

Table 4: SQL queries of micro benchmark queries.