

A Scalable Parallel XQuery Processor

E. Preston Carman, Jr.¹, Till Westmann^{2§}, Vinayak R. Borkar^{3*}, Michael J. Carey⁴, Vassilis J. Tsotras¹
¹University of California, Riverside ²Couchbase ³X15 Software, Inc. ⁴University of California, Irvine
 Email: ecarm002@ucr.edu

Abstract—The wide use of XML for document management and data exchange has created the need to query large repositories of XML data. To efficiently query such large data and take advantage of parallelism, we have implemented Apache VXQuery, an open-source scalable XQuery processor. The system builds upon two other open-source frameworks: Hyracks, a parallel execution engine, and Algebricks, a language agnostic compiler toolbox. Apache VXQuery extends these frameworks and provides an implementation of the XQuery specifics (data model, data-model dependent functions and optimizations, and a parser). We describe the architecture of Apache VXQuery, its integration with Hyracks and Algebricks, and the XQuery optimization rules applied to the query plan to improve path expression efficiency and to enable query parallelism. An experimental evaluation using a real 500GB dataset with various selection, aggregation and join XML queries shows that Apache VXQuery performs well both in terms of scale-up and speed-up. Our experiments show that it is about 3.5x faster than Saxon (an open-source and commercial XQuery processor) on a 4-core, single node implementation, and around 2.5x faster than Apache MRQL (a MapReduce-based parallel query processor) on an eight (4-core) node cluster.

I. INTRODUCTION

The widespread acceptance of XML as a standard for document management and data exchange has enabled the creation of large repositories of XML data. To efficiently query such large data collections, a scalable implementation of XQuery is needed that can take advantage of parallelism. While there are various native open-source XQuery processors (Saxon [1], Galax [2], etc.) they have been optimized for single-node processing and do not support scaling to many nodes. To create a scalable XQuery processor, one could: 1) add scalability to an existing XQuery processor, 2) start from scratch, or 3) extend an existing scalable query framework to support XQuery. Unfortunately, existing XQuery processors would require extensive rewriting of their core architecture features to add parallelism. Similarly, building an XQuery processor from scratch would involve the same complex scalable programming (some unrelated to the XML data model). The last option, extending an existing scalable framework to support XQuery, seems advantageous since it combines the benefits of proven parallel technology with a shorter time to implementation.

Among the several scalable frameworks available, one could use a relational parallel database engine and take

advantage of its mature optimization techniques. However, this entails the overhead of translating the data/queries to the relational model and back to XML; moreover, long XML path queries may result in many joins. Another approach is to build an XQuery processor on top of the MapReduce [3] framework. Examples include ChuQL [4], which is a MapReduce extension to XQuery built on top of Hadoop [5], and HadoopXML [6], which combines many XPath queries into a few Hadoop MapReduce jobs. Similarly, Apache MRQL [7] translates XPath queries into an SQL-like language implemented through MapReduce operators. However, these Hadoop-based approaches are limited in that they can only use the few MapReduce operators available (i.e. map, reduce, and combine).

Recently, frameworks have been proposed that generalize the MapReduce execution model by supporting a larger set of operators to create parallel jobs (including Hyracks [8], Spark [9], and Stratosphere [10]). Such 'dataflow' systems [11] typically include flexible data models supporting a wide range of data formats (relational, semi-structured, text, JSON, XML, etc.) which makes them easy to extend. In this paper, we utilize Hyracks as our parallel framework and use Algebricks [12], a language agnostic compiler toolbox, to implement XQuery.

Our implementation is available as open source at the ASF [13]. We have performed an experimental evaluation using a large (500GB) real dataset (a NOAA weather dataset from [14]) and various selection, aggregation, and join XML queries that show the efficiency of our XQuery processor, both in terms of speed up and scale up.

The rest of this paper is organized as follows: Section II reviews current approaches for querying large XML data repositories while Section III covers the Apache VXQuery software stack with details about the underlying framework (Hyracks and Algebricks) and how the data model, parser, and runtime were extended for XQuery support. Given the specifics of XQuery, we had to extend existing Algebricks rewrite rules and introduce new ones; this discussion appears in Section IV. Finally, Section V presents the results of our experiments on Apache VXQuery's performance as well as a comparison with two open-source XML processors – the single-threaded SaxonHE and the parallel Apache MRQL.

[§]work done while at Oracle Labs.

^{*}work done while at the University of California, Irvine.

II. RELATED WORK

Hadoop [15] provides a framework for distributed processing based on the MapReduce model. That leaves a significant implementation burden on the application programmer. As a result, a number of languages have been proposed on top of Hadoop (e.g. Hive [16], PigLatin [17], and Jaql [18]); however, popular high-level MapReduce languages do not support the XML data model. Recent approaches to close this gap include: ChuQL, Apache MRQL, HadoopXML, and Oracle XQuery for Hadoop.

ChuQL [4] extends XQuery to include MapReduce support for processing native XML on Hadoop. In ChuQL, a MapReduce expression is included as an XQuery function, allowing the query writer to specify the MapReduce job definition in XQuery. In contrast, VXQuery hides all parallel processing details from the query writer while still using standard XQuery constructs.

Apache MRQL (MapReduce Query Language) [19], [7] is a SQL-like language designed to run big data analysis tasks. The language supports parsing XML data from Hadoop through a *source* expression defining the XML parser, XML file, and XML tags. The XML parser processes the XML file and returns all elements matching these tags; Apache MRQL then translates these elements into the Apache MRQL data model. Each query is translated to an algebra expression for the Apache MRQL cost-based optimizer, which builds upon known relational query and MapReduce optimization techniques. The algebra uses a small number of physical operators to create a more efficient MapReduce job than directly writing it using the MapReduce operators.

HadoopXML [6] processes a single large XML file with a predetermined set of queries (each currently in a subset of XPath). The engine identifies the query commonalities (paths that are common) and executes those once; it then shares the common results and augments them with the non-common parts per query. This processing is performed using MapReduce jobs. When a query is executed, the query optimizer determines the optimum number of jobs to execute the requested query.

Recently, Oracle released Oracle XQuery for Hadoop (OXH) [20], which runs XQuery data transformations by translating them into a series of MapReduce jobs.

In summary, the above approaches share the MapReduce framework and are thus limited to using only the available MapReduce operators. Apache VXQuery differs in that it is built on top of a more general scalable framework (Hyracks) and can match XQuery computational tasks to Hyracks' richer existing operators (e.g. join); this in turn provides better performance. As will be seen in our experimental section, our rewrite rules, together with Hyracks' efficiency, provides over twice the performance of approaches that perform XML processing on top of MapReduce.

PAXQuery [21] implements XQuery top of Stratosphere [10] (a dataflow system that is similar to Hyracks). The

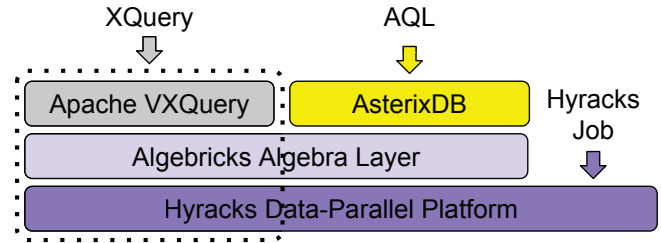


Figure 1. The layers of the Apache VXQuery stack.

system translates XQuery queries into an internal XQuery algebra and then into Parallelization Contracts (PACTs) while Apache VXQuery translates the query into a language agnostic algebra (Algebricks) and then into a Hyracks job for execution. PAXQuery builds on previous unnesting optimizations for tuple-based XQuery algebras [22], [23], [24], [25] Since Apache VXQuery also uses a tuple-based algebra, the same optimization techniques can be applied to the Algebricks query plans. PAXQuery was not available for comparison as of the writing of this paper. Similarly, the Apache MRQL group is currently working on supporting Apache MRQL on top of Apache Flink [26] (which evolved from the Stratosphere project) but at the time of this writing that implementation was still under development.

III. APACHE VXQUERY'S STACK

Apache VXQuery's software stack can be represented in three layers, as shown in Figure 1. The top layer, Apache VXQuery, forms an Algebricks logical plan based on parsing a supplied XQuery. The initial Algebricks logical plan is then optimized and translated into an Algebricks physical plan that maps directly to a Hyracks job. The Hyracks platform executes the job and returns the results. A brief explanation of each layer in the stack follows in the next subsection. Figure 1 also shows AsterixDB [27], another system that uses the Hyracks and Algebricks infrastructure,

A. Hyracks

Hyracks is a data-parallel execution platform that builds upon mature parallel database techniques and modern big data trends [8], [28]. This generic platform offers a framework to run dataflows in parallel on a shared-nothing cluster. The system was designed to be independent of any particular data model. Hyracks processes data in partitions of contiguous bytes, moving data in fixed-sized frames that contain physical records, and it defines interfaces that allow users of the platform to specify the data-type details for comparing, hashing, serializing and de-serializing data. Hyracks provides built-in base data types to support storing data on local partitions or when building higher level data types (first row of Table I).

A Hyracks job is defined by a dataflow DAG with operators (nodes) and connectors (edges). During execution, the operators allow the computation to consume an input

partition and produce an output partition while the connectors redistribute data among partitions. The dataflow among Hyracks operators is push-based – each source (producer) operator pushes the output frames to a target (consumer) operator. The extensible runtime platform provides a number of operators and connectors for use in forming Hyracks jobs. While each operator’s operation is defined by Hyracks, the operator relies on data-model specific functionality provided by the client of the platform.

B. Algebricks

Algebricks [12], [29] is a parallel framework providing an abstract algebra for parallel query translation and optimization. This language-agnostic toolbox complements the lower-level extensible Hyracks platform. Implementations of data-intensive languages can extend its model-agnostic algebraic layer to create parallel query processors on top of the Hyracks platform. A language developer is free to define the language and data model when using the Hyracks platform and the Algebricks toolkit. Algebricks features a rule-based optimizer and data model neutral operators that each allow for language specific customization.

A system that uses Algebricks for its query processing provides its own parser and translator to translate a query to a query plan that uses Algebricks’ logical operators as an intermediate representation. The Algebricks rule-based optimizer then transforms the query plan over three stages. The first is a Logical-to-Logical plan optimizer that creates alternate logical plans. Once the logical plan is finalized, the Logical-to-Physical plan optimizer converts the logical operators into a physical plan. Then, the physical optimizer considers the operator characteristics, partition properties, and data locality to choose the optimal physical implementation for the plan. Algebricks provides generic language-independent rewrite rules for each stage and allows for the addition of other rules. Finally, a Hyracks job is generated and submitted for execution on a Hyracks cluster.

Algebricks’ intermediate logical algebra uses logical operators that map onto Hyracks’ physical operators. A logical operator’s properties are considered when determining the best physical operator. For example, a join query that has an equijoin predicate allows a hash based join instead of the default nested loop join. The Algebricks logical operators exchange data in the form of logical tuples, each of which is a set of fields. The field names are represented by \$\$ followed by a number in remaining text. The following Algebricks logical operators are commonly used in VXQuery:

The DATASCAN operator reads from a data source and returns one tuple for each item in the data source.

The ASSIGN operator executes a scalar expression on a tuple and adds the result as a new field in the tuple.

The DISTRIBUTE-RESULT operator collects the local query results on each data node. Once the job is completed

the controller will request each local result and transfer it to the user to create a complete result.

The EMPTY-TUPLE-SOURCE operator contributes the first tuple without any fields. Algebricks uses this operator to start all DAG dataflow paths.

The JOIN operator matches and combines tuples from two streams of input tuples.

The AGGREGATE operator executes an aggregate expression to create a result tuple from a stream of input tuples. The result is held until all tuples are processed and then returned in a single tuple.

The UNNEST operator executes an unnesting expression for each tuple, creating a stream of single item tuples.

The SUBPLAN operator executes a nested plan for each tuple input.

The NESTED-TUPLE-SOURCE operator is used as the initial operator in nested plans. The operator links the nested plans with the input to the operator (such as a SUBPLAN) defining the nested plan.

The Algebricks operators are each parameterized with custom expressions. The expressions map directly to specific language functions or support runtime features. Each expression is an instance of one of the following expression types: scalar, aggregate, and unnesting. Most operators use a scalar expression, while the AGGREGATE and UNNEST operators have their own expression types. The three expression types differ in their input and output cardinalities. Scalar expressions operate on a single value and return a single value. Aggregate expressions consume many values to create a single result. Unnesting expressions consume a single (usually structured) value to create many new values. Correspondingly, the AGGREGATE and UNNEST operators change the cardinality of the tuple stream.

C. Apache VXQuery

Apache VXQuery extends the language agnostic layer provided by Algebricks to create a scalable XQuery processor. Apache VXQuery provides a binary representation of the XQuery Data Model (XDM) (an example can be found at the Apache VXQuery website [13]), an XQuery parser, an XQuery optimizer, and the data model dependent expressions. VXQuery can process data that is supplied in non-fragmented XML documents partitioned evenly throughout a cluster. A SAX based XML parser translates the XML documents at runtime into XDM instances. Hyracks base types were extended to build untyped XDM instances for the XQuery node types and the XQuery atomic types. (All XQuery types used are listed in Table I.)

Query evaluation proceeds through the usual steps. The query is parsed into an abstract syntax tree (AST) and is then analyzed, normalized, and translated into a logical plan. The logical plan consists of Algebricks data model independent operators parameterized with Apache VXQuery data model dependent expressions. The logical plan is then optimized

Hyracks Base	boolean, byte, short, integer, long, double, float, UTF8 string
XQuery Atomic	binary, decimal, date, datetime, time, duration, QName
XQuery Node	attribute, comment, document, element, processing instruction, text

Table 1
APACHE VXQUERY BUILDS ON THE HYRACKS BASE TYPES TO CREATE THE XQUERY ATOMIC AND NODE DATA TYPES.

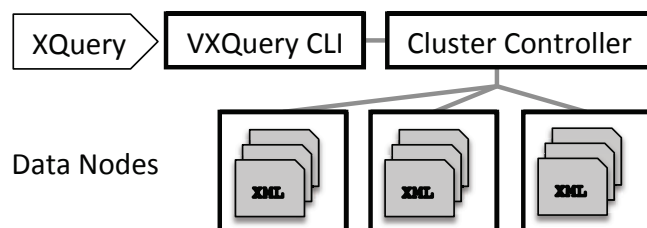


Figure 2. The VXQuery cluster configuration.

using both generic rewrite rules provided by Algebricks and XQuery specific rewrite rules provided by Apache VXQuery (discussed in Section IV). After rewriting the logical plan, it is translated into a physical plan and optimized further (physical optimization includes, e.g., the selection of join methods or the distribution of the plan). Finally the physical plan is translated into a Hyracks job that is executed. Similar to Algebricks operators that have physical representations based on Hyracks operators, Apache VXQuery provides executable functions that implement Apache VXQuery’s data model dependent expressions.

Special attention is required regarding how the XDM defines a set of items as a sequence. In Apache VXQuery, an XDM sequence can have two forms: a *sequence item* or a *tuple stream*. A sequence item holds all the values in a single tuple field; a tuple stream represents the sequence using a field with the same name in multiple tuples. To switch between these representations, we provide the *iterate* and the *create_sequence* expressions. The *iterate* unnesting expression works with Algebricks’ UNNEST operator to convert a tuple field that holds a sequence item into a stream of individual tuples. The *create_sequence* aggregate expression executes within Algebricks’ AGGREGATE operator to consume a tuple stream and create a sequence item for inclusion in a single output tuple. The two expressions are used during the logical rewrite process to switch between formats to enable further optimization rules to be applied to the query plan.

At runtime, the Apache VXQuery cluster processes a query using the Apache VXQuery Client Library Interface (CLI), a Hyracks Cluster Controller, and some Hyracks Data Nodes (as shown in Figure 2). The process starts with a user submitting an XQuery statement to the Apache VXQuery CLI for parallel execution. The CLI parses and optimizes the query and submits the generated Hyracks job

to the cluster controller, which manages and distributes tasks to each of the data nodes for evaluation. Each data node contains XML files, an XML parser and the XQuery runtime expressions used to evaluate the node’s tasks. Finally, the cluster controller collects the data nodes’ results and sends the result back to the Apache VXQuery CLI, which returns the result to the user.

IV. REWRITE RULES

Algebricks provides generic rules for both Logical-to-Logical and Logical-to-Physical plan optimizations. These rules include actions that consolidate, push down, and/or remove operators based on the operators’ properties and the query plan. In addition, to build the XQuery optimizer we needed to implement XQuery-specific rules; these rules fall into two categories. The *Path Expression Rewrite Rules* attempt to remove subplans that are introduced by the unnesting required to evaluate path expressions. The *Parallel Rewrite Rules* transform the plan to enable parallel evaluation for specific XQuery constructs (aggregation, join, and data access) using both pipelined and partitioned parallelism.

A. Path Expression Rewrite Rules

The normalization phase of query translation introduces explicit operations into the query plan that ensure the correctness of the plan (for example sorting to maintain document order). However, some of these operations may not be required based on knowledge of the structure of the plan and the implementation of operators and expressions. The following XML segment is based on the sample XML tree from the W3Schools tutorial [30] for XQuery and will be used to outline the path expression rewrite rules.

```

1 <?xml version="1.0" encoding="ISO-8859-1"?>
2 <bookstore>
3   <book id="1" category="COOKING">
4     <title lang="en">Everyday Italian</title>
5     <author>Giada De Laurentiis</author>
6     <year>2005</year>
7   </book>
8   <book id="2" category="CHILDREN">
9     <title lang="en">Harry Potter</title>
10    <author>J K. Rowling</author>
11    <year>2005</year>
12  </book>
13  ...
14 </bookstore>
  
```

Consider the following simple query.

```

1 doc("book.xml")/bookstore/book
  
```

The query reads data from the document *book.xml* located in the file system using the XQuery *doc* function. Next, the first child path step expression (“/bookstore”) is applied to the document node. Three stages are used when applying the child path step expression to a tuple: each input node is iterated over, any matching child nodes are put into a single sequence, and the sequence is then sorted in document order. The same path step process is applied to each resulting “bookstore” element node for the second child path step

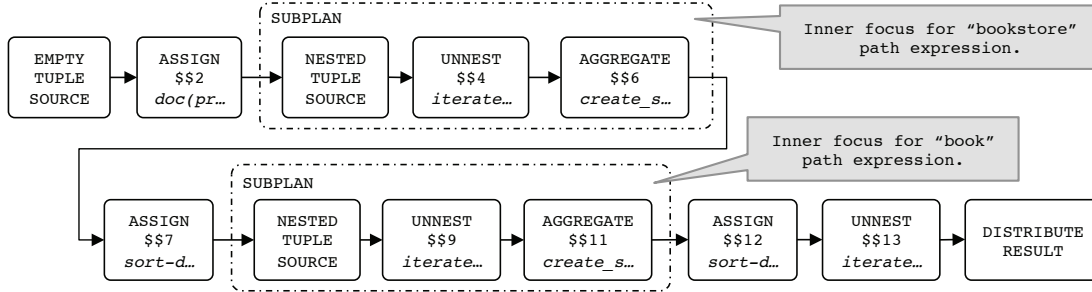


Figure 3. Example query dataflow DAG before applying rewrite rules.

expression (`"/book`). Finally, each `"book"` element node is then returned in the final query result.

VXQuery creates the initial plan shown below (after removing unused variables); here the curly braces represent nested plans that are executed for each of the SUBPLAN's input tuples. Schematically this plan, which is read bottom-up, corresponds to the dataflow DAG in Figure 3. The DAG is a single path of execution in this case. Each DAG is initialized with an EMPTY-TUPLE-SOURCE operator and collects its results into a DISTRIBUTE-RESULT operator.

```

1 DISTRIBUTE-RESULT( $$13 )
2 UNNEST( $$13:iterate($$12) )
3 ASSIGN( $$12:sort-distinct-nodes-asc-or-atomics($$11) )
4 SUBPLAN {
5   AGGREGATE( $$11:create_sequence(child(treat($$9,
6     element_node), "bookstore")) )
7   UNNEST( $$9:iterate($$7) )
8   NESTED-TUPLE-SOURCE
9 }
10 ASSIGN( $$7:sort-distinct-nodes-asc-or-atomics($$6) )
11 SUBPLAN {
12   AGGREGATE( $$6:create_sequence(child(treat($$4,
13     element_node), "bookstore")) )
14   UNNEST( $$4:iterate($$2) )
15   NESTED-TUPLE-SOURCE
16 }
17 ASSIGN( $$2:doc(promote(data("books.xml"), string) ) )
18 EMPTY-TUPLE-SOURCE

```

The plan's EMPTY-TUPLE-SOURCE operator creates the initial empty tuple. The `doc` expression in the ASSIGN operator (line 15) returns a document node using the string URI argument and adds a new field – `$$2:document node` – to the tuple. The `promote` and `data` expressions ensure the `doc` URI argument will be a string. The SUBPLAN operator (line 10) uses a nested plan to implement the first and second stages of the `/bookstore` path step. The subplan's nested plan ensures the correct dynamic context for the path step and provides an "inner focus" to evaluate the expression on each item in the sequence for the next step (if any). The NESTED-TUPLE-SOURCE operator (line 13) connects the nested plan to the SUBPLAN's input dataflow. The input tuple is passed on to the UNNEST operator (line 12) where each `$$2:document` item is iterated over and added as `$$4:document`. For the root path step expression, there is only one item in the sequence. The inner focus is closed with AGGREGATE (line 11) processing all SUBPLAN tuples using the `create_sequence(child(treat($$4, element_node),`

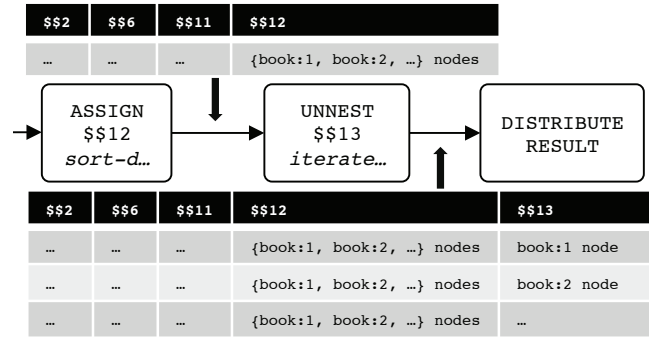


Figure 4. Dataflow segment for the last UNNEST operator.

`"bookstore"))` expression. The expression ensures that `child` expression's argument is of type `"element_node"` (`treat`), finds all `"bookstore"` child nodes (`child`), and creates a sequence of all the results (`create_sequence`). The resulting tuple now holds two fields: `$$2:document node` and `$$6:"bookstore" node`. All the SUBPLAN variables are discarded except the final operator's result (in this case, AGGREGATE `$$6`). The third stage of the path step is completed though the ASSIGN operator (line 9) with `sort-distinct-nodes-asc-or-atomics($$6)`. The expression creates a new field with nodes that are in document order and duplicate free from `$$6:"bookstore"` node. Since there is only one item, the `"bookstore"` node is copied over to `$$7`.

The next SUBPLAN (line 4) creates the inner focus for the `/book` path step expression. Similar to the `/bookstore` path step, the nested plan iterates over the input tuples and saves all child nodes `{book:1, book:2, ...}` to `$$11`. The ASSIGN (line 3) ensures document order in the child `"book"` node sequence by removing duplicates and sorting the sequence. Finally, each item in `$$12:{book:1, book:2, ...}` is unnested by UNNEST (line 2) to create a tuple stream for the DISTRIBUTE-RESULT operator (line 1). See Figure 4 for a graphic representation of the tuples before and after this UNNEST operator.

The initial query plan is inefficient and can be improved in several ways: we can (i) remove the computationally expensive sort operators (as document order is not changed by any of the other operators) and (ii) remove the SUBPLAN operators (since each SUBPLAN corresponds to a simple

step expression the inner focus is not required). After these optimization rules, the plan can be cleaned further by (iii) enabling unnesting (improves operator efficiency) and (iv) merging the path step unnesting operators (reduces number of operators). Due to space constraints, a detailed explanation of these path step expression rewrite rules can be found in our technical report [31].

After applying these rewrite rules recursively to the sample query plan, the resulting plan only uses only a single UNNEST operator to represent the two child path step expressions. The path expression rewrite rules create the following updated sample query plan:

```

1 DISTRIBUTE-RESULT( $$13 )
2 UNNEST( $$13:child(child($$2, "bookstore"), "book") )
3 ASSIGN( $$2:doc("books.xml") )
4 EMPTY-TUPLE-SOURCE

```

B. Parallel Rewrite Rules

After applying the path expression rewrite rules, the plan is optimized for parallel XQuery processing. Hyracks allows for both pipelined and partitioned parallelism. We thus introduce rules to enable the use of Hyracks’ parallel execution features. To take advantage of pipelining in Apache VXQuery we create fine grained data items. For example, the DATASCAN operator introduced next does not compute a whole collection at once, but instead computes chunks that can be fed to the remaining operators. As a side effect, the needed buffer size (Hyracks’ frame) is reduced between the operators in the pipeline. To introduce partitioned parallelism we use partitioned data access for physically partitioned data and we use partitioned parallel algorithms for join and aggregation.

1) *Introduce the DATASCAN Operator:* To query a collection of XML documents, XQuery defines a function called *collection* that maps a string to a sequence of nodes. Apache VXQuery interprets the string as a directory location, reads in data from the files in the directory, and returns all nodes as a single sequence value. Since the collection query considers many documents, it can produce a large number of query results. Instead of gathering all nodes into a single sequence, we would like to send one node at a time through the pipeline. To avoid this problem, we combine the *collection* expression with an *iterate* expression (typically inserted because of a path step or a for clause) to split the large document sequence into many single document tuples, thus reducing the size of the materialized result. Below is a sample *collection* query similar to the previous single document query example, followed by the query plan generated after the path expression rules have been applied.

```

1 collection("/books")/bookstore/book

1 DISTRIBUTE-RESULT( $$13 )
2 UNNEST( $$13:child( child($$4, "bookstore"), "book" ) )
3 UNNEST( $$4:iterate($$2) )
4 ASSIGN( $$2:collection(promote(data("/books"), string)) )
5 EMPTY-TUPLE-SOURCE

```

The path expression rules have conveniently moved an UNNEST *iterate* above the ASSIGN *collection*, creating a stream of XML document tuples. Algebricks offers a DATASCAN operator to directly create a stream of tuples based on a data source. Since *collection* already defines the data source, the DATASCAN operator can be used to replace UNNEST *iterate* and ASSIGN *collection*. The updated query plan is:

```

1 DISTRIBUTE-RESULT( $$13 )
2 UNNEST( $$13:child( child($$4, "bookstore"), "book" ) )
3 DATASCAN( collection("/books"), $$4 )
4 EMPTY-TUPLE-SOURCE

```

The finer grained tuples reduce the buffer size between operators during the query execution. Note that the above rewrite rule allows Apache VXQuery to process any amount of XML data provided that the largest XML document in the collection can fit in Hyracks’ frame size. This constraint can be further reduced to the largest subtree under the query path expression. This is possible when the UNNEST *child* expression is the consumer of a DATASCAN operator. The *child* expression can be merged into the DATASCAN operator to provide even smaller tuples. The query plan is updated to show that the DATASCAN operator has a third argument specifying the child path expression; the updated DATASCAN operator includes the path expression within the collection:

```

1 DISTRIBUTE-RESULT( $$4 )
2 DATASCAN( collection("/books"), $$4, "/bookstore/book" )
3 EMPTY-TUPLE-SOURCE

```

In addition to the improved pipeline, the DATASCAN operator offers a way to introduce partitioned parallelism simply by specifying Apache VXQuery’s partition details to this operator. In Apache VXQuery, data is partitioned among the cluster nodes. Each node has a unique set of XML documents stored under the same directory specified in the *collection* expression. The Algebricks’ physical plan optimizer uses these partitioned data properties details to distribute the query execution. For example, path *"/books"* defined in the *collection* expression is located on each node and represents a unique set of XML documents for the query. These partition properties are added to the DATASCAN operator although this is not shown in the query plan. Adding these properties allows Apache VXQuery to achieve partitioned parallel execution without any parallel programming.

2) *Replace Scalar with Aggregate Expressions:* The XQuery aggregate expressions (*avg*, *count*, *max*, *min*, and *sum*) use a default scalar implementation in a normalized query plan. This implies that the whole result is first stored in a sequence which is then processed to produce the aggregate. Instead of materializing the sequence, we can match the XQuery aggregate expression with an Algebricks AGGREGATE operator. When the Algebricks AGGREGATE operator is used with an XQuery aggregate expression, the result will be incremental aggregation instead of materializing all

records in the operator's buffer. Consider a query that counts the number of book elements in an XML collection and the query plan produced using the previous rules:

```

1 count (
2   for $x in collection("/books")/bookstore/book
3   return $x
4 )

```

```

1 DISTRIBUTE-RESULT( $$17 )
2 UNNEST( $$17:iterate($$16) )
3 ASSIGN( $$16:count($$15) )
4 SUBPLAN {
5   AGGREGATE( $$15:create_sequence($$4) )
6   DATASCAN( collection("/books"), $$4, "/bookstore/book" )
7   NESTED-TUPLE-SOURCE
8 }
9 EMPTY-TUPLE-SOURCE

```

The XQuery aggregate expression *count* is within an ASSIGN operator (line 3). The SUBPLAN finds the bookstore nodes and uses an AGGREGATE operator (line 5) to store them in a sequence. However, there is no UNNEST directly above the SUBPLAN (as shown in our technical report for the path expression rewrite rules) and thus the SUBPLAN cannot be removed. However, the scalar *count* expression applies its calculation on the produced XQuery sequence to create \$\$16's value. Instead, the aggregate *count* expression can replace the *create_sequence* within the Algebricks AGGREGATE operator, thus performing aggregation incrementally instead of first generating a large XQuery sequence. The updated query plan becomes:

```

1 DISTRIBUTE-RESULT( $$17 )
2 UNNEST( $$17:iterate($$16) )
3 SUBPLAN {
4   AGGREGATE( $$16:count($$4) )
5   DATASCAN( collection("/books"), $$4, "/bookstore/book" )
6   NESTED-TUPLE-SOURCE
7 }
8 EMPTY-TUPLE-SOURCE

```

The new plan keeps the pipeline granularity and enables partitioned aggregation processing. An additional Apache VXQuery rule annotates the AGGREGATE operator with local and global aggregate expressions, enabling the use of Algebricks' support for two-step aggregation – each partition calculates its local aggregate result on its data and then transmits the result to a central partition for the global computation. As a result, partitioning also reduces communication thus improving parallel processing efficiency.

3) *Introduce the JOIN Operator:* In XQuery, two distinct datasets can be connected (matched) through a nested *for* loop. The normalized query plan follows the same nested loop, which can be very expensive; we can do better by using a relational-style join. We note that Algebricks provides a JOIN operator as well as a set of language independent rewrite rules to optimize generic query plans. We can thus use these provided rewrite rules to translate the nested loop plan in to a join plan. Consider a query that takes two bookstores (Ann and Joe) and finds books with the same title, and its query plan below:

```

1 for $r in collection("/ann-books")/bookstore/book
2 for $s in collection("/joe-books")/bookstore/book
3 where $r/title eq $s/title
4 return $r

```

```

1 DISTRIBUTE-RESULT( $$32 )
2 UNNEST( $$32:iterate($$27) )
3 SELECT( boolean(value-eq($$27, $$28)) )
4 ASSIGN( $$28:data(child($$26, "title")) )
5 ASSIGN( $$27:data(child($$13, "title")) )
6 DATASCAN( collection("/joe-books"), $$26, "/bookstore/book" )
7 DATASCAN( collection("/ann-books"), $$13, "/bookstore/book" )
8 EMPTY-TUPLE-SOURCE

```

In this example each dataset is identified and accessed by a DATASCAN operator, while the SELECT operator contains the condition for connecting the two datasets (which effectively will become the join condition). The two ASSIGN operators (line 4 and 5) find the title child node and return the atomic value of the node. Three Algebricks language-independent rules are used to introduce the JOIN operator. The first rule converts the nested DATASCAN operator into a cross product; it identifies that each data source is independent and adds the JOIN operator with a condition of true (basically a cross-product). The second Algebricks rule manipulates the DAG to push down operators that only affect one side of the join branch (selection, assign etc). The third rule then merges the SELECT and JOIN operators so the join condition (from the SELECT) is within the JOIN operator. In the final plan, the JOIN operator has one branch from each data source, which allows each branch to be processed locally and then joined together globally.

```

1 DISTRIBUTE-RESULT( $$32 )
2 UNNEST( $$32:iterate($$27) )
3 JOIN( boolean(value-eq($$27, $$28)) )
4 {
5   ASSIGN( $$28:data(child($$26, "title")) )
6   DATASCAN( collection("/joe-books"), $$26, "/bookstore/book" )
7   EMPTY-TUPLE-SOURCE
8 } {
9   ASSIGN( $$27:data(child($$13, "title")) )
10  DATASCAN( collection("/ann-books"), $$13, "/bookstore/book" )
11  EMPTY-TUPLE-SOURCE
12 }

```

Going from here to the final logical plan does not require any custom Apache VXQuery rules, but the physical plan needs more information to choose the most efficient join algorithm. The equality comparison in our sample query allows the use of a more efficient partition-based algorithm. If Algebricks understands the condition characteristics, it can chose an optimal hash-based join. For Algebricks to identify the join condition, this condition must be represented by a boolean Algebricks expression, in this case the Algebricks' *equal* expression for a hash-based join. (Other Algebricks generic expressions include: *and*, *or*, *not*, *less than*, *greater than*, *less than or equal*, *greater than or equal*, *not equal*.) As the extraction of the XQuery's Effective Boolean Value of the value-comparison in the previous plan (`boolean(value-eq(...))`) is equivalent to Algebricks' *equal* expression, we convert one to the other – thus enabling Algebricks to

identify the join. After running the physical optimization rules the Algebricks expression is converted back to the original XQuery expressions for runtime evaluation. As a result, Hyracks will now use a Hybrid-Hash Join algorithm to achieve efficient partitioned parallelism.

V. APACHE VXQUERY PERFORMANCE

To examine the scalability of our XQuery implementation we have performed an experimental evaluation using publicly available weather XML data. We have also performed a comparison of Apache VXQuery with two open source XML processors: Saxon [1] and Apache MRQL [32], [7].

A. Weather Data

The NOAA website [14] offers weather data via an XML-based web service. For our queries, we chose the Global Historical Climatology Network (GHCN)-Daily dataset that includes daily summaries of climate recordings. The core data fields report high and low temperatures, snowfall, snow depth, and rainfall. The complete data definition and field list can be found on NOAA's site [14]. The date, data type, station id, value, and various attributes (i.e., measurement, source, and quality flags) are included for each weather report. In addition, a separate web service provides additional station data: name, latitude, longitude, and date of first and last reading. The datasets used had four different sizes, ranging from 500MB up to 500GB.

B. Queries

Here we consider three basic types of XQuery queries: selection, aggregation and join. The complete benchmark results include additional query variations, but due to space constraints some are shown only in our technical report [31]. For consistency, the queries below follow the same numbering as [31].

Selection: Query 2 finds all readings that report an extreme wind warning. Such warnings occur when the wind speed exceeds 110 mph. (The wind measurement unit, tenths of a meter per second, has been converted to miles per hour.)

```
1 for $r in collection("/sensors")/dataCollection/data
2 where $r/dataType eq "AWND"
3   and decimal(data($r/value)) gt 491.744
4 return $r
```

Query 2. Extreme Wind Warning

Aggregation: Query 4 finds the highest recorded temperature in the weather data set. The Celsius temperature is reported in tenths of a degree.

```
1 max(
2   for $r in collection("/sensors")/dataCollection/data
3     where $r/dataType eq "TMAX"
4     return $r/value
5 ) div 10
```

Query 4. Highest Recorded Temperature

Join: Query 6 finds the highest recorded temperature (TMAX) for each station for each day during the year 2000.

```
1 for $s in collection("/stations")/stationCollection/station
2 for $r in collection("/sensors")/dataCollection/data
3 where $s/id eq $r/station
4   and $r/dataType eq "TMAX"
5   and year-from-dateTime(dateTime(data($r/date))) eq 2000
6 return ($s/displayName, $r/date, $r/value)
```

Query 6. High Temperature per Station

Join and Aggregation: In Query 8 we join two large collections, one that maintains the daily minimum temperature per station and one that contains the daily maximum temperature per station. The join is on the station id and date and finds the daily temperature difference per station and returns the average difference over all stations.

```
1 avg(
2   for $r_min in collection("/sensors_min")/dataCollection/
3     data
4   for $r_max in collection("/sensors_max")/dataCollection/
5     data
6   where $r_min/station eq $r_max/station
7     and $r_min/date eq $r_max/date
8     and $r_min/dataType eq "TMIN"
9     and $r_max/dataType eq "TMAX"
10  return $r_max/value - $r_min/value
11 ) div 10
```

Query 8. Average Daily Temperature Differential

C. Experimental Results

Our performance study explores Apache VXQuery's ability to scale locally with the number of cores and then in a cluster with the number of nodes. In the single node tests, the number of data partitions is varied to demonstrate nodes scaling up to the number of available cores. For these tests, partitions represent data splits and each partition has a separate query execution thread. In the cluster tests, the number of partitions per node has been fixed (to one partition per core) and only the number of nodes is varied. The tests were all executed on an eight-node gigabit-connected cluster. Each node has two dual-core AMD Opteron(tm) processors, 8GB of memory, and two 1TB hard drives.

1) Single Node Experiments: Our single node experiments used one cluster node and repeated each query five times. The reported query time is an average of the last three runs. (In our setting, the first two executions are used to warm up the system.) The first single node experiment compares Apache VXQuery with Saxon [33], which is a highly efficient open source XQuery processor. The freely available Saxon Home Edition (SaxonHE 9.5) is typically limited to a single thread processing data that can fit into one fifth the size of the machine's memory. A group of weather stations were selected to create query results that fit these Saxon data restrictions. The 584MB data set has been partitioned on a single hard drive. The speed-up test keeps the total data set size constant while varying its number of data partitions and corresponding query processing threads.

Figure 5 shows the single node speed-up performance results for Apache VXQuery and Saxon. Only a single experiment is shown for Saxon since multi-threading is not available in the SaxonHE 9.5 version. Apache VXQuery

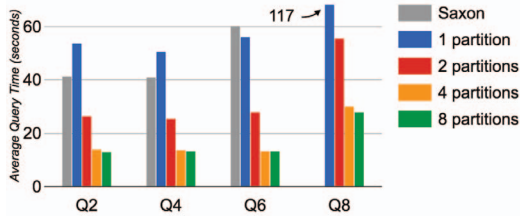


Figure 5. Single node speed-up comparison for Saxon and Apache VXQuery (584MB dataset; varying Apache VXQuery partitions).

outperforms Saxon when it uses two or more partitions. The single partition results are slower due to overhead introduced for parallel and distributed query processing. The join queries (Query 6 and 8) are translated into hash-based joins for Apache VXQuery, thus giving better performance than Saxon’s nested loop join. Saxon’s result for Query 8 is not reported since the large number of joined records caused it to never complete (and based on our other results, this query could take several months to complete). For the rest of the queries (Query 2, 4 and 6), when using 4 or more partitions, Apache VXQuery performed on average about 3.5x faster than Saxon. The Apache VXQuery performance for 8 partitions is similar to its 4 partitions performance, which is when the CPU becomes saturated.

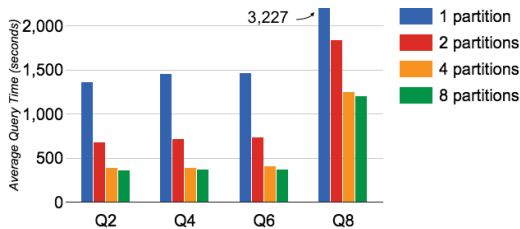


Figure 6. Single node Apache VXQuery speed-up (15.2GB dataset).

To further test single node speed-up for Apache VXQuery, we also used a dataset larger than the node’s memory (8GB). For this we used a XML weather data subset, the GCOS Surface Network (GSN) stations containing 15.2GB of XML data. The results appear in Figure 6. As with the previous figure, Apache VXQuery scales well up to the node’s number of cores (4). Similar to the single node 584MB experiments, the CPU is saturated when using 4 or more partitions. While profiling our experiments, we observed that Apache VXQuery is CPU bound here, despite the larger data size, due to the overhead of parsing the XML document for each query. This is also evident from the improvement in performance when increasing threads.

2) *Cluster Experiments*: Based on the single node speed-up results, the cluster experiments used eight nodes and four partitions per node. The first cluster tests used the U.S. Historical Climatology Network (HCN) stations dataset which holds 57GB of XML weather data. This dataset exceeds the available cluster memory when using less than eight nodes. For each experiment, the dataset was equally divided among the nodes participating in the experiment.

The cluster speed-up results for Apache VXQuery (as well as for Apache MRQL, to be discussed later) appear in Figure 7; the Apache VXQuery query times are depicted by the circles inside the corresponding bars (full bars represent Apache MRQL query times). As can be observed, adding nodes to the cluster proportionally lowers the query time. We next tested the scale-up characteristics of Apache VXQuery. We started by using a dataset that fits in the memory of each node (i.e., 7.2GB of data per node). The results appear in Figure 8 (again Apache VXQuery query times correspond to the circles). While nodes and data are added to the query, the query time remains comparable, that is, the additional data is processed in the same amount of time. Apache VXQuery thus scales up well for Queries 2, 4, 6 and 8 on XML data.

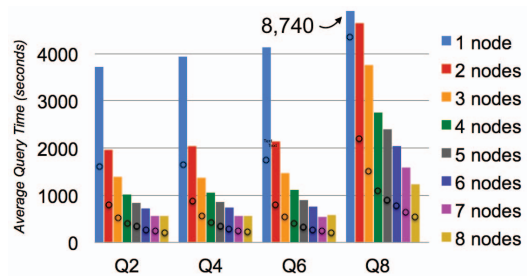


Figure 7. Apache VXQuery and Apache MRQL cluster speed-up (57GB dataset); circles mark the respective Apache VXQuery times.

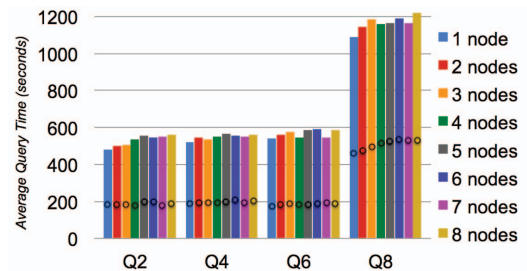


Figure 8. Apache VXQuery and Apache MRQL cluster scale-up (7.2GB per node); circles mark the respective Apache VXQuery times.

The next scale-up test utilizes all 528GB of weather data; here each node has 66GB of data split evenly on two local disks. The results appear in Figure 9; Apache VXQuery clearly scales-up well even for very large XML datasets.

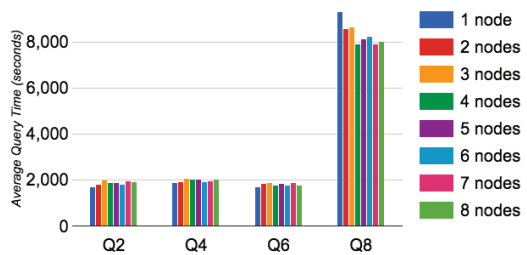


Figure 9. Apache VXQuery cluster scale-up (66GB per node).

Our final experiment sought to evaluate Apache VXQuery’s performance against other open source parallel XML processors. Among them we chose Apache MRQL [7]

as it was readily available. Using the HCN (57GB) dataset, we ran speed-up and scale-up tests for the same queries on Apache MRQL running on top of Hadoop 1.2.1 using MapReduce. Hadoop was configured with a 128MB block size and a replication factor of 1 (to reduce space on the cluster). Apache VXQuery outperforms Apache MRQL on all queries in terms of both scale-up and speed-up (Figures 7 and 8). Apache VXQuery’s performance advantage comes partly from reading and parsing XML about two times faster than Apache MRQL. In addition, its richer set of operators provides for better performance. For example, VXQuery utilizes a Hybrid Hash Join algorithm that can keep a partition in main memory. Being MapReduce-based, Apache MRQL divides the join responsibility: partitioning is done by the mapper, while the reducer joins the individual partitions. These two steps do not share state, yielding a traditional Grace Hash Join. On average over all experiments, VXQuery is 2.5x faster than Apache MRQL on Hadoop, validating the fact that building XQuery on top of a dataflow environment like Hyracks provides more opportunities for optimization and parallelism.

VI. CONCLUSIONS

Apache VXQuery is a scalable open-source XQuery processor that we have built on top of Hyracks and Algebricks. We have described its implementation, including the XML data model dependent rewrite rules. These rules facilitate existing, data model independent Algebricks optimizations that serve to create efficient and parallel Hyracks jobs. We have demonstrated using a real 500GB dataset that VXQuery can scale out to the number of nodes available on a cluster for various XML selection, aggregation, and join queries. Comparatively, Apache VXQuery is about 3.5x faster than Saxon on a single node and around 2.5x faster than Apache MRQL on a cluster in terms of scale-up and speed-up. The VXQuery source code is available at the Apache Software Foundation [13]; the current release contains approximately 100K LOC. We plan to add XQuery 3.0 features that support the analysis of Big Data, such as the `group by` and `window` clauses and to utilize indexing for increased query performance. Apache VXQuery developers are also adding support for large XML documents stored on a distributed file system and further optimizing the query compiler.

Acknowledgments: The work was supported in part by Google Summer of Code (2012-13) and by NSF grants IIS-0910859, IIS-0910989, CNS-1305253 and CNS-1305430.

REFERENCES

- [1] M. Kay, “SAXON: The XSLT and XQuery Processor,” 2004.
- [2] M. Fernández *et al.*, “Implementing XQuery 1.0: The Galax Experience,” in *VLDB*, 2003, pp. 1077–1080.
- [3] J. Dean and S. Ghemawat, “MapReduce: Simplified Data Processing on Large Clusters,” *Commun. ACM*, vol. 51, no. 1, pp. 107–113, 2008.

- [4] S. Khatchadourian *et al.*, “ChuQL: Processing XML with XQuery Using Hadoop,” in *CASCON*, 2011, pp. 74–83.
- [5] D. Borthakur, “The Hadoop Distributed File System: Architecture and Design,” *Hadoop Project Website*, vol. 11, p. 21, 2007.
- [6] H. Choi *et al.*, “HadoopXML: A Suite for Parallel Processing of Massive XML Data with Multiple Twig Pattern Queries,” in *CIKM*, 2012, pp. 2737–2739.
- [7] L. Fegarar *et al.*, “XML Query Optimization in Map-Reduce,” in *WebDB*, 2011.
- [8] V. Borkar *et al.*, “Hyracks: A flexible and extensible foundation for data-intensive computing,” in *IEEE ICDE*, 2011, pp. 1151–1162.
- [9] M. Zaharia *et al.*, “Spark: Cluster Computing with Working Sets,” in *HotCloud*, 2010.
- [10] A. Alexandrov *et al.*, “The Stratosphere Platform for Big Data Analytics,” *The VLDB Journal*, pp. 1–26, 2014.
- [11] S. Babu and H. Herodotou, “Massively Parallel Databases and MapReduce Systems,” *Foundations and Trends in Databases*, vol. 5, no. 1, pp. 1–104, 2013.
- [12] V. Borkar *et al.*, “Algebricks: A Framework for the Analysis and Efficient Evaluation of Data-Parallel Jobs,” *ACM SOCC*, 2015.
- [13] “Apache VXQuery,” <http://vxquery.apache.org/>.
- [14] “National Climate Data Center: Data Access,” <http://www.ncdc.noaa.gov/data-access/>.
- [15] “Apache Hadoop,” <http://hadoop.apache.org/>.
- [16] A. Thusoo *et al.*, “Hive: A Warehousing Solution over a Map-reduce Framework,” *PVLDB*, vol. 2, no. 2, pp. 1626–1629, 2009.
- [17] C. Olston *et al.*, “Pig Latin: A Not-so-foreign Language for Data Processing,” in *SIGMOD*, 2008, pp. 1099–1110.
- [18] K. S. Beyer *et al.*, “Jaql: A Scripting Language for Large Scale Semistructured Data Analysis,” *PVLDB*, vol. 4, no. 12, pp. 1272–1283, 2011.
- [19] L. Fegarar *et al.*, “An Optimization Framework for Map-Reduce Queries,” in *EDBT*, 2012, pp. 26–37.
- [20] “Using Oracle XQuery for Hadoop,” http://docs.oracle.com/cd/E53356_01/doc.30/e53067/oxh.htm.
- [21] J. Camacho-Rodríguez *et al.*, “Paxquery: Efficient parallel processing of complex xquery,” *IEEE TKDE*, vol. 27, no. 7, pp. 1977–1991, 2015.
- [22] P. Boncz *et al.*, “MonetDB/XQuery: a fast XQuery processor powered by a relational engine,” *ACM SIGMOD*, pp. 479–490, 2006.
- [23] A. Deutsch *et al.*, “The NEXT framework for logical XQuery optimization,” *VLDB*, pp. 168–179, 2004.
- [24] N. May *et al.*, “Strategies for query unnesting in XML databases,” *ACM TODS*, vol. 31, no. 3, pp. 968–1013, 2006.
- [25] C. Ré *et al.*, “A complete and efficient algebraic compiler for XQuery,” *IEEE ICDE*, p. 14, 2006.
- [26] “Apache Flink,” <http://flink.incubator.apache.org/>.
- [27] S. Alsubaiee *et al.*, “AsterixDB: A Scalable, Open Source BDMS,” *PVLDB*, vol. 7, no. 14, pp. 1905–1916, 2014.
- [28] “Hyracks,” <https://code.google.com/p/hyracks/>.
- [29] V. Borkar and M. J. Carey, “A Common Compiler Framework for Big Data Languages: Motivation, Opportunities, and Benefits,” *IEEE Data Eng. Bull.*, vol. 36, no. 1, pp. 56–64, 2013.
- [30] “XML Tutorial,” <http://www.w3schools.com/xml/>.
- [31] E. P. Carman, Jr. *et al.*, “Apache VXQuery: A Scalable XQuery Implementation,” *CoRR*, vol. abs/1504.00331, 2015.
- [32] “Apache MRQL,” <http://mrql.incubator.apache.org/>.
- [33] M. Kay, “Ten Reasons Why Saxon XQuery is Fast,” *IEEE Data Eng. Bull.*, vol. 31, no. 4, pp. 65–74, 2008.