AQuery: Query Language for Ordered Data, Optimization Techniques, and Experiments^{*}

Alberto Lerner[†]

Dennis Shasha

Ecole Nationale Superieure de Telecommunications Paris - France lerner@cs.nyu.edu New York University New York - USA shasha@cs.nyu.edu

Abstract

An order-dependent query is one whose result (interpreted as a multiset) changes if the order of the input records is changed. In a stock-quotes database, for instance, retrieving all quotes concerning a given stock for a given day does not depend on order, because the collection of quotes does not depend on order. By contrast, finding a stock's fiveprice moving-average in a trades table gives a result that depends on the order of the table. Query languages based on the relational data model can handle order-dependent queries only through add-ons. SQL:1999, for instance, has a new "window" mechanism which can sort data in limited parts of a query. Add-ons make order-dependent queries difficult to write and to optimize. In this paper we show that order can be a natural property of the underlying data model and algebra. We introduce a new query language and algebra, called AQuery, that supports order from-theground-up. New order-related query transformations arise in this setting. We show by experiment that this framework - language plus optimization techniques - brings orders-ofmagnitude improvement over SQL:1999 systems on many natural order-dependent queries.

Proceedings of the 29th VLDB Conference, Berlin, Germany, 2003

1 Introduction

Querying ordered data arises naturally in applications ranging from finance to molecular biology to network management. A financial analyst may be interested in moving averages within or correlations among price time series. A biologist may be interested in frequent nucleic acid motifs. A network manager may be interested in packet flow statistics. Several extensions to SQL have been suggested that are able to express such order-dependent queries [16, 2, 14, 11].

SQL:1999, through its OLAP amendment, is the first such language to gain commercial acceptance [11]. It provides this facility through the following new order-aware features: ordering in the SELECT clause (OVER...WINDOW construct), a notion of row numbering, and an ARRAY data type. Unfortunately, these extensions, however expressive, result in complex formulations of even simple queries. The complex formulation is then difficult to optimize.

1.1 Motivational Queries and Problems

Consider the schema Trades(ID, tradeDate, price, ts), where ID is the ticker symbol of a traded security, ts is short for timestamp, which identifies the date and time of a particular trade, tradeDate is a human-readable form of the day portion, and price is the price of the trade.

Consider the following query: for a given stock and a given date, find the best profit one could obtain by buying the stock and then selling it later that day (short selling – in which an item is sold before it is bought – is disallowed). Algorithmically, the solution

Work supported in part by U.S. NSF grants IIS-9988636 and N2010-0115586.

This work was done while this author was visiting NYU

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the VLDB copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Very Large Data Base Endowment. To copy otherwise, or to republish, requires a fee and/or special permission from the Endowment.

is straightforward: compute the profit resulting from selling at each time t by subtracting the price at t by the minimum price seen up until t. The answer to the query is the maximum of these profits. In SQL:1999 this query would look like:¹

[SQL:1999]

SELECT max(running_diff) FROM (SELECT ID, tradeDate, price - min(price) OVER (PARTITION BY ID, tradeDate ORDER BY ts ROWS UNBOUNDED PRECEDING) AS running_diff, FROM Trades) AS t1 WHERE ID = 'ACME' AND tradeDate = '05/11/03'

The nesting here is necessary because a windowed function (min(price) OVER ...) cannot be an argument of an aggregating function (max(running_diff)). One opportunity for optimization is to push the outer query's selection (ID='ACME' AND trade-Date='05/11/03') to the inner query. Although possible in this particular case, pushing a selection over a projection containing a generic expression involving windowed functions (SELECT ... price - min(price) OVER ...) requires deep analysis. As of this writing, the commercial optimizers we have tested do not do such an optimization.

As another example, consider the schema Packets(plD, src, dest, length, ts), where pID identifies a packet exchanged between a source (src) and a destination (dest) host. Length refers to the size of the packet and ts to the moment (timestamp) this packet was exchanged. A "flow" from a source s to a destination d ends if there is a 2-minute gap between consecutive packets from s to d [3]. In SQL:1999, a network administrator would issue the following query to know the count of packets and their average length within each flow.

[SQL:1999] WITH Prec (src, dest, length, ts, ptime) AS (SELECT src, dest, length, ts, min(ts) OVER (PARTITION BY src.dest ORDER BY ts **ROWS BETWEEN 1 PRECEDING** AND 1 PRECEDING) FROM Packets). Flow (src, dest, length, ts, flag) AS (SELECT src, dest, length, ts, CASE WHEN ts-ptime > 120 THEN 1 ELSE 0 END FROM Prec) FlowID (src, dest, length, ts, fID) AS (SELECT src, dest, length, ts,

	sum(flag) OVER
	(ORDER BY src, dest, ts
	ROWS UNBOUNDED PRECEDING)
FROM	Flow)
SELECT	<pre>src, dest, avg(length), count(ts)</pre>
FROM	FlowID
GROUP	BY src, dest, fID

Basically, this query needs to group packets into flows, and to count the number and average the length of packets within each flow. Finding the flows is very hard to express, though, because it involves order. The first sub-query, Prec, creates a new column, ptime, containing the previous packet's timestamp within each source and destination. Next, the Flow sub-query adds a flag column that is set to true (1) at each packet whose difference from the preceding one exceeds two minutes; otherwise the flag is set to false (0). Next, the FlowID sub-query sums these flags cumulatively, creating an auxiliary flow ID, fID. The main query uses these results.

Optimization of this query should seek to reduce the work required by PARTITION BY and ORDER BYs. That is hard because the windows defined in Prec and in FlowID have slightly different sliding parameters. The commercial optimizer we tested performed this query with two sorts before the grouping was done. Thus they did not perform these optimizations.

The problem with these queries is then two-fold. Their expression in SQL:1999 is complex and therefore both hard to read and difficult to optimize.

We have designed a data model, language, and system where expressions dependent on order are natural to write, fostering the exposure of idioms that the optimizer can exploit.

1.2 AQuery: First Look

In our data model, tables are not viewed as multisets, but rather as ordered entities that we call *arrables* (standing for array-tables). An arrable's ordering may be defined at creation time using an ORDERED BY clause and can later be altered.

Our query language, AQuery, is a semantic extension of the multiset relational model (i.e., SQL 92) that includes the classical SELECT-FROM-WHERE-GROUP BY-AGGREGATE-HAVING clauses. The main extensions are based on a new clause called ASSUMING ORDER which defines the order of the arrables identified in the FROM clause. Predicates and expressions, in whichever clause they are, can count on the order defined by the ASSUMING clause, leading to the natural expression of order-dependent queries as we show in section 2.

It is up to the optimizer to match the input arrables' existing order with a query's ASSUMING ORDER and to decide whether and when further sorting is necessary, as we show in section 3. This flexibility stems

¹The queries shown here use somewhat advanced SQL:1999 features. The unfamiliar reader is encouraged to refer to [11].

from the fact that in the AQuery algebra, each operator has an order-preserving and an order-cavalier variant. The transformations we suggest here may move the sort over other operators, possibly entailing a change to an alternate variant. We show that this schema is able to integrate new transformations with classical ones.

Our experiments, in section 4, show orders of magnitude differences between AQuery's and current commercial SQL:1999's renditions of queries and show that the transformations generate highly efficient plans.

In section 5, we identify several languages that have also considered order as a first-class concept and from which AQuery draws inspiration. In the same section, we also comment on the sources for some of our optimization techniques.

Finally, section 6 summarizes our contributions and describes future work.

2 Data Model and Algebra

2.1 Arrables and Order

Definition 1 (Arrables). Let \mathcal{T} be a set of types in which each $t \in \mathcal{T}$ corresponds to a basic type (e.g., integer, boolean, etc) or to a one-dimensional array made up of elements of a basic type. Let A be a finite, unbounded array of elements of a type $t \in \mathcal{T}$. The cardinality of A is its number of elements. The k-th element of A is denoted by A[k], and k is said to be an *index* or *position* in A. Indexes start at 0. An *arrable* r is a collection of named arrays A_1, \dots, A_n that have the same cardinality, and such that each A_i , $1 \leq i \leq n$, is of a type of \mathcal{T} . \Box

The Figure 1 shows examples of two well-formed arrables. Their schema corresponds to the table Trades described in Section 1.1. Observe that if A_1, \dots, A_n are all vectors (i.e., their elements are all scalars), arrables have the appearance of tables as we know them. That is the case for the arrable Trades in that figure. We will show shortly how arrables having vector elements, such as the arrable Series in the same figure, can be useful.

	Trades	ID	tradeDate	price	ts		
		ACME	05/11/03	12.02	1		
		WXYZ	05/11/03	43.23	2		
		ACME	05/11/03	12.04	5		
		ACME	05/11/03	12.05	9		
		WXYZ	05/11/03	43.22	13		
Series	ID	tradeDate	pric	e		ts	
	ACME	05/11/03	[12.02 1	2.04 12.0	05]	[1 5 9]	
	WXYZ	05/11/03	[43.23 4	3.22]		[2 13]	1

Figure 1: Example of two well-formed arrables

Definition 2 (Arrable Indexing). The *k*-th record

of an arrable r is formed by the k-th element of each of r's component arrays. This operation, denoted *indexing*, is represented as $r[k] = \langle A_1[k], \dots, A_n[k] \rangle$. \Box

For instance, Trades[0] corresponds to the record $\langle ACME, 05/11/03, 12.02, 1 \rangle$, Trades[1] to $\langle WXYZ, 05/11/03, 43.23, 2 \rangle$, and so on.

Because an arrable consists of arrays and arrays are ordered, an arrable is ordered.

Definition 3 (Ordered by). An arrable r may be (lexicographically) ordered by a subset of its arrays, $B_1, \dots, B_m \subseteq A_1, \dots, A_n$. If the ordering is ascending and k_1 and k_2 are two indexes of r and $k_1 < k_2$, then either (i) $B_1[k_1] = B_1[k_2], \dots, B_m[k_1] = B_m[k_2]$ or (ii) there exists a $i, 1 \leq i \leq m$, such that $B_i[k_1] < B_i[k_2]$ and if i > 1 then $B_1[k_1] = B_1[k_2], \dots, B_{i-1}[k_1] = B_{i-1}[k_2]$. The definitions are symmetric for descending orders, but for purposes of exposition, we will consider order to be ascending throughout this paper. \Box

For instance, the arrable Trades of Figure 1 could be defined as Trades(ID, tradeDate, price, ts) ORDERED BY ts.²

Definition 4 (Order-Equivalence). Let r and s be arrables over the same set of attributes. Suppose that r is ORDERED BY some attributes X_1, \dots, X_p , and s by Y_1, \dots, Y_q . Then r and s are *order-equivalent* with respect to attributes B_1, \dots, B_m , denoted $r \equiv_{B_1,\dots,B_m}$ s, if the following conditions hold: (i) r and s are multiset-equivalent (i.e., there exists a permutation of rows P^1, P^2 such that $P^1(r) = P^2(s)$); (ii) B_1, \dots, B_m is a prefix of both X_1, \dots, X_p and Y_1, \dots, Y_q . When r and s are simply multiset-equivalent, we say that $r \equiv_{\{1\}} s$. \Box

2.2 Column-Oriented Semantics

One problem in expressing order-dependent queries is that each resulting row may depend on a combination of values from more than one input row. For example, consider the Trades table and a query to find the difference between each price and its previous value, assuming a time order. The query needs to access two prices at once that are in distinct rows to calculate the pairwise difference. *Row-oriented languages* such as SQL can iterate over only one row at a time, though. Thus they need to resort to either a self-join or to an auxiliary construct to build a row that contains both prices. This operation has to be repeated for each pair.

In contrast, AQuery adopts a *column-oriented se*mantics in that a variable is bound to an entire array. Because variables in AQuery always refer to arrays, expressions always define mappings from a list of arrays to an array. For instance, the above pair-wise difference can be captured by a simple expression – price – prev(price). The function prev() over an array A is an

 $^{^2\}rm We$ are omitting the typing information here for the sake of simplicity. A complete definition would include also NULL and referential integrity information.

array such that $\operatorname{prev}_A[i] = A[i-1]$ if i > 0 and A[0] if i = 0. For two arrays A and B such that |A| = |B|, minus (-) is element-wise subtraction.

The function prev() is a sample of the set of *vector*to-vector functions that AQuery includes. These functions are classified according to their dependency on the input's array sort order and on the cardinality of the output they generate. For instance, prev() is orderdependent and size-preserving. The latter property indicates that it outputs vectors that have as many elements as the input array. Formally order-dependency can be defined as follows.

Definition 5 (Order-Dependency). An expression *e* that maps a list of arrays to an array is said to be *order-independent* if for all operand arrays A_i , $1 \leq i \leq m$, where m is the degree of the expression, and for any corresponding permutations A_i^{perm} , $e(A_1, \dots, A_m) \equiv_{\{\}} e(A_1^{perm}, \dots, A_m^{perm})$. For example, avg(price) is order-independent. An expression that is not order-independent is *order-dependent*. For example, price - prev(price), is order-dependent. \Box

Other functions in the order-dependent, sizepreserving category are the running aggregates. A running minimum over an array A, mins(A), is mins $_A[i] =$ min $(A[i], \min_A[i-1])$ for 0 < i < |A| or A[i] for i =0. Running aggregates use this "s"-as-suffix pattern. A running sum over an array A, denoted sums(A), is sums $_A[i] = A[i] + \text{sums}_A[i-1]$ for 0 <i < |A|, or A[i] for i = 0. Some running aggregates can be computed over sliding windows. For instance, a running average using a fixed-sized window of w positions over an array A is denoted avgs(w, A) and is defined as $\operatorname{avgs}_{w,A}[i] = \operatorname{sum}(A[i-(w-1)]..A[i])/w$, for $w-1 \leq i < |A|$ or $\operatorname{sum}(A[0]..A[i])/i$ for $0 \leq i < w-1$.³

Another category of vector-to-vector functions are those that are order-dependent but not sizepreserving. If they retain either the beginning or the end of an array, they are called *edge functions*. For instance, the first n positions of an array A, denoted first(n, A), is first $_{A,n} = A[0..n - 1]$. Similarly, $last_{A,n} = A[|A| - n..|A| - 1]$.

The classic SQL aggregate functions (min, max, avg, count) can be seen as non-order-dependent, non-size-preserving vector-to-vector functions.

2.3 An Algebra and Query Language

The AQuery algebra supports the operators of the relational algebra. But here each operator takes arraytyped expressions as arguments. If an expression is order-dependent, then the operator behaves in an order-preserving way. Otherwise the operator behaves in an order-cavalier way. The order-cavalier variant of an operator is simply one that is multiset equivalent to its order-preserving variant. In the remaining of the section we define the order-preserving variants of the relational algebra operators.

Definition 6 (Projection). Let r be an arrable and $e = e_1, \dots, e_m$ be a list of expressions involving r's arrays, such that $|e_1| = \dots = |e_m|$. An order-preserving projection of r over e, denoted $\pi_e^{op}(r)$, is defined as follows.

projection(e,r)

- 1. s:= empty arrable having the same schema as e
- 2. for i = 0 to |r|-1
- 3. append $\langle e_1[i], \cdots, e_m[i] \rangle$ to s
- 4. end for
- 5. output s

As mentioned before, if any e_i is order-dependent, the projection is said to be order-preserving, otherwise the projection is order-cavalier, denoted simply $\pi_e(r)$. \Box

Definition 7 (Selection). Let r be an arrable and p be a predicate mapping a list of r's arrays into an array of booleans, such that |r| = |p|. An order-preserving selection of r over p, denoted $\sigma_p^{op}(r)$, is defined as follows.

selection(p,r)

- 1. s:= empty arrable having the same schema as r
- 2. for i = 0 to |r|-1
- if p[i] is true
- 4. append r[i] to s
- 5. end if
- 6. end for
- 7. output s

As for a projection, a selection can be order-dependent, and either order-preserving or order-cavalier. \Box

An arrable's sort order is a property that can be manipulated on a per-query basis.

Definition 8 (Sort). Let $r(A_1, \dots, A_n)$ be an arrable and $B_1, \dots, B_m \subseteq A_1, \dots, A_n$. By a sort of r over B_1, \dots, B_m , we mean a permutation s of r that is OR-DERED BY B_1, \dots, B_m . \Box

Having defined these operators, it is now possible to show the AQuery's rendition of the best-profit query.

```
SELECT max(price - mins(price))
FROM Trades
ASSUMING ORDER ts
WHERE ID = 'ACME' AND tradeDate = '05/11/03'
```

AQuery clauses (SELECT, FROM, ...) are processed in the same order as SQL's. Semantically, ASSUMING ORDER is translated to a sort after the FROM clause is computed. It enforces the desired order for the query and it obliges future clauses (WHERE, GROUP BY, HAVING, and SELECT) to

³Such a definition is commonplace in financial applications. Other domains may require avgs() to return NULLs on positions where the window is incomplete. In any case, it is often convenient to have the running average return an array the same size as its argument.

be translated to order-preserving algebra variations. (This is the required semantics. Optimization may avoid performing the sort this early as we show later.)

Note that due to the column-oriented semantics of AQuery, the mins() function is called only once and takes the whole price vector as an argument. Sub-tracting a vector (mins(price)) from another (price) with the same cardinality is a standard array expression [1], as is taking the max() of the resulting vector.

The above query is translated to the AQuery algebra as follows, letting $e = \max(\text{price} - \min(\text{price}))$ and $p = (\text{ID} = '\text{ACME'}) \land (\text{tradeDate} = '05/11/03')$, as:

$$\pi_e^{op}(\sigma_p^{op}(\text{sort}_{ts}(\text{Trades}))))$$

Grouping in AQuery uses an arrable's facility to store array-valued fields. Intuitively, a grouping operation partitions the operand arrable into disjoint subarrables that share the same group value. It then transforms each sub-arrable into a single row by replacing each non-grouped column (in the sub-array) by its equivalent array-typed value. For instance, the arrable Series in Figure 1 shows the effect of grouping the arrable Trades in the same figure by ID and tradeDate.

Definition 9 (Grouping). Let r be an arrable and $g = G_1, \dots, G_m$ be a list of expressions over r's arrays such that $|G_1| = \dots = |G_m| = |r|$. That is, to each r[i] there must exist a group characterized by g[i]. The order-preserving group-by of r over g, denoted gby_g^{op} , is defined as follows.

group-by(g,r)

```
1. groups := empty arrable having the same schema as g
2. s:= empty arrable having the same schema as r
3. for i = 0 to |r|-1
4.
      if g[i] in groups
5.
          j:= index of g[i] in groups
6.
          for each column C in r
7.
             if C is not a grouped-by column
                 concat r[i].C to s[j].C
8.
             end if
9
10.
          end for
11.
      else
12.
          append g[i] to groups
13.
          append r[i] to s
14.
      end if
15. end for
16. output s
```

Step 13 above forms a single element list (or equivalently a vector). Step 8 concatenates to that list. The result is that fields may consist of vectors. As before, group by is order-dependent if any of its grouping expressions is. Group-by can also have an order-cavalier variation. Further, it is convenient to have an order-generating version of group-by. Semantically, such a group by delivers the results ordered by the grouping expression. \Box

The network management query benefits from AQuery's order-dependent grouping, as Figure 2 depicts. Recall that this query involves a group-by over source host, destination host, and a flow ID between them. Flow ID is order-dependent – a new flow between a pair of hosts starts whenever there is a 120-seconds gap between consecutive packets. In AQuery, such a grouping expression corresponds to src, dest, sums(deltas(ts)>120). Figure 2(a) shows how this expression is computed, supposing the Trades arrable is sorted over src, dest, and ts and that the boolean TRUE carries a value of 1, and the FALSE, 0.

Grouping and aggregation are independent operations in AQuery. The arrable we see in 2(b) shows the result of the grouping operation alone. Note that the non-grouped columns of Packets have arrays within fields. Because fields may be arrays (though not arrables), aggregate functions may apply over an entire column or over each field. To express the latter, AQuery provides an operator modifier called *each* that applies functions to each array-valued element of a column.

Definition 10 (Each Modifier). Let the array A be a parameter of a function F. The execution of F modified by 'each' is defined as follows:

each(F, A)

- 1. B := empty array of the same type of F's result
- 2. for i = 0 to |A|-1
- 3. append F(A[i]) to B
- 4. end for
- 5. output B

This definition is naturally extended for cases where F takes more than one argument. An "each-ed" operator is necessarily an order-preserving one. \Box

Each is a way of applying an operator to each field of a grouped column. In figure 2(c) we see that avg() was applied to each of the array-values of the column length, and similarly to count(ts). Having defined the operators involved in the network management query, we can now show its AQuery rendition.

```
      SELECT
      src, dest, avg(length), count(ts)

      FROM
      Packets

      ASSUMING ORDER src, dest, ts

      GROUP
      BY src, dest, sums(deltas(ts) > 120)
```

The algebraic version of the network management query, supposing that $e = \operatorname{src}$, dest, $\operatorname{each}(\operatorname{avg}(),\operatorname{length})$, $\operatorname{each}(\operatorname{count}(),\operatorname{ts})$ and $g = \operatorname{src}$, dest, $\operatorname{sums}(\operatorname{deltas}(\operatorname{timestamp}) > 120)$, looks like the following. We mark with a corresponding superscript the operations that have components modified by each.

 $\pi_e^{each}(gby_q^{op}(\text{sort}_{src,dest,ts}(\text{Packets})))$

Cross-product (\times) in AQuery is order-cavalier and hence has the same definition as in the relational algebra. By contrast, joins have both order-preserving and order-cavalier variations.



Figure 2: Grouping Trades over src, dest, sums(deltas(ts) > 120): (a) computing the groups; (b) replacing the groups by single records; and (c) aggregating using the each modifier

Definition 11 (Join). Consider the arrables $r(A_1, \dots, A_n)$ and $s(B_1, \dots, B_m)$. A *left-right order*preserving join of r and s over the join predicate p, denoted $r \bowtie_n^{lrop} s$, is defined in the following way.

join(p, r, s)

```
1. o:= empty arrable with schema \langle A_1, \cdots, A_n, B_1, \cdots, B_m \rangle
2. for i = 0 to |r| - 1
```

3. for j = 0 to |s| - 1

- 4. if p(r[i],s[j]) is true
- 5. append $\langle A_1[i], \cdots, A_n[i], B_1[j], \cdots, B_m[j] \rangle$ to o
- 6. end if
- 7. end for

8. end for

9. output o

A query's order may require that only one of the join operand arrables' order be preserved. In that case a simpler order-dependent variation of the join can be used. A *left order-preserving* join, $r \bowtie_p^{lop} s$, is one that is order-equivalent with respect only to A_1, \dots, An to a left-right order-preserving join of the same two arrables. \Box

Consider the arrable Portfolio(ID, tradedSince) ORDERED BY ID, which stores information about the stocks that make one analyst's portfolio. It is a subset of the stocks that appear in Trades. To extract the ten last quotes for each stock in the portfolio, then the following query does the job.

SELECT FROM	t.ID, last(10, price) Trades t, Portfolio p ASSUMING ORDER ts
WHERE	t.ID = p.ID
GROUP BY	t.ID

Semantically, the query first performs a crossproduct (\times) between Trades and Portfolio. Crossproduct in AQuery is order-cavalier. Next, the AS-SUMING clause imposes the desired sort order and the join predicate is applied. Note that the order is imposed on the result of the Cartesian product. Then, the resulting arrable is partitioned into groups according to ID values. The assumed order is preserved within each group. The last() function "trims" each array-valued priced column to a maximum of the ten last positions of each price array. Letting e= ID, each(last(),10,price) and p= Trades.ID=Portfolio.ID, this query can be represented as:

$$\pi_e^{each}(gby_{ID}^{op}(\sigma_p^{op}(\text{sort}_{ts}(\text{Trades} \times \text{Portfolio}))))$$

Our purpose here was to show that by incorporating order from the ground up, AQuery can express orderdependent queries naturally. More advanced arrable manipulation features are described in [10].

3 Query Transformations

Sort elimination [17] and sort move-around [15] are known techniques that can be applied to AQuery optimization. AQuery's semantics and, in particular, the edge built-in functions (e.g., first, last) allow other aggressive order-related optimizations as well. We introduce these new techniques through examples.

3.1 Implicit Selections and Sort-edge

Let Connections(host, port, client, timestamp) OR-DERED BY timestamp be an arrable that stores the clients' addresses that accessed a network's services (port, host) and when did they do so. Suppose one wants to find the last client that connected to server "atlas." In AQuery:

SELECT	last(1, client)
FROM	Connections
	ASSUMING ORDER timestamp
WHERE	host = 'atlas'

We are going to show plans in the usual diagrammatic way. The above query's initial plan is depicted in Figure 3(a). We introduce some auxiliary notation as follows. A single arc between a pair of operators means that the producer operator is outputting records in an order-cavalier fashion (i.e., in the most efficient or simple way possible, without guaranteeing any order). Double-arcs mean it is doing so in an order-preserving way. Arrows represent the net effect of the application of a transformation. Each arrow is annotated with the



Figure 3: Implicit selection and sort-edge optimizations

corresponding transformation number. The formal descriptions of the transformations are given in Table 1. We say pos(r) = i when we refer to the record r[i]. The special indexes for an arrable r, FIRST and LAST, are 0 and |r|-1, respectively. Finally, order(r) returns the list of attributes the arrable r is ORDERED BY.

A regular selection such as $\sigma_{host='atlas'}$ can be pushed down over a sort [15]. Transformation 2 in Table 1 is a slight variation of that in [15] where orderpreservation is made explicit. There are two advantages to commuting the sort and the selection here: the selection can benefit from the existing order on Connections (host), and delaying the sort reduces the number of records that would need to be sorted.

The projection $\pi_{last(1,client)}$ includes an implicit selection, i.e., it is interested in only one client. This is a peculiarity of AQuery's column-oriented semantics – a projection over a function that itself does a selection. The transformation 3 in Table 1 is a new transformation that replaces a projection with indexing by a pure projection plus a selection of the desired positions. If such positions are on one end of the operand array, we call this selection an *edge selection*. The result of applying this transformation is seen in Figure 3(b).

The advantage of isolating the edge-selection from the original projection is that while the latter can't be moved around easily, the former can. In this example, the existence of an edge selection after a sort suggests that there is no need to sort all the input just to use some of the elements. In AQuery, the physical operator *sort-edge* implements the logical pattern $\sigma_{edge-condition}^{op}(sort(r))$. The sort-edge uses a modified heapsort to keep the top (or bottom) *n* elements, as appropriate. This is similar to the approach used in [2] except that we modify the heapsort to make it stable.⁴ The final plan is shown in Figure 3(c).

3.2 Sort Splitting

There are situations in which the arrable's existing order facilitates the evaluation of part of a query, even though it does not match the query's ASSUMING ORDER. Consider again the arrable Connections OR-DERED BY host. The following query finds all the



Figure 4: Sort-splitting optimization

clients that connected to the last host to which a client hooked in.

SELECT	client
FROM	Connections
	ASSUMING ORDER timestamp
WHERE	host = last(1,host)

Note that the predicate host = last(1,host) makes sense as an array expression. The array host is compared to the single element array (treated as a scalar) last(1, host) resulting in an array of booleans. Positions that map to false are eliminated by the WHERE clause.

An initial plan for this query appears in Figure 4(a). Timestamp is not a prefix of order(Connections), thus the sort over timestamp is required. However, host is a prefix of order(Connections), and therefore the selection $\sigma_{host=last(1,host)}$ may take advantage of it. This is where the split-sort technique comes in.

If A and B are arrays of an arrable r, a selection $\sigma_{A=(B[i])}(r)$ can be replaced by a semi-join as described by transformation 6 in table 1. The benefit of the semijoin is that we can now manipulate order on each of the semi-join's arguments differently.

Figure 4(b) shows the result of applying that transformation. Note that last(1, host) = host[LAST]. Let's analyze each side of the semi-join in turn. On the right-hand side we have the pattern edge-selection / sort, which can be efficiently implemented, as we have discussed. By contrast, the left-hand-side sort changes what could be an interesting order to the semi-join operation. We can thus defer it until after the join. The transformation 5 in Table 1 commutes a semi-join and a sort. This transformation can be derived from [15] and says that sorting the result of a semi-join is equivalent to sorting its left stream and then performing an order-preserving semi-join assuming the conditions stated in Table 1 hold. The net effect here is that computing the semi-join predicate is facilitated by an existing order and that sorting over the timestamp needs to be done only for the records generated by the semijoin – much cheaper than the original join over the whole arrable. The resulting plan appears in Figure 4(c).

A contrasting technique where a big sort is exchanged by a number of smaller ones is described next.

⁴A stable sort is one that does not change the original order of records having identical value on the sorted key. Heapsort is not naturally stable. It becomes stable if one concatenates a tuple ID to the key.

Sort Reduction/Elimination	
(1) $\operatorname{sort}_A(r) \stackrel{'}{\equiv}_{order(r)} r$	if A is a prefix of $order(r)$
Selection	
(2) $\sigma_p^{op}(\operatorname{sort}_A(r)) \equiv_A \operatorname{sort}_A(\sigma_p(r))$	if p is <i>not</i> order-dependent
Projection	
$(3) \pi_{e[i]}^{op}(r) \; \equiv_{_{order(r)}} \; \pi_{e}^{op}(\; \sigma_{pos()=i}(r) \;)$	e is an expression over r 's arrays
Join and Semi-Join	
(4) $\operatorname{sort}_A(r \bowtie_{A=B} s) \equiv_A \operatorname{sort}_A(r) \bowtie_{A=B}^{lop} s$	if $A, B \in$ schema of r, s , resp.
(5) $\operatorname{sort}_A(r \Join_{A=B} s) \equiv_A \operatorname{sort}_A(r) \Join_{A=B}^{lop} s$	if $A, B \in$ schema of r, s , resp.
$(6) \sigma_{A=(B[i])}^{op}(r) \equiv_{order(r)} r \Join_{A=B}^{lop} \sigma_{pos()=i}(r)$	if $A, B \in$ schema of r
(7) $\sigma_p^{op}(r \bowtie_{A=B}^{lop} s) \equiv_{order(r)} \sigma_p^{op}(\sigma_p^{each}(gby_A(r)) \bowtie_{A=B}^{lop} s)$	if $A, B \in$ schema of r, s , resp.,
	p is 'pos() = FIRST' or 'pos() = LAST',
	and B is unique
Group-By	
(8) $\operatorname{gby}_{A}^{op}(\operatorname{sort}_{A,B}(r)) \equiv_{A,B} \operatorname{sort}_{B}^{each}(\operatorname{gby}_{A}^{og}(r))$	

Table 1: A subset of the equivalences between sort and remaining algebra operators

3.3 Sort Embedding

Consider the arrable Trades(ID, tradeDate, price, timestamp), this time with no determined ORDERED BY. (Often trades arrive in a "near-timestamp" order.) The query seeks the ten most recent prices for each security ID. In AQuery:

SELECT	ID, last(10,price)
FROM	Trades
	ASSUMING ORDER ID, timestamp
GROUP	BY ID

An initial plan for this query is shown in Figure 5(a). We can separate the implicit selection from the projection as we did before. The resulting plan appears in Figure 5(b).



Figure 5: Sort-embedding optimization

It is possible to delay sort until after the GROUP BY ID is done. If delayed, sort would have to be applied only within each group. This is what we call sort embedding. Moreover, for this particular query the smaller sorts would be then followed by edge selections – sort-edge would apply. The transformation 8 in table 1 allows commuting a sort with a group-by. Note that (a) the group by results must be ordered over the grouping list (i.e., an order-generating operator); and (b) grouping must be over a prefix of sort's arguments. The result of this transformation is shown in Figure 5(c). Note how a double-arc connects groupby and sort-each, because this instance of group by is order generating.

3.4 Edgeby and Early Edge Selection

Let us look at another scenario where an edge selection may reduce cardinality early in a query. We use the arrable Trades here as well, but we now assume it is ORDERED BY timestamp. The arrable Portfolio(ID, name, tradedSince) ORDERED BY ID stores the subset of securities with which an analyst deals. Name is a unique identifier of securities in Portfolio, and so is ID. To retrieve the last price of a security named "DataOrder" one would do:

SELECT	last(1, price)
FROM	Trades, Portfolio
	ASSUMING ORDER timestamp
WHERE	Trades.ID=Portfolio.ID
	AND name = "DataOrder"

An initial plan for this query is depicted in Figure 6(a). A common heuristic to improve performance is to push down the regular selection over the sort, and further over the join. The transformation 2 in Table 1, which comes from [15] allows us to do so. The result is seen in Figure 6(b).

Upon detection of an appropriate existing order (i.e., order(Trades) matches the ASSUMING ORDER of the query), the optimizer would try to eliminate the sorting completely. Transformation 4 in table 1 commutes a join with a sort while still keeping track of order. That is a slight variation of a transformation in [15] in which order-preservation is made explicit. As Trades is already ORDERED BY timestamp, that sort may be eliminated, as transformation 1 in Table 1 determines. This transformation comes from [17]. The result is seen in Figure 6(c).



Figure 6: Early edgeby optimization

This query also contains a projection-withselection, and again we can break them apart. The consequent presence of the edge selection after the join suggests that we may not need to perform the edge selection in its entirety. Portfolio.ID is a key and so it guarantees that each record in Trades will match at most one record in Portfolio. (Foreign key joins are among the most frequent of equijoins.) In these conditions we could push down this edge selection in the following way. For each ID in Trades, find its last record. This can be done by grouping Trades by ID and edgeselecting each last record. Because of the edgeby cardinality reduction, the join would be performed over far fewer records. The final selection would then pick the desired price. This is what transformation 8 in table 1 does. The final plan is shown in Figure 6(d).

An edge-selection applied to groups is an idiom, called *Edgeby*, that can be highly optimized. Edgeby is a physical operator capable of implementing the logical pattern $\sigma_{edge-condition}^{each}(\text{Gby}(r))$. Instead of separating all elements of an arrable into groups just to use a slab of them (e.g., first n, last n, drop n, etc), edgeby discards, on-the-fly, elements for groups that already fulfill the edge condition. In our example, we need the end of an array (i.e., last() over ascending order) of prices for each ID. A backwards edgeby ID on the sorted Trades keeps a record if it belongs to a heretofore unseen ID or a group having less than ten records.

4 Experimental Results

To evaluate the relative performance of (a) AQuery vs. SQL:1999 queries and (b) syntactic-directed translations vs. optimized plans for AQuery queries, we have conducted an experimental study.

All our experiments were executed on a Pentium III-

M 1.13Mhz with 1 Gb of memory running Linux with no special setting of process priorities. The timings reported here correspond to wall clock timings.

4.1 Performance Measurements

AQuery's rendition of the best-profit and the networkmanagement queries contained no nesting and were straightforward to optimize. The figure 7 shows the relative performance improvement of AQuery plans versus an SQL:1999 commercial optimizer's. AQuery results were between eight and twenty one times faster for the best-profit query, and between two and three time faster for the network management query.



Figure 7: AQuery's vs. SQL:1999 Relative Performance

The best-profit query numbers were generated using Trades arrables/tables with varying number of securities from 200 to 1000, and using 1000/trades per security. Recall that this query was interested in profits for a given security for a given date. The time difference is due to AQuery's ability to push down the selection predicate and to use an index to evaluate it. The subsequent sort reorders only the trades for the relevant security. Because the SQL:1999 representation used a complicated nesting structure (see Section 1.1), its optimizer could not move the selection. That plan sorted tuples that ended up being discarded.

For the network management query we used a Packets arrable/table with 100 sessions and varying number of packets for each session from 2K to 10K. AQuery's plan was faster because it required only one sort to be done, the one enforcing the ASSUMING ORDER clause. Group-by in AQuery depends – and thus benefits – from that ordering. By contrast, the SQL:1999 optimizer had to figure out how to deal with two unrelated WINDOWS specifications (see Section 1.1). This resulted in having two distinct sorts before the processing of the group by, which did not benefit from them.

Moral: AQuery's structural simplicity helps in finding better plans.

In most cases, an edgeby requires a small fraction of the time required to perform the associated groupby, if done in its entirety as shown in Figure 8(a). We used the arrable **Trades** with 1 million records divided evenly among 10, 100, 1000, and 10000 securities. An edgeby over security ID with varying slab sizes is tested. The more records edgeby can discard, the faster its response time. For instance, when only a few distinct securities are used, groups are large, and therefore most records fall off the slabs even for the biggest slab sizes tested, greatly improving performance. As the groups get smaller (i.e., more distinct securities are used), highly selective slabs gets better performance. A degenerate case is seen where a 100-slab is taken from groups that are themselves 100 records wide (i.e., 10000 securities). Edgeby doesn't improve performance here – but doesn't hurt either.



Figure 8: Efficiency of work reduction techniques. Vertical axis are time in milliseconds. (a) Several sized edgeby ID slabs of Trades varying the number of distinct securities. (b) Embedding sorts within varying sized groups.

The idea behind the sort embedding technique is that a sort can be delayed until after a group by, and can be replaced by several sorts over the grouped columns (sort-each's). Figure 8(b) characterizes the performance gains of sort-eaches as compared to the entire sort they replace. We used arrables of 1 million records and varied the number of groups in which the arrable was divided. When only one group exists, there's no point in applying the technique – but, again, there's no penalty in doing so. Trading one big sort by several smaller ones starts to payoff whenever more than 10 groups exist.

Sort-edge presented similar results as those of sortstop [2] and we omit these results.

Moral: Reduction of work due to edge selections or sort handling techniques is significant.

AQuery's optimizer integrates sort with the standard relational operators such as selection and group by. For instance, in the query of section 3.1 it was able to apply selection push-down over a sort. Figure 9(a)shows that this technique was particularly efficient for arrable instances where the number of distinct hosts was greater than 10. But plans using a plain, regular sort on instances with less than 10 distinct hosts would still perform poorly. By identifying the implicit edge selection of that query and by using it to reduce the number of records to sort, AQuery generated an optimized plan that performed best in all the instances tested.

Figure 9(b) shows the performance gains of applying the sort splitting technique to the example query of section 3.2. The efficiency of the optimized plan stems from delaying the enforcement of the ASSUMING order up until after the semi-join reduces the number of records to be sorted. The gains stabilized at instances with 100 or more distinct hosts because at this point the cost of the query is dominated by the semi-join itself as opposed to the sort of its results. Note that application of this technique whenever the number of hosts is too low (e.g. just one) may represent an unnecessary overhead – although a small one. So, this technique depends on the data distribution, underscoring the need for cost-based optimization.

Figure 9(c) shows the comparative performance of plans for the example query of section 3.3. The naive plan sorts the whole arrable, groups the entire result and applies the edge-selection only at the end. Cost remains rather high, even when the edge selection removes most records. By contrast, the optimized plan trades one big sort for several smaller ones – sort-edges, in fact. Thus, even in the degenerate case where each group has only one record (i.e., number of distinct hosts is equal to the cardinality of the arrable), the optimized plan saves the cost of a big sort. The curves show order of magnitudes difference at instances with small number of distinct hosts.

Finally, Figure 9(d) shows the results for the naive and the optimized plans for the example query of section 3.4. By applying an edgeby early in the plan, the number of records that have to be joined is considerably reduced. The optimized plan also takes advantage of the existing order, eliminating any sort altogether. The result is consistently faster response times.

Moral: AQuery transformations bring substantial performance improvements, especially when used with cost-based query optimization.

5 Related Work

Whereas SQL:1999 is the commercially most significant implementation of order-dependent queries, other systems, both commercial and research, have proposed many excellent ideas.

5.1 Optimization Techniques

Sort order has always been treated in the optimization process as a physical property to be included in a plan (if not specified by SQL's ORDER BY clause) to support an efficient algorithm like the merge join. Mechanisms such as Starburst's "glue" [9] or Volcano's



Figure 9: Optimized vs. Non-optimized plans. Vertical axis is time in milliseconds. Horizontal axis is the number of hosts for (a) and (b) over which the 1 million connections are divided and the number of securities for (c) and (d) over which the 1 million trades are divided. Plans for (a) are shown in figure 3, for (b) in 4, for (c) in 5, and for (d) in 6.

"enforcer" [5] made sure a sort step was added whenever an efficient algorithm required it.

In [17], the authors added an order-optimization step before plans were enumerated in the context of DB2's optimization process. This step may dramatically improve queries that have order requirements due to the clause ORDER BY, GROUP BY, or the DIS-TINCT modifier. But order optimization in [17] was still a separate optimization step in that transformations involving sort elimination or reduction (i.e., sort over fewer attributes) were not considered at the same time transformations involving other algebraic operators were. By contrast, because we consider sort operators with other operators in the spirit of [15], we are able to discover techniques such as sort splitting or sort embedding (i.e., transformations 6 and 8, respectively, in Table 1).

In [2] a clause STOP AFTER was suggested which was capable of limiting the resulting cardinality of queries whose results were ordered (by an ORDER BY). When only the *top-k* tuples of a query's result need be consumed, queries can run much faster. Our sort-edge is similar to stop-sort. A difference between that work and ours is that, again, their order optimization process happens as a separate, isolated phase. Because we integrate the two, AQuery uses the stop-after idea within groups as well. After all cardinality restriction is a generally useful technique.

The first work, to the best of our knowledge, to provide optimization of order in an integrated framework was [15]. Most of their list-based transformations apply to AQuery. The transformations we presented here extend their framework with several new techniques (e.g., transformations 3 and 7 in Table 1).

Other optimization rules related to order or ordered structures were suggested in the context of their corresponding query languages, which we discuss next.

5.2 Languages

AQuery is a descendent of KX systems's KSQL [7] from which AQuery takes its arrable notion – a fully vertically partitioned implementation of tables, each of whose columns is an array. AQuery differs from KSQL by trying to preserve the SQL flavor to a much greater extent than KSQL, by the use of the ASSUM-ING clause to make the use of order declarative, by the introduction of transformations, and by the exploitation of a cost-based framework for optimization.

The sequence query languages SEQUIN [16] and SRQL [14] are precursors of SQL:1999 in that they are SQL dialects that handle order-based queries. They are true to the spirit of [13] which showed that order was a much needed feature in queries. SEQUIN treats sequences as an extended abstract data type, although a sequence can serve as the sole source of data for a query. SRQL was inspired by SEQUIN and treats tables as ordered relations. AQuery borrowed from these languages the early introduction in a query of an order defining clause. Both SEQUIN and SRQL keep the tuple semantics of SQL, as opposed to the vector (column) processing of AQuery. A consequence is that several valid vector expressions in AQuery are invalid in these languages, e.g. $\max(\text{price} - \min(\text{price}))$. Another difference is in the way non-1NF relations/arrables are handled. In [16], if a table has a sequence attribute, SEQUIN is used to express predicates over the sequence, while SQL is used over the table. That can lead difficultto-read queries and complex optimization transformations. Further, SRQL did not pursue non-1NF ordered relations, whereas we have found them very useful.

Non-SQL efforts at query languages and data models based on arrays have been suggested before [8, 12]. Neither AQL [8] nor AML [12] provides a declarative mechanism to define the order in which the queries manipulate data. Queries process data in the order it is stored. While that makes sense for databases of raster images [12] or scientific data in CDF format [8], it makes less sense in general data processing. Nevertheless, one could argue that these languages could easily incorporate a sort function and express most, if not all, of the queries here. We agree. We would welcome such extensions because the implementors may then arrive at interesting optimizations that would be complementary to ours. Finally, we have a bias for pragmatic reasons for an SQL dialect, but reasonable people can differ on this point.

A particularly inspiring feature of the AQL optimizer is that it has the powerful capability of optimizing operators (or newly added functions) on the calculus level, i.e., by application of variations of λ -calculus reductions over the operators definitions. Reductions help find syntactically simpler forms of an expression while keeping its semantics intact. We have not yet fully exploited that ability in AQuery. On the other hand, we have shown that, for instance, the sort splitting technique requires more than simplifying an expression. It involved transforming what was one sort plus a selection in a semi-join plus two sorts plus a selection – and that resulted in sorting fewer tuples than the simpler expression. AML by contrast uses a fixed set of transformation rules aimed generically at function application on array slabs. A complete fusion of these ideas requires more exploration.

6 Conclusion

AQuery builds on previous language and query optimization work to accomplish the following goals:

- 1. Incorporate order in a declarative fashion to a query language (using the ASSUMING clause) built on SQL 92.
- 2. Introduce a query semantics that, while keeping compatibility to SQL's, supports inter-tuple calculations without query nesting.
- 3. Add order-dependent functions (e.g. sums, first) that are natural to express and flexible, for instance, to allow querying top-k elements within groups.
- 4. Create a simple, yet powerful optimization framework that results in performance that is an order of magnitude faster than commercial SQL:1999 systems for natural queries. Edge optimization and sort splitting and embedding seem to be particularly promising for order-dependent queries.

Acknowledgments

We would like to thank David Tanzer and Jennifer Whelpley for their many helpful comments on earlier drafts.

References

- Budd, T., "An APL Compiler." Springer-Verlag, April, 1988.
- [2] Carey, M.J., and Kossmann, D., "On Saying 'Enough Already!' in SQL." SIGMOD Int'l Conf. on Management of Data, 219–230, 1997.

- [3] Cranor, C., Gao, Y., Johnson, T., Shkapenyuk, V., and Spatscheck, O., "Gigascope: High Performance Network Monitoring with an SQL Interface." SIGMOD Int'l Conf on Management of Data, 623–623, 2002.
- [4] Claussen, J., Kemper, A., Kossmann, D., and Wiesner, C. "Exploiting early sorting and early partitioning for decision support query processing." VLDB Journal, 9(3): 190–213, 2000.
- [5] Graefe, G., and McKenna, W.J., "The Volcano Optimizer Generator: Extensibility and Efficient Search." ICDE Int'l Conf on Data Engineering, 209–218, 1993.
- [6] Knuth, D.E., "The Art of Computer Programming." V3, Sorting and Searching, 2nd Ed,. Addison Wesley Longmand, 1998
- [7] KX Systems. "KSQL Reference Manual, Version 2.0." http://www.kx.com
- [8] Libkin, L., Machlin, R., and Wong, L., "A Query Language for Multidimensional Arrays: Design, Implementation, and Optimization Techniques." SIGMOD Int'l Conf on Management of Data, 228–239, 1996.
- [9] Lohman, G.M., "Grammar-like Functional Rules for Representing Query Optimization Alternatives." SIGMOD Int'l Conf on Management of Data, 18–27, 1988.
- [10] Lerner, A. "Querying Ordered Databases with AQuery." Ph.D. Thesis, Ecole Nationale Superieure de Telecommunications, ENST-Paris, 2003.
- [11] Melton, J., "Advanced SQL:1999 Understanting Object-Relational and Other Advanced Features." Morgan Kaufmann Publishers, September, 2002.
- [12] Arunprasad P. Marathe and Kenneth Salem. "Query Processing Techniques for Arrays." The VLDB Journal, 11:68–91, 2002.
- [13] Maier, D., and Vance, B., "A Call to Order." PODS Symp on Principles of Database Systems, 1–16, 1993.
- [14] Ramakrishnan, R., Donjerkovic, D., Ranganathan, A., Beyer, K.S., and Krishnaprasad, K., "SRQL: Sorted Relational Query Language." SSDBM Int'l Conf on Scientific and Statistical Database Management, 84–95, 1998.
- [15] Slivinskas, G., Jensen, C.J., and Snodgrass, R.T., "Bringing Order to Query Optimization." SIG-MOD Record, 31(2):5–14, 2002.
- [16] Seshadri, P., Livny, M., and Ramakrishnan, R. "The Design and Implementation of a Sequence Database System." VLDB Int'l Conf on Very Large Data Bases, 99–110, 1996
- [17] Simmen, D.E., Shekita, E.J., and Malkemus, T., "Fundamental Techniques for Order Optimization." EDBT Int'l Conf on Extending Database Technology, 625–628, 1996.