

# Thèse

Présentée pour obtenir le grade de  
**Docteur de l'Ecole Nationale Supérieure des  
Télécommunications**

(Spécialité : Informatique et Réseaux)

**Alberto LERNER**

Interrogation de Données Ordonnées avec  
AQuery

Soutenue le 11 juillet 2003, devant le jury composé de :

**Philippe Poucheral**

Président du Jury

**Peter Buneman**

Rapporteurs

**Eric Simon**

**Dennis Shasha**

Directeurs de Thèse

**Talel Abdessalem**



# Doctoral Thesis

Submitted to the Department of Infomatics  
in partial fulfillment of the requirements for the degree of  
**Docteur de l'Ecole Nationale Supérieure des  
Télécommunications**  
(Specialization: Informatics and Networks)

**Alberto LERNER**

Querying Ordered Databases with  
AQuery

Presented at July 11, 2003 to the jury members:

<b>Philippe Poucheral</b> (INRIA)	Jury's President
<b>Peter Buneman</b> (University of Edinburgh) <b>Eric Simon</b> (INRIA)	Rapporteurs
<b>Dennis Shasha</b> (New York University) <b>Talel Abdessalem</b> (ENST)	Thesis Advisors



# Contents

<b>I</b>	<b>Résumé en français - Long abstract in french</b>	<b>11</b>
<b>1</b>	<b>Interrogation de données ordonnées avec AQuery</b>	<b>12</b>
1.1	Introduction . . . . .	12
1.2	Requêtes Motivées et Problèmes . . . . .	12
1.2.1	AQuery: Premier aperçu . . . . .	14
1.3	Modèle de Données et Algèbre . . . . .	15
1.3.1	Arrables et Ordre . . . . .	15
1.3.2	Sémantiques Orientées-Colonne . . . . .	17
1.3.3	Un Langage de Requêtes et une Algèbre . . . . .	18
1.4	Transformations de Requêtes . . . . .	23
1.4.1	Sélections implicites et Tri de bord ( <i>Sort-edge</i> ) . . . . .	23
1.4.2	Partage de tri ( <i>Sort Splitting</i> ) . . . . .	25
1.4.3	Enchâssement de tris ( <i>Sort Embedding</i> ) . . . . .	26
1.4.4	<i>Edgeby</i> (“ Par bords ”) et Sélection de Bord Précoce ( <i>Early Edge Sélection</i> ) . . . . .	27
1.5	Résultats Expérimentaux . . . . .	29
1.5.1	Mesures de Performance . . . . .	29
1.6	Travaux apparentés . . . . .	32
1.6.1	Techniques d’Optimisation . . . . .	32
1.6.2	Langages . . . . .	33
1.7	Conclusion . . . . .	34
<b>II</b>	<b>Mémoire en anglais</b>	<b>37</b>
<b>2</b>	<b>Introduction</b>	<b>43</b>
2.1	Order-Dependent Queries . . . . .	43
2.2	Principles and Goals . . . . .	44
2.3	Thesis Overview . . . . .	45
<b>3</b>	<b>State of the Art</b>	<b>47</b>
3.1	Introduction . . . . .	47
3.2	Standard SQL with Late Order . . . . .	50
3.3	SQL Dialects over Ordered Structures . . . . .	52
3.3.1	SEQUIN . . . . .	52
3.3.2	SRQL . . . . .	53

3.4	Array-Based Querying Systems . . . . .	54
3.4.1	AQL . . . . .	54
3.4.2	KSQL . . . . .	55
3.5	Discussion . . . . .	56
<b>4</b>	<b>AQuery Syntax and Semantics</b>	<b>59</b>
4.1	An Array-Based Data Model . . . . .	59
4.2	Column-Oriented Semantics . . . . .	61
4.3	Relational Manipulation of Arrables . . . . .	62
4.3.1	Projection . . . . .	63
4.3.2	Selection . . . . .	63
4.3.3	Group By . . . . .	64
4.3.4	Flatten . . . . .	66
4.3.5	Cross Product and Join . . . . .	67
4.4	Positional Manipulation of Arrables . . . . .	68
4.4.1	Querying with Arrable Indexing . . . . .	69
4.4.2	Querying with Row Direct Addressing . . . . .	69
4.5	Comparing AQuery to Other Order-Aware Languages . . . . .	70
4.6	Conclusion . . . . .	71
<b>4</b>	<b>AQuery Optimization</b>	<b>73</b>
4.1	Introduction . . . . .	73
4.2	Optimization of Edge Selections . . . . .	76
4.2.1	Implicit Selections and Sort-Edge . . . . .	76
4.2.2	Sort Splitting . . . . .	78
4.2.3	Early Edge Selection and Edgeby . . . . .	79
4.2.4	Sort Embedding . . . . .	81
4.3	Related Work . . . . .	83
<b>5</b>	<b>System Design and Implementation</b>	<b>85</b>
5.1	A Column-Oriented Execution Model . . . . .	85
5.2	Implementing the Execution Model . . . . .	89
5.3	From Text to Execution: the entire flow . . . . .	90
5.3.1	Parsing . . . . .	90
5.3.2	Semantics Step . . . . .	91
5.3.3	Relational Optimization Support . . . . .	91
5.3.4	'K'-Code Generation . . . . .	93
5.4	Conclusion . . . . .	93
<b>6</b>	<b>Performance Analysis</b>	<b>95</b>
6.1	Introduction . . . . .	95
6.2	The Best Profit Query . . . . .	97
6.3	Network Management Query . . . . .	98
6.4	Conclusion . . . . .	101

<b>7</b>	<b>Conclusion</b>	<b>103</b>
7.1	Summary . . . . .	103
7.2	Ongoing Work . . . . .	103
7.3	Future Work . . . . .	104

# List of Figures

1.1	Exemple de deux arrables correctement constituées. . . . .	16
1.2	Regroupement de Trades sur src, dest, sums(deltas(ts) > 120) . .	21
1.3	Sélection implicite et optimisations de tri de bord . . . . .	24
1.4	Optimisation de partage de tri . . . . .	26
1.5	Optimisation d’enchâssement de tris . . . . .	27
1.6	Optimisation précoce “ edgeby ” . . . . .	28
1.7	Les performances relatives de AQuery contre SQL:1999 . . . . .	29
1.8	Efficacité des techniques de réduction de travail. . . . .	30
1.9	Plan optimisé contre plan non-optimisé. . . . .	31
3.1	A Sales table instance and the result of the delta sales query . . .	48
3.2	A stock’s price curve and its running minimum . . . . .	49
4.1	Example of two well-formed arrables . . . . .	59
4.2	Intermediate arrables in the Newtork Management Query . . . . .	66
4.1	An initial QEP and the application of a sort elimination transfor- mation . . . . .	74
4.2	An initial QEP and the application of a selection push-down trans- formation . . . . .	75
4.3	Implicit selection and sort-edge optimization . . . . .	76
4.4	Efficiency of sort-edge technique . . . . .	77
4.5	Sort-splitting optimization . . . . .	78
4.6	Efficiency of sort-splitting technique . . . . .	79
4.7	Early edgeby optimization . . . . .	80
4.8	Efficiency of the early edgeby technique . . . . .	81
4.9	Sort-embedding optimization . . . . .	82
4.10	Efficiency of the sort embedding technique . . . . .	83
5.1	Two example arrables and their rows indexes . . . . .	87
5.2	Effective index array during a query execution . . . . .	87
5.3	Example of a K plan . . . . .	90
5.4	Two possible MEMO configurations for a given query . . . . .	93
6.1	A table to be used on a window definition . . . . .	95
6.2	A running minimum window . . . . .	96
6.3	A previous row window . . . . .	96
6.4	Plans for the best-profit query . . . . .	98

6.5	Best profit query relative improvement . . . . .	99
6.6	Plans for the network management query . . . . .	100
6.7	Network management query relative improvement . . . . .	101

# List of Tables

1.1	Un sous-ensemble des équivalences entre le tri et les opérateurs restants de l'algèbre . . . . .	24
3.1	Comparative table of languages with order constructs . . . . .	58
4.1	Equivalences between sort and remaining algebra operators . . . . .	75
5.1	Cardinality of the addition operation . . . . .	91

# Part I

Résumé en français

Long abstract in french

# Chapter 1

## Interrogation de données ordonnées avec AQuery

### 1.1 Introduction

La recherche de données ordonnées est un problème qui se pose naturellement dans des applications allant de la finance à la biologie moléculaire, en passant par la gestion des réseaux. Un analyste financier peut s'intéresser à des moyennes mobiles ou des corrélations dans des séries chronologiques de prix. Un biologiste peut s'intéresser à des répétitions de motifs d'acides nucléiques. Un gestionnaire de réseaux peut s'intéresser aux statistiques de flux de paquets. Plusieurs extensions au SQL, capables d'exprimer des requêtes dépendant ainsi de l'ordre, ont été proposées [32, 8, 30, 27].

A travers son amendement OLAP, le SQL:1999 est le premier langage du genre à gagner l'approbation commerciale [27]. Il réalise les fonctions décrites précédemment grâce à de nouvelles fonctionnalités qui prennent l'ordre en compte: triage dans la clause SELECT (OVER WINDOW construct), une notion de numérotation des lignes et l'introduction de données de type ARRAY (matrice). Malheureusement ces extensions, bien qu'expressives, ont pour résultat des formulations complexes, même pour des requêtes simples. Cette formulation complexe est alors difficile à optimiser.

### 1.2 Requêtes Motivées et Problèmes

Considérons le schema Trades(ID, tradeDate, price, ts), où ID est le code SICOVAM d'une action échangée, ts est le "timestamp", qui Identifie la date et l'heure d'un échange particulier, tradeDate est un l'affichage de la partie du jour dans un format lisible et enfin price est le prix d'échange.

Considérons la requête suivante: pour une action donnée et pour une date donnée, trouver le meilleur bénéfice que l'on aurait pu obtenir en achetant l'action et en la vendant plus tard dans la journée (la vente à découvert, qui consiste à vendre une action avant de l'avoir achetée, n'est pas autorisée). Algorithmiquement, la solution est simple : il suffit de calculer le profit réalisé en vendant à

chaque instant  $t$  en soustrayant au prix à  $t$  le minimum qu'aient connus les prix avant  $t$ . La réponse à la requête est le maximum de ces profits. Dans SQL:1999, cette requête aurait l'allure suivante :<sup>1</sup>

```
[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'
```

L'imbrication est ici nécessaire puisqu'une fonction fenêtre (`min(price) OVER ...`) ne peut être l'argument d'une fonction d'agrégation (`max(running_diff)`). Une possibilité d'optimisation consiste à déplacer la sélection de la requête supérieure (`ID='ACME' AND tradeDate='05/11/03'`) jusqu'à la requête inférieure. Bien que réalisable dans ce cas particulier, le déplacement d'une sélection sur une projection contenant une expression générique faisant intervenir des fonctions fenêtre (`SELECT ... price - min(price) OVER ...`) requiert une analyse approfondie. A partir de cette expression, les optimiseurs commerciaux que nous avons testés ne font pas de telle optimisation.

Prenons comme autre exemple le schema `Packets(pID, src, dest, length, ts)`, où `pID` représente un paquet échangé entre une source (`src`) et un hôte destinataire (`dest`). `Length` représente la taille du paquet et `ts` l'instant de l'échange (timestamp). Un "flux" d'une source  $s$  vers une destination  $d$  cesse s'il y a un intervalle de 2 minutes entre deux paquets consécutifs de  $s$  vers  $d$  [9]. Dans SQL:1999, un administrateur réseau émettrait la requête suivante afin de compter le nombre de paquets et leur longueur moyenne dans chaque flux.

```
[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
```

---

<sup>1</sup>Les requêtes montrées ici utilisent des fonctionnalités de SQL:1999 relativement avancées. Le lecteur en difficulté est encouragé à se référer à [27].

```

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 src, dest, avg(length), count(ts)
FROM   FlowID
GROUP  BY src, dest, fID

```

Dans les grandes lignes, cette requête doit regrouper les paquets en flux, et au sein de chaque flux compter le nombre de paquets et moyenner leur longueur. Cependant, la recherche des flux est très difficile à exprimer puisque cela implique l'ordre. La première sous-requête, *Prec*, va créer une nouvelle colonne, *ptime*, qui contient les *timestamp* des paquets précédents au sein de chaque source et destination. Ensuite, la sous-requête *Flow* ajoute une colonne de variables booléennes dont l'état est vrai (1) pour chaque paquet dont la différence avec le précédent excède deux minutes et faux (0) sinon. Enfin, la sous-requête *FlowID* ajoute ces drapeaux, créant un flux ID auxiliaire, *fID*. La requête principale utilise ces résultats intermédiaires.

L'optimisation de cette requête vise à réduire la quantité de travail requise par les *PARTITION BY* et *ORDER BY*. Ceci est difficile dans la mesure où les fenêtres définies dans *Prec* et *FlowID* ont des paramètres de glissement légèrement différents. L'optimiseur commercial que nous avons testé a exécuté la requête avec deux tris avant que le regroupement ne soit fait. Ainsi, il n'a pas effectué ces optimisations.

Le problème de ces requêtes est alors double. Leur expression dans SQL:1999 est complexe, et par conséquent à la fois difficile à lire et difficile à optimiser.

Nous avons conçu un modèle de données, un langage et un système où les expressions dépendantes de l'ordre s'écrivent de manière naturelle et mettant en avant des idiomes que l'optimiseur peut utiliser.

### 1.2.1 AQuery: Premier aperçu

Dans notre modèle de données, les tables ne sont pas vues comme des ensembles, mais plutôt comme des entités ordonnées que nous appellerons *arrables* (de l'anglais "array-table", signifiant "matrice-table"). Le tri d'une *arrable* peut être défini dès sa création en utilisant une clause *ORDERED BY* et peut être modifié par la suite.

Notre langage de requêtes, *AQuery*, est une extension sémantique du modèle rationnel d'ensembles (à savoir le SQL 92), qui comprend les clauses classiques *SELECT-FROM-WHERE-GROUP BY-AGGREGATING-HAVING*. Les principales extensions sont basées sur une nouvelle clause appelée *ASSUMING ORDER* qui définit l'ordre des *arrables* identifiées dans la clause *FROM*. Les prédicats et les expressions, quelles que soient leurs clauses d'appartenance, peuvent s'appuyer

sur l'ordre défini par la clause ASSUMING BY, menant à l'expression naturelle de requêtes dépendant de l'ordre telles que nous les montrons dans la section 1.3.

Il revient à l'optimiseur de faire correspondre l'ordre existant des arrables d'entrée avec l'ASSUMING ORDER d'une requête, et de décider si et quand un tri supplémentaire est nécessaire, comme nous le montrons dans la section 1.4. Cette flexibilité provient du fait que dans l'algèbre AQuery, chaque opérateur possède une variante qui préserve l'ordre et une variante qui l'ignore. Les transformations que nous suggérons ici peuvent déplacer le tri sur d'autres opérateurs, impliquant éventuellement un changement vers une autre variante. Nous montrons que ce schéma est capable d'intégrer les nouvelles transformations aux transformations classiques.

Dans la section 1.5, nos expériences montrent des différences de plusieurs ordres de grandeur dans l'interprétation des requêtes entre AQuery et la version commerciale actuelle de SQL:1999, ce qui montre que les transformations produisent des plans très efficaces.

Dans la section 1.6, nous identifions plusieurs langages ayant également considéré l'ordre comme étant de première importance, et desquels AQuery s'inspire. Nous y faisons aussi quelques commentaires sur l'origine de certaines de nos optimisations techniques.

Enfin, la section 1.7 résume nos contributions et décrit notre travail à venir.

## 1.3 Modèle de Données et Algèbre

### 1.3.1 Arrables et Ordre

**Définition 1.1 (Arrables)** – Soit  $\mathcal{T}$  un ensemble de types tel que pour tout  $t \in \mathcal{T}$ ,  $t$  correspond à un type basique (entier, booléen, etc) ou à une matrice unidimensionnelle d'éléments de type basique. Soit  $A$  une matrice infinie d'éléments finis de type  $t \in \mathcal{T}$ . Le cardinal de  $A$  est le nombre de ses éléments. On note  $A[k]$  le  $k$ -ième élément de  $A$ ;  $k$  s'appelle *index* ou *position* dans  $A$ . Les indexes partent de 0. Une *arrable*  $r$  est une collection de matrices  $A_1, \dots, A_n$ , nommés, de même cardinal et telles que tout  $A_i$ ,  $1 \leq i \leq n$ , est d'un type de  $\mathcal{T}$ .  $\square$

La Figure 1.1 montre des exemples deux arrables correctement constituées. Leur schéma correspond au à la table Trades décrite dans le paragraphe 1.2. Notons que si  $A_1, \dots, A_n$  sont tous des vecteurs (c'est-à-dire que leurs éléments sont des scalaires), les arrables ont l'apparence de tables. C'est le cas de l'arrable Trades de cette figure. Nous allons montrer brièvement comment les arrables ayant des éléments vectoriels, tels que l'arrable Series de cette même figure, peuvent être utiles.

**Définition 1.2 (Indexation d'Arrables)** – Le  $k$ -ième registre d'une arrable  $r$  est formé du  $k$ -ième élément de chaque matrice contenue dans  $r$ . Cette opération, que l'on appelle *indexing* (indexage), est représentée par  $r[k] = \langle A_1[k], \dots, A_n[k] \rangle$ .  $\square$

Ticks	ID	price	date	ts
	ACME	12.02	05/11/03	1
	WXYZ	43.23	05/11/03	2
	ACME	12.04	05/11/03	5
	ACME	12.05	05/11/03	9
	WXYZ	43.22	05/11/03	13

(a)

Series	ID	price	date	ts
	ACME	[12.02 12.04 12.05]	05/11/03	[1 5 9]
	WXYZ	[43.23 43.22]	05/11/03	[2 13]

(b)

Figure 1.1: Exemple de deux arrables correctement constituées.

Par exemple,  $\text{Trades}[0]$  correspond au registre  $\langle ACME, 05/11/03, 12.02, 1 \rangle$ ,  $\text{Trades}[1]$  à  $\langle WXYZ, 05/11/03, 43.23, 2 \rangle$ , et ainsi de suite.

Puisqu'une arrable est faite de matrices et que les matrices sont ordonnées, une arrable est ordonnée.

**Définition 1.3 (Ordered by)** – Une arrable  $r$  peut être (lexicographiquement) ordonnée par un sous-ensemble de ses matrices,  $B_1, \dots, B_m \subseteq A_1, \dots, A_n$ . Si l'ordre est croissant et que  $k_1$  et  $k_2$  sont deux indices de  $r$ , tels que  $k_1 < k_2$ , alors soit (i)  $B_1[k_1] = B_1[k_2], \dots, B_m[k_1] = B_m[k_2]$  soit (ii) il existe un  $i$  tel que  $B_i[k_1] < B_i[k_2]$  pour  $1 \leq i \leq m$  et  $B_1[k_1] = B_1[k_2], \dots, B_{i-1}[k_1] = B_{i-1}[k_2]$  pour  $i > 1$ . Les définitions sont symétriques pour les ordres décroissants, mais pour des raisons de clarté nous considérerons l'ordre comme croissant dans tout le document.  $\square$

Par exemple, les arrables  $\text{Trades}$  de la Figure 1.1 pourraient être définies par  $\text{Trades}(\text{ID}, \text{tradeDate}, \text{price}, \text{ts})$  ORDERED BY  $\text{ts}$ .<sup>2</sup>

**Définition 1.4 (équivalence du point de vue de l'ordre)** – Soient  $r$  et  $s$  deux arrables créées à partir du même jeu d'attributs. Supposons que  $r$  est ordonné par (ORDERED BY) des attributs  $X_1, \dots, X_p$  et  $s$  par  $Y_1, \dots, Y_q$ .  $r$  et  $s$  sont alors *order-equivalent* (équivalents du points de vue de l'ordre) par rapport aux attributs  $B_1, \dots, B_m$ , ce qu'on note  $r \equiv_{B_1, \dots, B_m} s$ , si les conditions suivantes sont respectées : (i)  $r$  et  $s$  sont des ensembles équivalents (c'est-à-dire qu'il existe une permutation de lignes  $P^1, P^2$  tel que  $P^1(r) = P^2(s)$ ) ; (ii)  $B_1, \dots, B_m$  est un préfixe de  $X_1, \dots, X_p$  et de  $Y_1, \dots, Y_q$ . Quand  $r$  et  $s$  sont simplement des ensembles équivalents, on les note  $r \equiv_{\{\}} s$ .  $\square$

<sup>2</sup>Pour des raisons de simplicité, nous omettons ici les informations dactylographiques. Une définition complète inclurait aussi NULL et des informations sur l'intégrité référentielle.

### 1.3.2 Sémantiques Orientées-Colonne

Un problème qui existe dans les requêtes dépendant de l'ordre est que chaque ligne du résultat peut dépendre de toute une combinaison de valeurs issues de plus d'une ligne d'entrée. Considérons par exemple la table Trades et une requête visant à trouver la différence entre chaque prix et sa valeur précédente, en supposant un ordre établi. Afin de calculer la différence, la requête doit accéder à deux prix à la fois, qui se trouvent dans des lignes distinctes. Cependant, les *langages orientés-ligne* tels que le SQL ne peuvent opérer que sur une ligne à la fois. C'est pourquoi ils vont devoir avoir recours à une jointure sur eux-même ou à un constructeur auxiliaire pour construire une ligne contenant les deux prix. Cette opération est à répéter pour chaque paire.

A l'opposé, AQuery adopte une *sémantique orientée-colonne* dans laquelle *une variable est liée à toute une matrice*. Puisque dans AQuery les variables font toujours référence à des matrices, les expressions font toujours référence à des mappings d'une liste de matrices vers une matrice. Par exemple, la différence entre les paires telles qu'étudiée ci-dessus peut être saisie par une expression simple `price - prev(price)`. La fonction `prev()` appliquée à une matrice  $A$  est une matrice, telle que  $\text{prev}_A[i] = A[i - 1]$  si  $i > 0$  et  $A[0]$  si  $i = 0$ . Pour deux matrices  $A$  et  $B$  telles que  $|A| = |B|$ , `minus (-)` (moins) est la soustraction appliquée aux éléments.

La fonction `prev()` est un échantillon de l'ensemble des fonctions *vector-to-vector* (vecteur-à-vecteur) qu'inclut AQuery. Ces fonctions sont classées selon leur dépendance vis-à-vis de l'ordre de tri de la matrice d'entrée et du cardinal du résultat qu'elles produisent. Par exemple, `prev()` est *order-dependant* (dépendant de l'ordre) et *size-preserving* (qui préserve la taille). Cette dernière propriété indique que les vecteurs en sortie possèdent autant d'éléments que la matrice d'entrée. Formellement, la dépendance vis-à-vis de l'ordre peut être utilisée comme suit.

**Définition 1.5 (Dépendance vis-à-vis de l'ordre)** – Une application  $e$  allant d'une liste de matrices vers une matrice est dite *order-dependant* si pour chaque matrice opérande  $A_i$ ,  $1 \leq i \leq m$ , où  $m$  est le rang de l'expression, et pour chaque permutation correspondante  $A_i^{perm}$ ,  $e(A_1, \dots, A_m) \equiv_{\Omega} e(A_1^{perm}, \dots, A_m^{perm})$ . Par exemple, `avg(price)` est indépendant de l'ordre. Une expression qui n'est pas indépendante de l'ordre est *order-dependant*. Par exemple, `price - prev(price)` est dépendante de l'ordre.  $\square$

D'autres fonctions dans la catégorie "dépendant de l'ordre" et "préservant la taille" sont les agrégats mouvants. Un minimum courant sur une matrice  $A$ , `mins(A)`, est  $\text{mins}_A[i] = \min(A[i], \text{mins}_A[i - 1])$  pour  $0 < i < |A|$  ou  $A[i]$  pour  $i = 0$ . Les agrégats courants se distinguent par le suffixe "-s". Une somme courante sur une matrice  $A$ , notée `sums(A)`, est  $\text{sums}_A[i] = A[i] + \text{sums}_A[i - 1]$  pour  $0 < i < |A|$ , ou  $A[i]$  pour  $i = 0$ . Quelques agrégats mouvants peuvent être calculés sur des fenêtres glissantes. Par exemple, une moyenne mobile utilisant une fenêtre de taille fixe de  $w$  positions sur une matrice  $A$  est notée `avgs(w, A)` et se définit par  $\text{avgs}_{w,A}[i] = \text{sum}(A[i - (w - 1)]..A[i])/w$ , pour  $w - 1 \leq i < |A|$  ou  $\text{sum}(A[0]..A[i])/i$

pour  $0 \leq i < w - 1$ .<sup>3</sup>

Une autre catégorie de fonctions vecteur-à-vecteur est celle des fonctions dépendant de l'ordre mais ne préservant pas la taille. Si elles conservent soit le début, soit la fin d'une matrice on les appelle *edge functions* (fonctions de bord). Par exemple, les  $n$  premières positions d'une matrice  $A$ , ce qui se note  $\text{first}(n, A)$ , s'écrira  $\text{first}_{A,n} = A[0..n - 1]$ . De même,  $\text{last}_{A,n} = A[|A| - n..|A| - 1]$ .

Les fonctions d'agrégat classiques du SQL ( $\text{min}$ ,  $\text{max}$ ,  $\text{avg}$ ,  $\text{count}$ ) peuvent être vues comme des fonctions vecteur-à-vecteur ne dépendant pas de l'ordre et ne préservant pas la taille.

### 1.3.3 Un Langage de Requêtes et une Algèbre

L'algèbre AQuery gère les opérateurs de l'algèbre relationnelle. Mais ici, chaque opérateur prend comme argument des expressions de type matriciel. Si une expression est dépendante de l'ordre, alors l'opérateur se comporte de manière à préserver l'ordre. Sinon, l'opérateur se comporte d'une manière à l'ignorer. La variante ignorant l'ordre d'un opérateur est tout simplement une variante équivalente en ensemble à sa variante préservant l'ordre. Dans la suite de la section nous définissons les variantes préservant l'ordre des opérateurs de l'algèbre relationnelle.

**Définition 1.6 (Projection)** – Soit  $r$  une arrable et  $e = e_1, \dots, e_m$  une liste d'expressions faisant intervenir les matrices de  $r$ , tel que  $|e_1| = \dots = |e_m|$ . Une projection préservant l'ordre de  $r$  sur  $e$ , notée  $\pi_e^{op}(r)$ , est défini comme suit.

projection(e,r)

1. s:= matrice vide ayant le meme schema que e
2. for i = 0 to |r|-1
3.   append <e<sub>1</sub>[i], ..., e<sub>m</sub>[i]> to s
4. end for
5. output s

Comme mentionné précédemment, si n'importe quel  $e_i$  est dépendant de l'ordre, la projection est dite préservant l'ordre, sinon elle l'ignore, ce que l'on note simplement  $\pi_e(r)$ .  $\square$

**Définition 1.7 (Selection)** – Soient  $r$  une arrable et  $p$  un prédicat faisant un mapping de matrices de  $r$  en une matrice de booléens, tels que  $|r| = |p|$ . Une sélection de  $r$  sur  $p$  qui préserve l'ordre, notée  $\sigma_p^{op}(r)$ , est définie comme suit.

selection(p,r)

1. s:= arrable vide ayant le meme schema que r
2. for i = 0 to |r|-1

---

<sup>3</sup>Une telle définition est monnaie courante dans une application financière. D'autres domaines peuvent nécessiter que la fonction  $\text{avgs}()$  renvoie NULL sur certains positions où la fenêtre est incomplète. Dans tous les cas, il est souvent pratique que la moyenne courante renvoie une matrice de même taille que ses arguments.

3. if p[i] is true
4.     append r[i] to s
5. end if
6. end for
7. output s

Comme pour une projection, une sélection peut être dépendante de l'ordre, et soit préserver l'ordre soit l'ignorer.  $\square$

L'ordre de tri d'une arrable est une propriété qui peut être manipulée indépendamment pour chaque requête.

**Définition 1.8 (Tri)** – Soit  $r(A_1, \dots, A_n)$  une arrable et  $B_1, \dots, B_m \subseteq A_1, \dots, A_n$ . Le tri de  $r$  sur  $B_1, \dots, B_m$  est une permutation  $s$  de  $r$  qui est ORDERED BY (ordonnée par)  $B_1, \dots, B_m$ .  $\square$

Ayant défini ces opérateurs, il est désormais possible de montrer le résultat que donne AQuery sur la requête de meilleur profit.

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

Les clauses (SELECT, FROM, ...) de AQuery sont traitées dans le même ordre que dans le SQL. Sémantiquement, ASSUMING ORDER est interprété comme un tri après que la clause FROM ait été traitée. Cela assure que la requête traite l'ordre comme désiré et oblige les clauses suivantes (WHERE, GROUP BY, HAVING et SELECT) à être traduites en des variations algébriques préservant l'ordre. (C'est la sémantique voulue. l'optimisation peut éviter d'effectuer ce tri si tôt, comme nous allons le montrer plus loin.)

Notons qu'à cause de la sémantique orientée-colonne de AQuery, la fonction mins() est appelée une seule fois et prend pour argument l'ensemble du vecteur price. Soustraire un vecteur (mins(price)) à un autre (price) de même cardinal est une opérations matricielle standard [5], de même que prendre le max() du vecteur résultat.

La requête ci-dessus se traduit en algèbre AQuery de la manière suivante, avec  $e = \max(\text{price} - \text{mins}(\text{price}))$  et  $p = (\text{ID} = \text{'ACME'}) \wedge (\text{tradeDate} = \text{'05/11/03'})$  :

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

Dans AQuery, les regroupements utilisent les capacités d'une arrable à stocker des variables de type matriciel. Intuitivement, une opération de regroupement va partitionner l'arrable opérande en différentes sous-arrables disjointes partageant la même valeur de groupe. Elle va ensuite transformer chaque sous-arrable en une ligne unique en remplaçant (dans la sous-matrice) chaque colonne qui n'est pas groupée par son équivalent matriciel. Par exemple, l'arrable Series de la Figure

1.1 montre l'effet du regroupement des arrables Trades de la même figure selon ID et tradeDate.

**Définition 1.9 (Regroupement)** – Soit  $r$  une arrable et  $g = G_1, \dots, G_m$  une liste d'expressions portant sur les matrices de  $r$  et telles que  $|G_1| = \dots = |G_m| = |r|$ . C'est-à-dire que pour tout  $r[i]$  il doit exister un groupe caractérisé par  $g[i]$ . Le “group-by” (regroupement par) de  $r$  sur  $g$ , préservant l'ordre, noté  $gby_g^{op}$ , est défini comme suit.

```

group-by(g,r)
1. groups := arrable vide ayant le meme schema que g
2. s:= arrable vide ayant le meme schema que 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
8.         concat r[i].C to s[j].C
9.       end if
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

```

L'étape 13 ci-dessus forme une liste d'éléments seuls (équivalent à un vecteur). L'étape 8 concatène à cette liste. Le résultat est que les champs peuvent être des vecteurs. Comme précédemment, le “group by” est dépendant de l'ordre si n'importe laquelle de ses expressions de regroupement l'est. Group-by peut aussi avoir une variante ignorant l'ordre. De plus, il est pratique d'avoir une version de group-by générant l'ordre. Sémantiquement, un tel group-by donne le résultat ordonné par l'expression de regroupement.  $\square$

Les requêtes de gestion des réseaux vont bénéficier des regroupements dépendants de l'ordre de AQuery, comme le décrit la Figure 1.2. Rappelons que cette requête fait appel à un *group-by* sur un hôte source et un hôte destination et un *flow ID* (flux) entre eux. Le flux ID est dépendant de l'ordre – un nouveau flux entre un couple d'hôtes est initialisé dès qu'il y a un intervalle de 120 secondes entre deux paquets consécutifs. Dans AQuery, une telle expression de regroupement correspond à `src, dest, sums(deltas(ts)>120)`. La Figure 1.2(a) montre comment cette expression est calculée, en supposant que l'arrable Trades est trié par rapport à `src, dest` et `ts` et que le booléen TRUE est associé à la valeur 1 et FALSE à 0.

Dans AQuery, le regroupement et l'agrégation sont des opérations indépendantes. L'arrable que nous voyons dans 1.2(b) montre le résultat de l'opération de regroupement seule. Notons que certains champs des colonnes non groupées de Packets sont des matrices. Puisque certains champs peuvent être des matrices

Packets	src	dest	length	ts	deltas(ts)>120	sums(deltas(ts)>120)	
	s1	s2	250	1	F	0	] g1
	s1	s2	270	20	F	0	
	s1	s2	235	141	T	1	□ g2
	s2	s1	330	47	F	1	] g3
	s2	s1	280	150	F	1	
	s2	s1	305	155	F	1	

(a)

Packets'	src	dest	length	ts
	s1	s2	[[250, 270]]	[[1, 20]]
	s1	s2	[[235]]	[[141]]
	s2	s1	[[330, 280, 305]]	[[47, 150, 155]]

(b)

Packets''	src	dest	avg(length)	count(ts)
	s1	s2	260	2
	s1	s2	235	1
	s2	s1	305	3

(c)

Figure 1.2: Regroupement de Trades sur src, dest, sums(deltas(ts) &gt; 120)

(mais pas des arrables), les fonctions d'agrégation peuvent s'appliquer sur une colonne entière ou sur chaque champ. Pour exprimer ces derniers, AQuery fournit un modificateur d'opérateur appelé *each* (chaque) qui applique les fonctions à chaque élément de type matriciel d'une colonne.

**Définition 1.10 (Modificateur *Each*)** – Soit  $A$  une matrice paramètre d'une fonction  $F$ . l'exécution de  $F$ , affectée du modificateur *each* se définit comme suit :

each( $F$ ,  $A$ )

1.  $B :=$  matrice vide de meme type que le resultat de  $F$
2. for  $i = 0$  to  $|A|-1$
3.     append  $F(A[i])$  to  $B$
4. end for
5. output  $B$

Cette définition s'étend naturellement aux cas où  $F$  prend plus d'un seul argument. Un opérateur affecté de *each* préserve nécessairement l'ordre.  $\square$

*Each* est une manière d'appliquer un opérateur à chaque champ d'une colonne groupée. Dans la figure 1.2(c) nous pouvons voir que avg() a été appliqué à

chacune des valeurs matricielles de la colonne length, et de même pour count(ts). Ayant défini les opérateurs impliqués dans les requêtes de gestion de réseaux, nous pouvons maintenant montrer les résultats que donne AQuery.

```
SELECT  src, dest, avg(length), count(ts)
FROM    Packets
        ASSUMING ORDER src, dest, ts
GROUP   BY src, dest, sums(deltas(ts) > 120)
```

La version algébrique de la requête de gestion de réseaux, avec  $e = \text{src, dest, each(avg(),length), each(count(),ts)}$  et  $g = \text{src, dest, sums(deltas(timestamp)>120)}$  ressemble à ce qui suit. Nous notons avec un exposant correspondant les opérations qui ont des composants modifiés par each.

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

Dans AQuery, le produit vectoriel ( $\times$ ) ignore l'ordre et a donc la même définition que dans l'algèbre relationnelle. A l'opposée, les opérateurs ensembliste ont des variantes à la fois ignorant l'ordre et préservant l'ordre. A l'opposé, les jointures possèdent à la fois des variantes préservant l'ordre et l'ingorant.

**Définition 1.11 (Jointure)** – Considérons les arrables  $r(A_1, \dots, A_n)$  et  $s(B_1, \dots, B_m)$ . Une jointure *left-right order-preserving* (préservant l'ordre de gauche à droite) de  $r$  et  $s$  sur le prédicat de jointure  $p$ , noté  $r \bowtie_p^{\text{lr op}} s$ , est noté de la manière suivante.

join(p, r, s)

1. o := arrables vide avec le schema  $\langle A_1, \dots, A_n, B_1, \dots, B_m \rangle$
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  $\langle A_1[i], \dots, A_n[i], B_1[j], \dots, B_m[j] \rangle$  to o
6. end if
7. end for
8. end for
9. output o

L'ordre d'une requête peut nécessiter que seul un des ordres de l'arrable de l'opérateur de jointure soit préservé. Dans ce cas une variation de la jointure, dépendant de l'ordre et plus simple, peut être utilisée. Une jointure *left order-preserving* (préservant l'ordre à gauche),  $r \bowtie_p^{\text{lop}} s$ , est une jointure qui est équivalente en ordre à une jointure préservant l'ordre de gauche à droite de ces mêmes deux arrables, mais seulement par rapport à  $A_1, \dots, A_n$ .  $\square$

Considérons l'arrable Portfolio(ID, tradedSince) ORDERED BY ID qui stocke des informations sur les actions qui composent le portefeuille d'un analyste. C'est un sous-ensemble des actions apparaissant dans Trades. Afin d'extraire les dix dernières cotes pour chaque action du portefeuille, nous utilisons la requête suivante :

```

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

```

Sémantiquement, la requête effectue tout d’abord un produit vectoriel ( $\times$ ) entre Trades et Portfolio. Dans AQuery, le produit vectoriel ignore l’ordre. Ensuite, la clause ASSUMING impose l’ordre de tri désiré et le prédicat de jointure est appliqué. Notons que l’ordre est imposé suite au produit Cartésien. Ensuite, l’arrable produite est partitionnée en groupes selon les valeurs de ID. L’ordre induit par assumed est préservé au sein de chaque groupe. La fonction last() va “rogner” chaque colonne de prix, à valeurs matricielles, pour ne conserver au plus que les dix dernières positions de chaque matrice price. Avec  $e = \text{ID}$ ,  $\text{each}(\text{last}(), 10, \text{price})$  et  $p = \text{Trades.ID} = \text{Portfolio.ID}$ , cette requête peut se représenter comme suit :

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

Notre but ici était de montrer qu’en incorporant l’ordre de bout en bout, AQuery permet d’exprimer des requêtes dépendant de l’ordre de manière naturelle.

## 1.4 Transformations de Requêtes

L’élimination par tri [34] et le déplacement par tri [36] sont des techniques connues qui peuvent être appliquées à l’optimisation d’AQuery. La sémantique d’AQuery et, en particulier, les fonctions de bord intégrées (*first*, *last*, ..) permettent aussi d’autres optimisations agressives liées à l’ordre. Nous introduisons ces nouvelles techniques à travers des exemples.

### 1.4.1 Sélections implicites et Tri de bord (*Sort-edge*)

Soit Connections(host, port, client, timestamp) ORDERED BY timestamp une arrable qui stocke les adresses des clients ayant accédé aux services du réseau (port, host) et quand cela a été fait. Supposons que quelqu’un souhaite trouver le dernier client s’étant connecté au serveur “atlas”. Dans AQuery, cela donne :

```

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

```

Nous allons montrer les plans de la manière schématique usuelle. Le plan initial de la requête ci-dessus est représenté dans la Figure 1.3(a). Nous introduisons quelques notions auxiliaires comme suit. Un arc simple entre deux d’opérateurs signifie que l’opérateur “producteur” donne en sortie des registres qui ignorent l’ordre (c’est-à-dire de la manière la plus efficace ou la plus simple, sans garantie d’ordre). Les arcs doubles signifient qu’il donne des registres en préservant l’ordre.

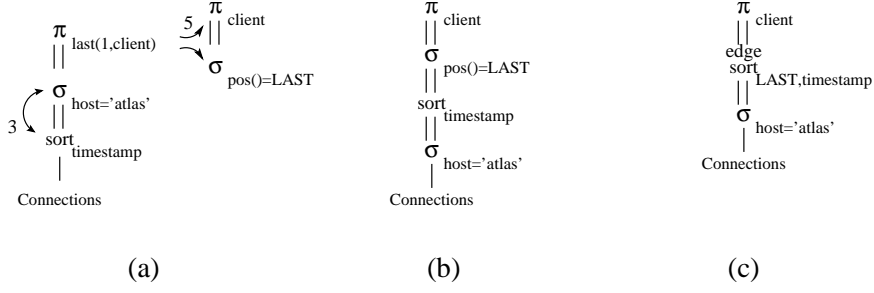


Figure 1.3: Sélection implicite et optimisations de tri de bord

Les flèches représentent l'effet de réseau de l'application d'une transformation. Chaque flèche est affectée du numéro de la transformation correspondante. Les descriptions formelles des transformations sont données dans la Table 1.1. On dit  $\text{pos}(r) = i$  lorsque l'on se réfère au registre  $r[i]$ . Les indices spéciaux pour une matrice  $r$ , FIRST (premier) et LAST (dernier), sont respectivement 0 et  $|r| - 1$ . Enfin,  $\text{order}(r)$  renvoie la liste des attributs par lesquels l'arrable  $r$  est ORDERED BY.

Une sélection régulière telle que  $\sigma_{\text{host}='atlas'}$  peut être effectuée sur un tri [36]. La transformation 2 de Table 1.1 est une légère variation de celle de [36] où la conservation de l'ordre est rendue explicite. Il y a ici deux avantages à permuter le tri et la sélection : la sélection peut bénéficier de l'ordre existant sur Connections (host), et retarder le tri réduit le nombre de registres qui auraient nécessité un tri.

Réduction/Élimination de Tri	
(1) $\text{sort}_A(r) \equiv_{\text{order}(r)} r$	si $A$ est un préfixe de $\text{order}(r)$
Sélection	
(2) $\sigma_p^{op}(\text{sort}_A(r)) \equiv_A \text{sort}_A(\sigma_p(r))$	si $p$ n'est pas dépendant de l'ordre
Projection	
(3) $\pi_{e[i]}^{op}(r) \equiv_{\text{order}(r)} \pi_e^{op}(\sigma_{\text{pos}()=i}(r))$	$e$ est une expression portant sur les matrices de $r$
Jointure and Semi-Jointure	
(4) $\text{sort}_A(r \bowtie_{A=B} s) \equiv_A \text{sort}_A(r) \bowtie_{A=B}^{lop} s$	si $A, B \in$ schéma de $r, s$ , resp.
(5) $\text{sort}_A(r \ltimes_{A=B} s) \equiv_A \text{sort}_A(r) \ltimes_{A=B}^{lop} s$	si $A, B \in$ schéma de $r, s$ , resp.
(6) $\sigma_{A=(B[i])}^{op}(r) \equiv_{\text{order}(r)} r \ltimes_{A=B}^{lop} \sigma_{\text{pos}()=i}(r)$	si $A, B \in$ schéma de $r$
(7) $\sigma_p^{op}(r \bowtie_{A=B}^{lop} s) \equiv_{\text{order}(r)} \sigma_p^{op}(\sigma_p^{each}(\text{gby}_A(r)) \bowtie_{A=B}^{lop} s)$	si $A, B \in$ schéma de $r, s$ , resp., $p$ est 'pos()=FIRST' ou 'pos()=LAST', et $B$ est unique
Group-By	
(8) $\text{gby}_A^{op}(\text{sort}_{A,B}(r)) \equiv_{A,B} \text{sort}_B^{each}(\text{gby}_A^{og}(r))$	

Table 1.1: Un sous-ensemble des équivalences entre le tri et les opérateurs restants de l'algèbre

La projection  $\pi_{\text{last}(1,\text{client})}$  inclut une sélection implicite, c'est-à-dire qu'elle s'intéresse à un seul client. C'est là une particularité de la sémantique ori-

entée colonne de AQuery : une projection sur une fonction qui fait elle même une sélection. La transformation 3 de Table 1.1 est une nouvelle transformation qui remplace une projection indicée par une un projection pure et une sélection des positions désirées. Si de telles positions sont à une extrémité de la matrice opérande, nous appelons cette sélection *edge selection* (sélection de bord). Les résultats obtenus suite à l'application de cette transformation peuvent être vus sur la Figure 1.3(b).

L'avantage qu'il y a à isoler les sélections de bord des projections originales est que tandis que cette dernière ne peut être déplacée aisément, la première peut. Dans cet exemple, l'existence d'une sélection de bord après un tri suggère qu'il n'est pas nécessaire de trier tout le résultat juste pour pouvoir utiliser quelques-uns des éléments. Dans AQuery, l'opérateur physique *sort-edge* (tri de bord) implémente le modèle logique  $\sigma_{edge-condition}^{op}(sort(r))$ . Le tri de bord utilise un tri par tas modifié pour conserver les  $n$  premiers (ou derniers) éléments, comme approprié. Ceci est similaire à l'approche utilisée en [8], sauf que nous modifions le tri par tas pour le rendre stable.<sup>4</sup> Le plan final est montré dans la Figure 1.3(c).

### 1.4.2 Partage de tri (*Sort Splitting*)

Il y a des situations dans lesquelles l'ordre qui existe dans une arrable facilite l'évaluation d'une partie d'une requête, même s'il ne correspond pas à l'ASSUMING ORDER de la requête. Considérons de nouveau l'arrable Connections ORDERED BY hosts. La requête suivante trouve tous les clients qui se sont connectés au dernier hôte auquel un client s'est connecté.

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

Notons que le prédicat `host = last(1,host)` a un sens en tant qu'expression matricielle. La matrice `host` est comparée à la matrice à un seul élément (traitée comme un scalaire) `last(1, host)`, ce qui donne une matrice de booléens. Les positions renvoyant le booléen faux sont éliminées par la clause WHERE (où).

Un plan initial pour cette requête apparaît dans la Figure 1.4(a). Timestamp n'est pas un préfixe de  $order(Connections)$ , et par conséquent le tri sur timestamp est nécessaire. Cependant, `host` est un préfixe de  $order(Connections)$ , C'est pourquoi la sélection  $\sigma_{host=last(1,host)}$  peut en profiter. C'est là qu'intervient la technique de partage de tri.

Si  $A$  et  $B$  sont des matrices d'une arrable  $r$ , une sélection  $\sigma_{A=(B[i])}(r)$  peut être remplacée par une demi-jointure comme le décrit la transformation 6 de la table 1.1. Le but de la demi-jointure est que nous pouvons désormais manipuler l'ordre différemment sur chacun des arguments de la celle-ci.

---

<sup>4</sup>Un tri stable est un tri qui ne change pas l'ordre original des registres possédant des valeurs identiques sur la ligne triée. Le tri par tas n'est pas naturellement stable. Il le devient si l'on concatène un tuple ID à la ligne.

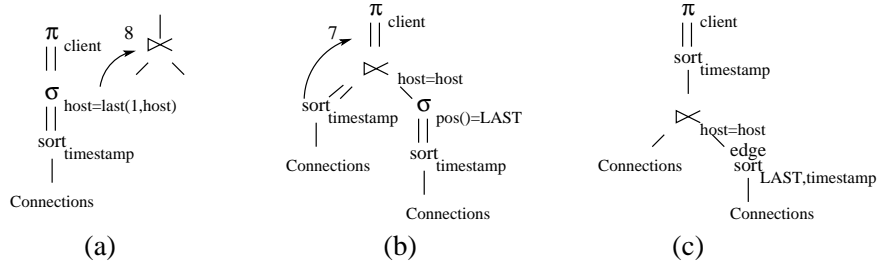


Figure 1.4: Optimisation de partage de tri

La Figure 1.4(b) montre le résultat de l'application de la transformation. Notons que  $last(1, host) = host[LAST]$ . Analysons chaque côté de la demi-jointure à tour de rôle. Sur le côté droit, nous avons le modèle sélection de bord / tri, qui peut être implémenté de manière efficace, comme nous avons pu le voir. À l'opposé, le tri du côté gauche change ce qui aurait pu être un ordre intéressant pour l'opération de semi-jointure. Par conséquent, nous pouvons la différer jusqu'à après la jointure. La transformation 5 de la Table 1.1 permute une semi-jointure et un tri. Cette transformation peut être dérivée de [36], et montre que trier le résultat d'une semi-jointure est équivalent à trier sa partie gauche et ensuite d'effectuer une semi-jointure préservant l'ordre, en supposant que les conditions énumérées dans Table 1.1 sont valables. L'effet de réseau est ici que le calcul du prédicat de la semi-jointure est facilité par un ordre existant et que le tri sur le timestamp n'a besoin d'être fait que pour les registrés générés par la semi-jointure — beaucoup moins coûteux que la jointure de départ effectuée sur toute l'arrable. Le plan ainsi produit apparaît dans la Figure 1.4(c).

Une technique opposée où un gros tri est remplacé par plusieurs plus petits est décrite dans ce qui suit.

### 1.4.3 Enchâssement de tris (*Sort Embedding*)

Considérons l'arrable `Trades(ID, tradeDate, price, timestamp)`, cette fois-ci sans `ORDERED BY` prédéterminé. (Souvent, les "trades" arrivent dans un ordre proche de celui des timestamp.) La requête recherche les dix prix les plus récents pour chaque titre ID. Dans AQuery :

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

Un plan initial pour cette requête est montré dans la Figure 1.5(a). Nous pouvons séparer la sélection implicite de la projection comme nous l'avons fait auparavant. Le plan ainsi obtenu apparaît dans la Figure 1.5(b).

Il est possible de retarder le tri jusqu'à la fin du `GROUP BY ID`. S'il était retardé, le tri devrait seulement être appliqué dans chaque groupe. C'est ce que nous appelons l'enchâssement de tris. De plus, pour cette requête particulière les tris plus petits seraient alors suivis de sélections de bord ; le tri de bord serait

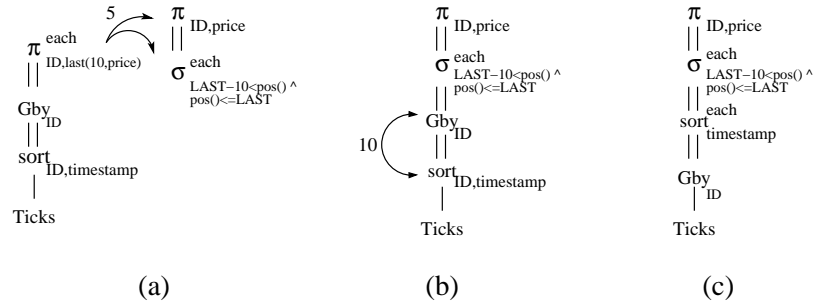


Figure 1.5: Optimisation d'enchâssement de tris

approprié. La transformation 8 de la table 1.1 permet la permutation d'un tri et d'un group by. Notons que (a) les résultats du group by doivent être triés sur la liste de regroupement (c'est-à-dire un opérateur générateur d'ordre), et (b) le regroupement doit être fait sur un préfixe des arguments du tri. Le résultat de cette transformation est montré dans la Figure 1.5(c). Remarquons qu'un arc double relie group by et sort-each, puisque cette instance de group by est génératrice d'ordre.

#### 1.4.4 *Edgeby* (“ Par bords ”) et Sélection de Bord Précoce (*Early Edge Sélection*)

Intéressons nous à un autre scénario dans lequel une sélection de bord peut réduire le cardinal tôt dans une requête. Nous utilisons les arrables Trades ici également, mais supposons qu'elles sont désormais ORDERED BY timestamp. L'arrable Portfolio(ID, name, tradedSince) ORDERED BY ID ORDERED BY ID stocke le sous-ensemble des titres avec lequel un analyste travaille. Name est un identifiant unique de titres dans Portfolio, ainsi que ID. Pour trouver le dernier prix d'un titre appelé “DataOrder”, il faudrait faire :

```

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

```

Un plan initial pour cette requête est décrit dans la Figure 1.6(a). Une heuristique commune pour améliorer les performances est de déplacer la sélection régulière sur le tri, et ensuite sur la jointure. La transformation 2 de la Table 1.1, qui vient de [36], nous permet de le faire. Le résultat est présenté dans la Figure 1.6(b).

A la détection d'un ordre existant approprié (c'est-à-dire order(Trades) satisfait l'ASSUMING ORDER de la requête), l'optimiseur essaierait d'éliminer complètement le tri. La transformation 4 de la table 1.1 permute une jointure avec un tri, tout en gardant une trace de l'ordre. Il s'agit d'une légère variation d'une transformation de [36] dans laquelle la préservation de l'ordre est rendue explicite. Puisque Trades est déjà ORDERED BY timestamp, ce tri peut être

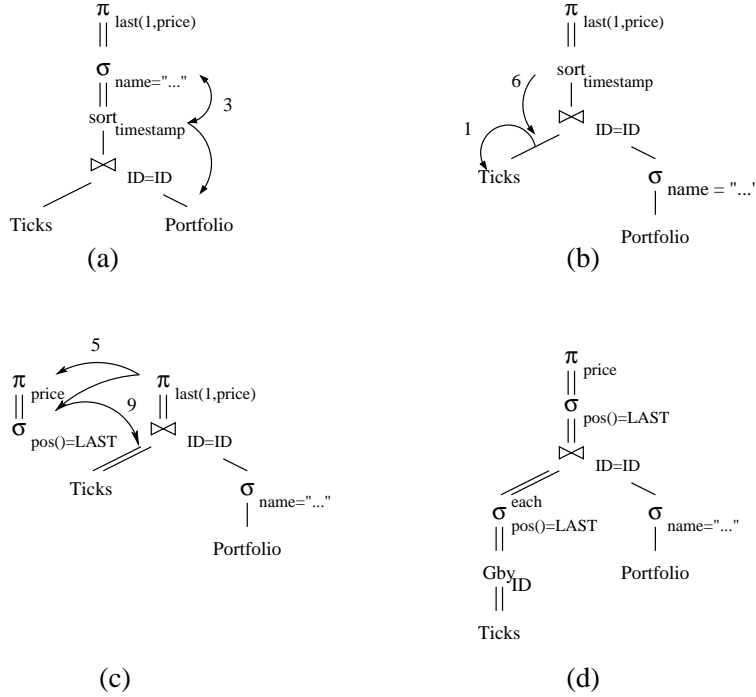


Figure 1.6: Optimisation précoce “ edgeby ”

éliminé, comme le détermine la transformation 1 de la Table 1.1. Cette transformation vient de [34]. Le résultat apparaît dans la Figure 1.6(c).

Cette requête contient également une “ projection+sélection ”, que nous pouvons encore une fois séparer en deux. La présence de la sélection de bord due à la jointure suggère que nous n’aurons peut-être pas besoin d’effectuer la sélection de bord dans son intégralité. Portfolio.ID est une clé et donc garantit que chaque registre de Trades correspondra à au plus un registre de Portfolio. (Les jointures avec une clé extérieure sont parmi les équi-jointures les plus fréquentes). Dans ces conditions, nous déplacerions cette sélection de bord de la manière suivante. Pour chaque ID de Trades, trouvons son dernier registre. Ceci peut être fait en regroupant Trades par ID et en “ sélectionnant par bord ” chaque dernier registre. A cause de la réduction de cardinal induite par le edgeby, la jointure serait effectuée sur beaucoup moins de registres. La sélection finale choisirait alors le prix souhaité. C’est ce que fait la transformation 8 de la table 1.1. Le plan final est montré dans la Figure 1.6(d).

Une sélection de bord appliquée à des groupes est un idiome, appelé *Edgeby*, qui peut être fortement optimisé. Edgeby est un opérateur physique capable d’implémenter le schéma logique  $\sigma_{edge-condition}^{each}(Gby(r))$ . Au lieu de séparer tous les éléments d’une arrable en groupes et de n’en utiliser qu’un slab (par exemple first  $n$ , last  $n$ , drop  $n$ , etc), edgeby se débarrasse, à la volée, des éléments de groupes qui remplissent déjà les conditions de bord. Dans notre exemple, nous avons besoin de la fin d’une matrice (c’est-à-dire last() d’un ordre croissant) de prix pour chaque ID. Un edgeby ID à l’envers sur les Trades triés garde un registre s’il appartient à une ID inconnue jusqu’à présent, ou à un groupe possédant moins de dix registres.

## 1.5 Résultats Expérimentaux

Afin d'évaluer les performances relatives de (a) les requêtes de AQuery contre celles de SQL:1999 et (b) les traductions induites par la syntaxe contre les plans optimisés pour les requêtes d'AQuery, nous avons mené une étude expérimentale.

Toutes nos expériences ont été exécutées sur un Pentium III-M 1,13 MHz équipé de 1 Go de mémoire, tournant sous Linux, sans réglages particuliers ou priorités de processus. Les minutages consignés ici correspondent à des mesures à la seconde près.

### 1.5.1 Mesures de Performance

Les rendus d'AQuery du meilleur profit et des requêtes de gestion de réseaux ne contenaient pas d'enchassements et étaient simples à optimiser. La Figure 1.7 montre les améliorations de performances des plans de AQuery comparées à celle d'un optimiseur commercial de SQL:1999. Les résultats d'AQuery ont été entre huit et vingt-et-une fois plus rapides pour la requête du meilleur profit, et entre deux et trois fois plus rapides pour la requête de gestion de réseaux.

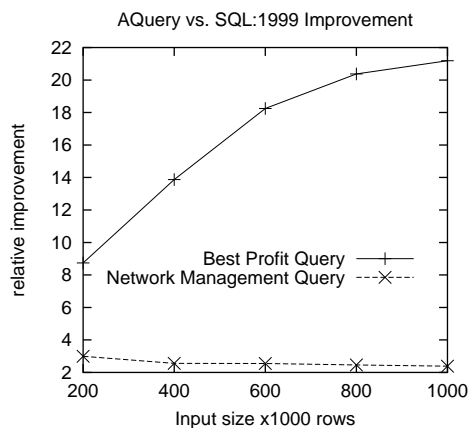


Figure 1.7: Les performances relatives de AQuery contre SQL:1999

Les chiffres de la requête du meilleur profit ont été générés en utilisant des arrables/tableaux Trades avec un nombre variable de titres allant de 200 à 1000, et en utilisant 1000/trades par titre. Rappelons que cette requête s'intéressait à un profit pour un titre donné et une date donnée. La différence de temps est due à la capacité d'AQuery à déplacer le prédicat de sélection et à utiliser un index pour l'évaluer. Le tri suivant ne réorganise que les trades pour le titre concerné. Puisque la représentation SQL:1999 a utilisé une structure d'enchassements compliquée (voir Section 1.2), son optimiseur n'a pas pu déplacer la sélection. Ce plan a trié des tuples qui en fin de compte ont été jetés.

Pour la requête de gestion de réseau nous avons utilisé une arrable/table Packets avec 100 sessions et un nombre variable de paquets pour chaque session, allant de 2K à 10K. Le plan d'AQuery était plus rapide, puisqu'il ne nécessitait qu'un tri pour être réalisé, celui appliquant la cause ASSUMING ORDER. Le Group-by de

AQuery dépend — et donc bénéficie — de cet ordre. A l’opposé, l’optimiseur de SQL:1999 a du trouver comment comment s’occuper de deux spécifications non liées WINDOWS (voir Section 1.2). Cela a amené à avoir deux tris distincts avant le traitement du group-by, ce qui ne leur a pas été bénéfique.

*Morale : La simplicité structurelle d’AQuery aide à trouver de meilleurs plans.*

Dans la plupart des cas, un edgeby requiert une petite fraction du temps nécessaire à l’exécution du group-by associé, s’il est fait dans son intégralité comme le montre la Figure 1.8(a). Nous avons utilisé l’arrable Trades avec 1 million de registres répartis équitablement parmi 10, 100, 1000 et 10000 titres. Un edgeby sur le titre ID, de taille de slab variable, est testé. Plus edgeby peut mettre des titres de côté, plus rapide est son temps de réponse. Par exemple, lorsque seuls quelques titres distincts sont utilisés, les groupes sont larges, et donc la plupart des registres ne tiennent pas dans le slab, même pour les tailles de slab les plus grandes que nous ayons testées, ce qui améliore beaucoup les performances. Alors que les groupes deviennent plus petits (c’est-à-dire que plus de titres différents sont utilisés), les slab hautement sélectifs obtiennent de meilleures performances. Un cas dégénéré est vu lorsqu’un 100-slab est pris dans des groupe ayant eux même 100 registres de largeur (soit 10000 titres). Ici edgeby n’améliore pas les performances — mais n’y nuit pas non plus.

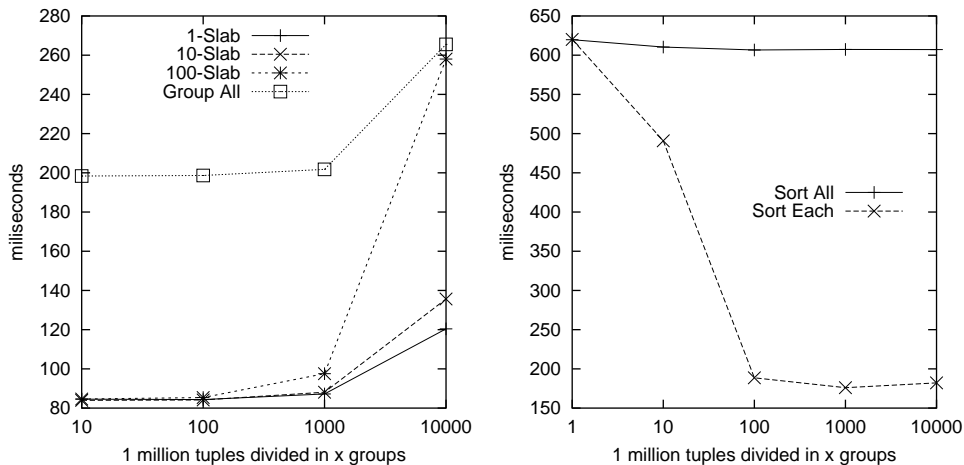


Figure 1.8: Efficacité des techniques de réduction de travail.

L’idée qui repose derrière la technique d’enchâssement de tris est qu’un tri peut être retardé jusqu’à après un group-by, et peut être remplacé par plusieurs tris effectués sur les colonnes groupées (des *sort-each*). La Figure 1.8(b) caractérise les gains de performance de chaque sort-each en comparaison du tri de l’ensemble qu’ils remplacent. Nous avons utilisé des arrables de 1 million de registres et avons fait varier le nombre de groupes dans lesquels l’arrable était divisée. Quand il n’existe qu’un groupe, il n’y a aucun intérêt à appliquer la technique — mais, encore une fois, il n’y a aucun inconvénient à le faire. Échanger un gros tri pour plusieurs petits n’est rentable qu’à partir du moment où il existe plus de 10 groupes. Le *sort-edge* a présenté des résultats similaires à ceux de *sort-stop* [8] et nous omettons donc ces résultats.

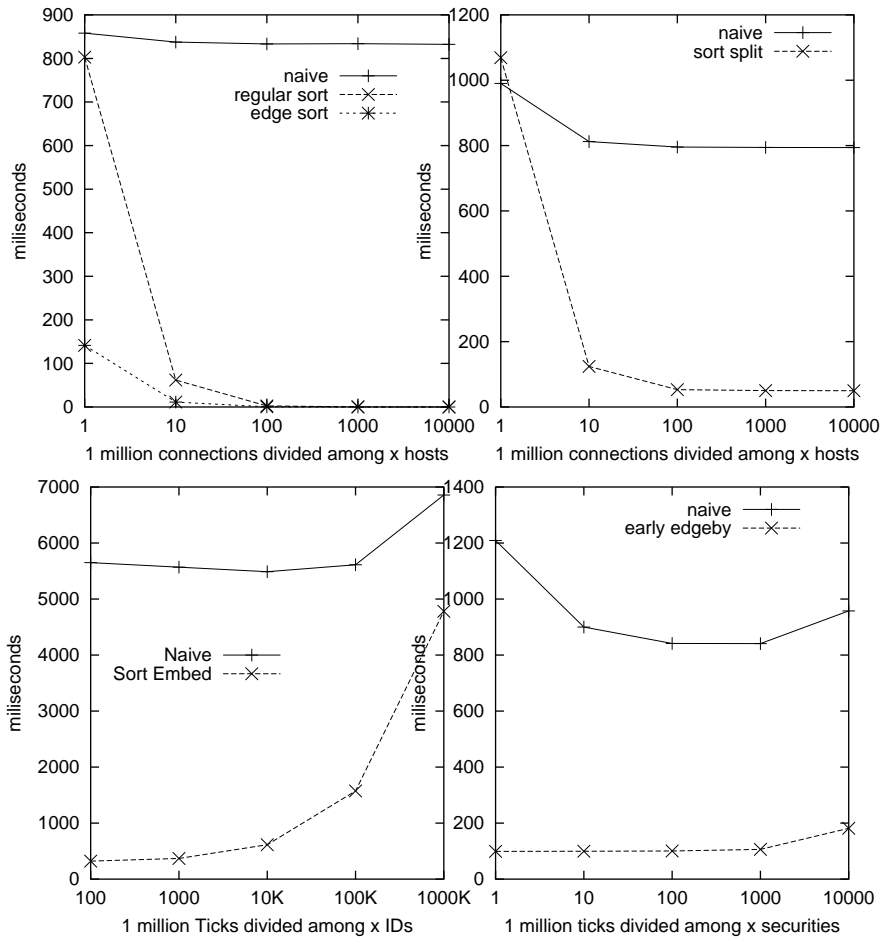


Figure 1.9: Plan optimisé contre plan non-optimisé.

*Morale : La réduction de travail due aux sélections de bord ou aux techniques de maniement de tri est conséquente.*

L'optimiseur de AQuery intègre le tri avec les opérateurs relationnels standard, tels que " sélection " et group-by. Par exemple, dans la requête de la section 1.4.1 il était capable d'appliquer un *push-down* de la sélection sur un tri. La Figure 1.9(a) montre que cette technique était particulièrement efficace pour des cas d'arrables où le nombre d'hôtes distincts était supérieur à 10. Mais les plans utilisant un tri simple sur des cas de moins de 10 hôtes différents s'effectueraient toujours mal. En identifiant la sélection de bord implicite de cette requête et en l'utilisant pour réduire le nombre de registres à trier, AQuery a généré un plan optimisé qui s'est avéré le meilleur dans tous les cas testés.

La Figure 1.9(b) montre le gain de performance qu'il y a à appliquer la technique de partage de tri à la requête donnée en exemple dans la section 1.4.2. L'efficacité du plan optimisé provient de ce qu'il retarde l'exécution de l'ordonnancement ASSUMING jusqu'à après que la semi-jointure ait réduit le nombre de registres à trier. Les gains se sont stabilisés dans les cas de 100 ou plus hôtes différents, parce qu'à ce point le coût de la requête est dominé par le la semi-jointure elle-même, au lieu du tri de ses résultats. Notons que l'application

de cette technique dans le cas où le nombre de sommets est trop faible (par exemple un seul) peut nécessiter un préliminaire inutile — quoique petit. Ainsi, cette technique dépend de la répartition des données, soulignant le besoin des optimisations basées sur le coût.

La Figure 1.9(c) montre les comparaisons de performances des plans de la requête donnée en exemple dans la section 1.4.3. Le plan naïf trie la totalité de l'arrable, regroupe tout le résultat et n'applique la sélection de bord qu'à la fin. Le coût reste relativement élevé, même lorsque la sélection de bord retire la plupart des registres. A l'opposé, les plans d'optimisation échange un grand tri contre plusieurs plus petits – des sort-edge, en fait. Ainsi, même dans le cas dégénéré où chaque groupe n'a qu'un registre (c'est-à-dire que le nombre d'hôtes distincts est égal au cardinal de l'arrable), le plan optimisé évite le coût d'un grand tri. Les courbes montrent des différences de plusieurs ordres de grandeur sur des cas avec un petit nombre d'hôtes distincts.

Enfin, la Figure 1.9(d) montre les résultats des plans naïfs et optimisés pour la requête donnée en exemple dans la Section 4.2.3. En appliquant un edgeby tôt dans le plan, le nombre de registres ayant besoin de la jointure est considérablement réduit. Le plan d'optimisation prend aussi avantage de l'ordre existant, éliminant complètement tout tri. Le résultat est des temps de réponse sans cesse plus brefs.

*Morale : Les transformations d'AQuery apportent des améliorations de performances conséquentes, surtout lorsque utilisées avec des optimisations de requêtes basées sur le coût.*

## 1.6 Travaux apparentés

Alors que le SQL:1999 est l'implémentation de requêtes dépendant de l'ordre la plus conséquente d'un point de vue commercial d'autres systèmes, à la fois commerciaux et expérimentaux, ont proposé beaucoup d'idées excellentes.

### 1.6.1 Techniques d'Optimisation

Dans le processus d'optimisation, l'ordre de tri a toujours été traité comme une propriété physique qui devait être incluse dans un plan (quand elle n'était pas spécifié par la clause `ORDER BY` de SQL:1999) afin d'aider un algorithme efficace tels que la jointure-fusion. Des mécanismes tels que la *glue* (colle) [23] de Starburst ou le *enforcer* [13] de Volcano s'assurèrent qu'une étape de tri était ajoutée dès qu'un algorithme efficace le requérait.

Dans [34], les auteurs ont ajouté une étape d'optimisation de l'ordre avant que les plans ne furent énumérés dans le contexte du processus d'optimisation de DB2. Cette étape peut améliorer radicalement des requêtes ayant des critères d'ordre dus aux clauses `ORDER BY`, `GROUP BY` ou au modificateur `DISTINCT`. Mais l'optimisation de l'ordre [34] était toujours une étape d'optimisation séparée dans la mesure où les transformations impliquant des réductions ou des éliminations de tris (c'est-à-dire des tris s'appliquant à moins d'éléments) n'étaient pas con-

sidérées en même temps que celles impliquant d'autres opérations algébriques. A l'opposé, parce que nous considérons l'opérateur de tri avec les autres opérateurs dans l'esprit de [36], nous sommes capables de découvrir des techniques telles que le partage de tri ou l'enchâssement de tris (c'est-à-dire respectivement les transformations 6 et 8 de la table 1.1).

Dans [8] une clause `STOP AFTER` a été suggérée, elle était capable de limiter le cardinal des résultats des requêtes dont les résultats étaient ordonnés (par un `ORDER BY`). Lorsque seulement les *k-meilleurs* tuples du résultat d'une requête ont besoin d'être consommés, les requêtes peuvent tourner beaucoup plus vite. Notre sort-edge est proche du stop-sort. Une différence entre ce travail et le nôtre est que, encore une fois, le processus d'optimisation de l'ordre apparaît dans un endroit à part. Parce que nous intégrons les deux, AQuery utilise également le concept de stop-after au sein des groupes. Après tout la restriction de cardinaux est une technique généralement utile.

Au meilleur de notre connaissance, le premier travail pour fournir une optimisation de l'ordre dans une infrastructure intégrée était [36]. La plupart de leurs transformations basées sur les listes s'appliquent à AQuery. Les transformations que nous avons présentés ici étendent le cadre de leur travail avec plusieurs nouvelles techniques (telles que les transformations 3 et 7 de la Table 1.1).

D'autres règles d'optimisation liées à l'ordre ou aux structures ordonnées furent suggérées dans le contexte de leur langage de requêtes correspondant, ce que nous discutons dans la suite.

## 1.6.2 Langages

AQuery est un descendant du système KSQL de KX [20], duquel AQuery tire sa notion d'arrable — une implémentation de tables partitionnés totalement verticalement, dont chaque colonne est une matrice. AQuery diffère de KSQL en essayant de conserver l'esprit de SQL bien plus que KSQL, ceci par l'utilisation de la clause `ASSUMING` pour rendre l'usage de l'ordre déclaratif, par l'introduction de transformations, et par l'exploitation d'une structure d'applications basées sur le coût.

Les langages de requêtes séquentielles SEQUIN [32] et SRQL [30] sont des précurseurs de SQL :1999 dans le sens où ils sont des langages de SQL qui gèrent des requêtes basées sur l'ordre. Ils restent fidèles à l'esprit de [24] qui a montré que l'ordre était un paramètre crucial dans les requêtes. SEQUIN traite les séquences comme un type de donnée abstraite étendu, bien qu'une séquence puisse servir comme unique source de données pour une requête. SRQL a été inspiré par SEQUIN et traite les tables comme des relations ordonnées. AQuery a emprunté à ces langages l'introduction précoce dans une requête d'une clause définissant l'ordre. A la fois SEQUIN et SQL gardent la sémantique de tuple du SQL, à l'opposé de l'exécution par vecteurs (colonnes) de AQuery. En conséquence, plusieurs expressions vectorielles valides dans AQuery sont invalides dans ces langages, comme par exemple `max(price - mins(price))`. Une autre différence est dans la manière dont les relations/arrables non-1NF sont traitées. Dans [32], si une table a un attribut séquentiel, SEQUIN a pour habitude d'exprimer des prédicats sur la séquence,

tandis que SQL les exprime sur la table. Cela peut conduire à des requêtes difficiles à lire et à de complexes transformations d'optimisation. De plus, SRQL n'a pas poursuivi les relations non ordonnées 1NF, tandis que nous les avons trouvées très utiles.

Des efforts non-SQL pour les langages de requêtes et les modèles de données basées sur les matrices ont été suggérées auparavant [22, 26]. Ni AQL [22] ni AML [26] ne fournissent de mécanisme déclaratif pour définir l'ordre dans lequel les requêtes manipulent les données. Les requêtes traitent les données dans l'ordre dans lequel elles sont stockées. Tandis que cela a un sens pour des bases de données d'images raster [26] ou des données scientifiques de format CDET [22], ça en a beaucoup moins dans le calcul général de données. Néanmoins, on pourrait soutenir que ces langages pourraient intégrer facilement une fonction de tri et exprimer la plupart, si ce n'est tout, des requêtes décrites ici. Nous sommes d'accord. De telles extensions seraient les bienvenues, parce que les implémenteurs obtiendraient peut-être des optimisations intéressantes qui pourraient être complémentaires aux nôtres. Enfin, nous avons pour des raisons pragmatiques une préférence pour un dialecte SQL, mais certains peuvent ne pas être d'accord sur ce point.

Une particularité particulièrement inspirante des optimiseurs de AQL est qu'ils ont de puissantes capacités à optimiser les opérateurs (ou des fonctions récemment ajoutées) au niveau du calcul, c'est-à-dire par l'application de variations de  $\lambda$  réductions de calcul sur les définitions des opérateurs. Les réductions aident à trouver des formes d'une expression syntaxiquement plus simples, tout en gardant intacte sa sémantique. Nous n'avons pas encore complètement exploité cette capacité dans AQuery. D'autre part nous avons montré que, par exemple, la technique de partage de tri requiert plus que simplifier une expression. Elle implique de transformer ce qui était un tri plus une sélection en une semi-jointure plus deux tris plus une sélection – et cela nous a amené à trier moins de tuples qu'avec l'expression simplifiée. A l'opposé, AML utilise un jeu de transformations déterminé, destiné génériquement à l'application de fonctions à des slabs de matrices. Une fusion complète de ces idées requiert plus d'explications.

## 1.7 Conclusion

AQuery s'appuie sur des travaux de langage de d'optimisation de requêtes déjà effectués pour remplir les objectifs suivants :

1. Incorporer l'ordre de manière déclarative à un langage de requêtes (en utilisant la clause ASSUMING) basé sur le SQL 92.
2. Introduire une sémantique de requêtes qui, tout en restant compatible avec les SQL, gère des calculs inter-tuples sans enchâssement de requêtes.
3. Ajouter des fonctions dépendantes de l'ordre (par exemple *sums*, *first*) qui sont naturelles à exprimer et flexibles, par exemple, pour permettre la requête de  $k$ -meilleurs éléments au sein d'un groupe.

4. Créer une infrastructure d'optimisation simple, quoique puissante, qui amène des gains en vitesse d'au moins un ordre de grandeur par rapport aux systèmes commerciaux de SQL:1999 pour les requêtes naturelles. Les optimisations de bord, le partage et l'enchâssement de tris paraissent particulièrement prometteurs pour les requêtes dépendant de l'ordre.



**Part II**

**Mémoire en anglais**

# Acknowledgements

The preparation of this thesis spawned three countries and I am indebted to quite some people in all of them.

Prof. Sergio Carvalho and I met while I was at PUC-Rio, in Brazil. He was a source of vitality – despite being ill – inspiration, and motivation. This image of him in his office, door opened, welcoming me with a “What are you doing here? Why are you here?” and a smile, marked the beginning of my thesis work. At first, I was caught by surprise by the question, and simply answered something along the lines of “I thought we had a meeting.” But then he insisted: “Who, in his or her right state of mind pursues a Doctorate degree?” and went on about adversities of being a graduate student. We followed this ritual several times. We always found good reasons. “Are you sure? Ok then. Let’s work.”

Prof. Dennis Shasha and I started to work together – one could say by coincidence – before we met personally. It all started with this e-mail, out of the blue, inquiring about my progress. I was in France and he in the US. Why was this person contacting me? It took me a while to recognize Dennis from the articles I was reading and linking the e-mail writer to the professor that was in sabbatical at INRIA just before I arrived. Dennis showed me by example what is to do academic research. What a fortunate coincidence, having picked his topics to work at INRIA; I would not have found a better source of talent, inspiration, and generosity.

Thanks to Dr Eric Simon and Dr Talel Abdessalem I had great opportunities of pursuing my work in excellent environments. Eric welcomed me warmly to his group at INRIA and gave me the environment I was looking for to develop my graduate work. Talel smoothed a very tricky transition between universities when I was already in mid-flight and made it so that the disruption was practically imperceptible. That was not a small feat.

How lucky I feel of having crossed these people’s paths.

I’m also thankful to all the following people for their support and for the opportunities we had to interact during these years: Adriana Alvim, Elisabeth Baque, Pascal Baque, Luc Bouganim, Peter Buneman, Harriet Barry, Celso Carneiro, Tatiana Dutra, Georges Gardarin, Olivier Gardarin, Helena Galhardas, Rosalba Giugno, Ted Johnson, Sergio Lifschitz, Fernanda Lima, Ioana Manolescu, Lee Phillips, Joao Pereira, Marcus Poggi, Fabio Porto, Philippe Poucheral, Lourdes Santana, Jean-Marc Saglio, David Tanzer, Patrick Valduriez, Khaled Yagoub, and Yunyue Zhu.

No dream seem too far fetched and no change seem too unsettling when I have the warmth and support of Gimi, Fany, Branca, Iona, and the kids. Without them nothing would make much sense anyway.



# Abstract

An order-dependent query is one whose result (interpreted as a multi-set) 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 the five price moving average in a trade 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 example, uses a data ordering mechanism called a “window.” However, order-dependent queries become difficult to write in those languages, and optimization techniques for these features, applied as pre- or post-enumerating phases, are generally crude.

In this thesis we show that when order is a property of the underlying data model and algebra, writing order-dependent queries in a language can be natural as is their optimization. We introduce AQuery, an SQL-like query language and algebra that explicitly support for ordering of records. We also introduce a framework for optimization of order-dependent queries. The framework is able to take advantage of the large body of existing query transformations while incorporating new ones described here. We show by experiment that the resulting system is orders of magnitude faster than current SQL:1999 systems on many natural order-dependent queries.



# Chapter 2

## Introduction

### 2.1 Order-Dependent Queries

An order-independent query is one for which the results (interpreted as a multiset) do not change if the order of the input records change. In a stock-quotes database, for instance, calculating the maximum price of a stock in a given day is order-independent. Regardless of the order in which records are examined, their maximum is the same. Relational databases support order-independent queries extremely well.

By contrast, finding the price changes of a stock over many days depends on order. Such a query is therefore *order-dependent*. Order-dependent queries arise naturally in many application domains. In finance, an analyst often looks at *n-moving averages* over price time series, which is the average of a price and its  $n$  predecessors, calculated for each price in the series [19]. The analyst may also be interested in correlations among time series, which requires prices to be in time order [40]. In network management, an administrator may want to analyze packet logs for statistics or security purposes. Statistics may involve breaking sessions between any pairs of hosts down into “flows” (sub-sessions), a flow-separator occurring whenever a packet and its predecessor are more than a given time interval apart [9]. A security check of the log may look for a port scanning attempt, in which a same client sends a succession of packets to different ports on a given host [10]. Again, packet ordering is relevant. In the relational storage of XML, the order of XML elements and attributes need to be encoded [37]. In Biology, frequent nucleic acid motifs are of interest. In Linguistics, texts are scanned linearly. In epidemiology, unusual spikes in emergency room visits may suggest the start of an epidemic. The reader may imagine many other applications.

Many queries whose natural formulation requires order can be expressed using query languages based on multisets. For instance, suppose a stock’s quotes and their timestamps are stored in a table `Quotes` and that one wants to obtain each quote’s predecessor. Joining `Quotes` with itself (using as a predicate the maximum timestamp that is less or equal to the current quote’s timestamp) would give the desired result. If done often enough, a reasonably skilled SQL writer would recognize the predecessor idiom at once. As a practical matter however, the more structurally complex a query’s rendition is (joins or nested sub-queries), the more

difficult it is to optimize (join elimination or query un-nesting).

On the other hand, multiset query languages have a long and illustrious history. Order should therefore be inserted in a language through careful design. This thesis is the result of our studies into building such a language. We present here not only the language itself, but also the underlying data model that makes it coherent, the optimization techniques that make it efficient, and the system that implements it all. We call this framework AQuery.

## 2.2 Principles and Goals

AQuery is a language in which order-dependent queries can be expressed naturally. AQuery’s design began with these principles:

- **Declarative Order** – A query is able to define the order it requires records to be processed, regardless of the way the records are stored. One implication is that a query’s results can be made independent of the underlying storage strategies. For instance, stocks’ quotes are naturally generated in time order and could conceivably be stored so. But queries may require quotes to be in stock ID and time order; queries may simply declare so.
- **Ubiquitous Order** – The assumption that data is in a declared order is valid everywhere in a query, be it in a calculation involved in the result, a filter that depends on order, or still a grouping or aggregation operation. Order in AQuery can simply be counted upon everywhere.
- **Conciseness and Lucidity** – The most common order idioms are easily expressed – without auxiliary language constructs – in AQuery. The order idioms are not buried under or masked by a query’s structure. For instance, if a filter is dependent on the order declared in its query and order is needed in no other context, then this is evident from the query structure alone. The rationale here is that if a query looks simple, then it should be easy to understand and optimize.
- **SQL Compatibility** – SQL practitioners should be able to familiarize themselves with AQuery with very little effort.

These principles provide a natural basis for optimization. By declaring the order it requires, a query allows the optimizer to determine whether data is already organized in a convenient order. If it is, sorting may be eliminated altogether.

If data is not in an appropriate order, then the AQuery optimizer may still have some space to maneuver. For instance, if a join followed by an aggregation are involved in a query, but only the latter is order-dependent, then the optimizer may pick among several alternatives. It may perform the join in an “order-cavalier” fashion and then sort immediately before the aggregation. Or it may sort the relations in an appropriate way, perform an order-preserving join, and then finally aggregate. Or it may take advantage of the underlying order the relations are stored in and propagate that order until the aggregation. This choice is a cost-based one.

In many ways, we integrate the results from previous work: some focussed on linguistic aspects alone without considering optimization [30]; others did consider optimization but would not allow order handling in a declarative fashion [32, 22]; still others focused on order management in query optimization but for non-order-aware languages [34] or for languages with simple order extensions [36].

To the best of our knowledge, AQuery is the first comprehensive effort at addressing order-dependent queries' needs from the data model to the query language to the optimization process. As a result, AQuery expresses most order-dependent queries evaluated in a more concise way than other languages, and has delivered orders of magnitude performance gains for order-dependent queries as compared to robust commercial SQL:1999 query optimizers.

## 2.3 Thesis Overview

The first chapter of this thesis investigates the difficulties that arise in writing order-dependent queries in SQL. It then brings a comparative study of new languages that address that difficulty. From the study, we draw a powerful set of order-manipulation mechanisms but we conclude that such a set is not entirely available in any given language. We also observe that the most successful attempts at supporting order were those that replaced multisets in the data model by some form of array.

Chapter 3 defines a data model built around arrays that not only supports all the relational operators, but also order-preserving variations of them. The model provides several other features that are essential for expressing order-manipulations of realistic queries. The chapter also describes the AQuery language syntax and semantics, the latter by mapping queries into the algebra defined. The languages studied previously are then compared against AQuery on the basis of simplicity and conciseness.

The optimization of order-dependent queries is presented in Chapter 4. The chapter starts by reviewing the techniques applied to order (sort) management and shows that they also apply to AQuery. However, order optimization is usually done in a pre- or post-plan enumeration step. We advocate considering sort as any other operator in the enumeration process and present a set of new query transformations using this approach. The new transformations bring orders of magnitude improvement to the performance of some plans. For each new transformation we conduct a performance evaluation study.

Chapter 5 describes the architecture and design of a system that implements the AQuery language and model. It is a fully functional system. Its most salient design aspect is that it is fully based on vector-processing techniques. The system manipulates data in a vertically partitioned fashion and translates queries into sequences of vector-operations over columns. The system itself is written in the very same vector-oriented language to which it translates queries.

A performance comparison, both qualitative (query plans) and quantitative (response times), between AQuery and an commercial SQL:1999 optimizer is described in Chapter 6. SQL:1999 gained order-manipulation capabilities through

a late amendment [17] and therefore constitutes the first order-aware language to gain commercial acceptance. The comparison shows that AQuery is more amenable to order optimization, for its order idioms are more compact and easier to identify. Its plans were simpler than and in some cases almost two orders of magnitude faster than SQL:1999's.

Finally, chapter 7 concludes by describing our ongoing work and future research possibilities.

# Chapter 3

## State of the Art

### 3.1 Introduction

This chapter opens our investigation into querying ordered databases by addressing the following questions: Does a query language need specific features in order to express order-dependent queries? If yes, which ones? Is there any already existing query language with such features?

The best way to understand how hard – and why – it is to formulate queries that involve order is to try a few examples in SQL:92. Take a table Sales(month, sales) that stores for each month the amount sold in that month. For simplicity, assume that months are represented by integers starting with 1. Now, try to find the difference in sales between each month and its predecessor. (We encourage the reader to spend a minute trying to write this query.)

The first problem lies in the “predecessor” part of the query. How does one express such a property in a language that is based on (multi-) sets? Elements in a set are not ordered unless an ordering relationship is provided. Here the order can be obtained through months numbers. Again, for simplicity, let us assume that months are numbered consecutively with no gaps. Therefore one can find a previous month by simply subtracting 1 from the current month. Let us leave the boundary issue of month 1 aside for a brief moment.

The second problem is that *row-oriented languages* such as SQL – in which variables iterate over rows – cannot easily access two sales amounts (rows) at once. To subtract a previous month’s sales from the current’s, both sales must be in the same row. Producing such a configuration requires a self-join such as the following.

```
[SQL:92] – Delta Sales Query (naive)
SELECT t1.month, t1.sales - t2.sales AS delta
FROM   Sales t1, Sales t2
WHERE  t1.month - 1 = t2.month
```

Note that the first month has no predecessor thus it was eliminated by the join. Fixing it requires replacing the join by an outer join and handling the case where t2.sales would be null.

Sales	month	sales
	1	100
	2	120
	3	140
	4	140
	5	130

Result	month	delta
	2	20
	3	20
	4	0
	5	-10

Figure 3.1: A Sales table instance and the result of the delta sales query

**[SQL:92]** – Delta Sales Query

```
SELECT t1.month,
       t1.sales - CASE
                   WHEN t2.sales is null THEN 0
                   ELSE t2.sales
                   END
FROM   Sales t1 LEFT OUTER JOIN Sales t2
       ON t1.month - 1 = t2.month
```

We call such a manipulation of column values at different rows *Inter-row Operations*. In this particular case we are executing a *running delta* operation. It is one of the most common order manipulations.

Another quite frequent operation category is that of the *Running Aggregate Operations*, which brings together a number of cumulative operations such as moving averages. A *moving average* over a column divides its elements into groups of rows that are not necessarily disjoint (“windows”) and calculates the average of each group. For instance, if one wanted to compute a 3-month sales moving average, one would find the two previous sales for each month, would form a group with the current month, and would then calculate the average of that group. This operation would be repeated for each month.

The previous query can be used as the basis for this one, except that an extra outer join would be required so to have the 3-sales window. Managing the null values becomes more complex, too. The new query looks like

**[SQL:92]** – 3-month Sales Moving Average Query

```
SELECT t1.month,
       t1.sales + CASE
                   WHEN t2.sales is null AND
                        t3.sales is null
                        THEN 2*t1.sales
                   WHEN t2.sales is not null AND
                        t3.sales is null
                        THEN t2.sales + (t1.sales+t2.sales)/2
                   ELSE t2.sales + t2.sales
                   END
FROM   Sales t1 LEFT OUTER JOIN Sales t2
       ON t1.month - 1 = t2.month
LEFT OUTER JOIN Sales t3
       ON t1.month - 2 = t3.month
```

Despite the complex syntax of their SQL renditions, the queries so far are quite simplistic. To show a more realistic example based on the operation categories presented so far, we turn our attention to stock trading. An increasingly popular way of trading in a very dynamic market is “day-trading.” It consists of buying and selling shares within very short periods of time – often in the very same day – aiming at immediate albeit modest profits. To verify whether a trader made a profitable enough transaction, one may wish to know what would be the best profit that could be done by buying and then selling a given stock in the same day.

In this query we use the schema Ticks(ID, date, volume, price, timestamp), which stores all quotes and trades for securities negotiated in a trade floor or trading system. ID is the ticker symbol of a security, the code by which a security is identified; timestamp is the date and time of a particular quote or trade; date is the human-readable form of the day portion of the timestamp; volume is the number of shares that are being negotiated; and price is the value per share.

A strategy to solve this query would be to calculate, for each row of Ticks, what would have been the minimum price seen for that security up until that tick. This operation is called a running minimum – another running aggregate – and requires ticks to be sorted by ID and timestamp order. Subtracting each Tick’s price by its running minimum yields the profit made by buying at the lowest price up until that tick and selling at that time. The maximum of such differences is the best profit. The graph in Figure 3.2 illustrates this approach by showing the price curve of a stock on a given date along with its price running minimum curve. The difference between the curves depict profits at each tick. For brevity, we omit the solution of this query in SQL:92.

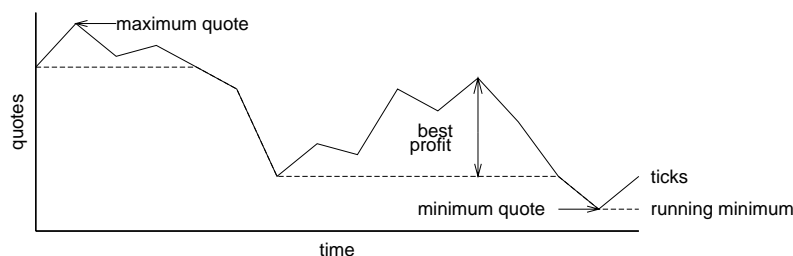


Figure 3.2: A stock’s price curve and its running minimum

The above examples show that writing queries that deal with order may result in complex SQL. That is not without consequences to the query’s optimization process – upon which declarative query languages depend so heavily to be efficient. From a strictly algorithmic point of view, all the queries shown here could use solutions (query plans) with linear time complexity in the database size. Had the tables been in an appropriate order, running deltas and running minimums would have required a single pass on the data. Yet whenever expressed in SQL, even state-of-the-art commercial optimizers are not able to find such plans and eliminate the (algorithmically unnecessary) joins.<sup>1</sup>

<sup>1</sup>As of this writing, the DBMS products we tested picked join-based (non-linear) solutions.

We can now go back to our first question. Do query languages need specific mechanisms to express order operations? If one wants order-dependent queries to be easy to read and to optimize, then yes. Now, which mechanisms would accomplish that?

Supporting order-dependent queries means supporting an ordered data structure and providing a query language and algebra equipped to deal with it [24]. There is not consensus on whether that structure is a list, an array, an arbitrary combination thereof, or still some other structure or on whether that query language should be based on SQL.

Several alternatives have been investigated and we can find inspiring insights in the database literature. We divide these works in three categories: those that added ordered structures and adapted SQL accordingly, those that designed a whole data model and a language around a given structure, and finally the standard SQL approach itself, which we describe next.

## 3.2 Standard SQL with Late Order

Order support in SQL:1999 [16] comes indirectly from a language mechanism whose purpose is to define a *sliding window*. This mechanism was in fact a late addition to the standard [17] and was conceived more to deal with analytical queries (OLAP) needs than to address more general order-dependent queries.

One way to present SQL:1999's sliding windows is as a means of transforming aggregate functions into running aggregates. For example, to compute the 3-month sales moving average query stated previously, we define a window whose size is of 3 positions (rows), starting two months before the current one and finishing in the latter. This definition requires months to be sorted on ascending order. Applying the aggregate function `avg()` over each distinct window yields the desired running average. This query in SQL:1999 would look like

```
[SQL:1999] – 3-month Sales Moving Average Query
SELECT month, sales, avg(sales) OVER ( ORDER BY month
                                     ROWS BETWEEN 2 PRECEDING
                                     AND CURRENT ROW)
FROM Sales
```

This query may not look familiar to the SQL practitioner. Having in a `SELECT` list the aggregate function `avg()` and columns that were not `GROUP BY` may look at first as an error. It is not; this query returns the same result as its SQL:92 counterpart. The `OVER` clause specifies a sliding window which modifies the function `avg()` to a running aggregate. The change is less in the way an average is computed than in how many times this function is called. The `OVER` clause associates to each row a window and then causes the aggregate function to be called on that window. The net effect is that `avg()` is called as many times as there are rows.

Although being quite expressive, window usage is subject to a few syntactical restrictions [27]. For one, they can be used only in the `SELECT` clause. For

another, they can modify only aggregate functions. Thus, for instance, to find a “predecessor” element as was the case in the delta sales query specified previously, some creativity is involved. Defining a single-element sliding-window consisting of the element that precedes the current and applying `min()` to such a window would return that very element. The delta sales query in SQL:1999 looks like

```
[SQL:1999] – Delta Sales Query
SELECT month, sales - min(sales) OVER ( ORDER BY month
                                      ROWS BETWEEN 1 PRECEDING
                                      AND 1 PRECEDING)
FROM Sales
```

Some DBMS vendors provide non-standard “windowed functions” that make the above query more intuitive. For instance, Oracle 9i has a pair of functions, `lag()` and `lead()`, that retrieve values by their relative position [29]. The standard itself defined other additional functions specifically to be used over windows, but the presentation of those is beyond the scope of this discussion.

Of interest to our study is to evaluate whether the window mechanism itself helps or hinders order-dependent queries. One way to assess it is to write the best profit query in SQL:1999. Recall that it relied on finding a running minimum of prices for each stock. The difficulty here is to define a sliding window that takes care not to mix ticks of different stocks. Enforcing a sort order or a window size alone would not suffice. We therefore resort to a more advanced way to define a window which involves partitioning data. Our point becomes clearer if we try to find the best profit of all stocks rather than just of a single one. This query in SQL:1999 looks like

```
[SQL:1999] – Best Profit Query
SELECT ID, max(running_diff)
FROM ( SELECT ID, date,
              price - min(price) OVER ( PARTITION BY ID, date
                                       ORDER BY timestamp
                                       ROWS UNBOUNDED PRECEDING)
        AS running_diff
      FROM Ticks ) AS t1
WHERE date = '05/11/2003'
GROUP BY ID
```

To transform the aggregate function `min()` in the desired running minimum, we partition data by ID and date. This means that the current row’s stock ID and quoting date values are considered when forming that row’s window. Only after eliminating all elements that do not belong to the current row’s partition does the sort order get enforced. We also used a cumulative way to form a window in which it starts at the first row of the current partition seen and ends at the current row (`ROWS UNBOUNDED PRECEDING`).

Because an aggregate function (`max()`) cannot take a running aggregate function (`min(price) OVER ...`) as an argument, this query is a nested one. The inner

block of the query calculates the running difference while the main query performs grouping and aggregation. The impact of nesting is felt both on readability and on optimization. We show in Chapter 6 how a commercial optimizer missed optimization opportunities due to the structure of the query.

Sliding windows are not the only mechanism in which order is involved in SQL:1999. The language incorporated a new array type; columns can now hold arrays as opposed to just scalars. The manipulation of “array-fields” in SQL:1999 is very limited, though. There were academic prototypes that exploited arrays fields better than SQL:1999. These will be discussed in the next section.

### 3.3 SQL Dialects over Ordered Structures

Academic prototypes had defined dialects of SQL with order features before SQL:1999 was amended with OLAP functions. SEQUIN, or the PREDATOR system to be precise, suggested a pacific cohabitation of sequences and relations through the use of “enhanced” ADTs [32]. SRQL and its underlying algebra did not make distinctions between relations and sequences, treating the former as a degenerate case of the latter [30]. Let us discuss each language in turn.

#### 3.3.1 SEQUIN

The PREDATOR database system introduced the concept of “enhanced” abstract data types (E-ADT) [32]. By enhanced it meant that each new data type carried more than a collection of methods that could manipulate it. It carried a particular query language and even an optimizer of its own. SEQUIN is PREDATOR’s language designed to deal with its sequence E-ADT.

A sequence in a SEQUIN query occupies the same place a table would in SQL. The main difference is that sequences are ordered [31] and thus the FROM clause supports special joins based on the ordering domain of the participating sequences. This feature comes in handy to express the delta sales query. We show the SEQUIN rendition of the query below

```
[SEQUIN] – Delta Sales Query
PROJECT t1.month, (t1.sales - t2.sales)
FROM Sales AS t1, PREVIOUS(Sales) AS t2
```

When more than one sequence appears in the from clause an implicit natural join on their positions takes place, i.e. first element joins with first element, second element with second element, and so on. The modifiers NEXT, PREVIOUS, and OFFSET in the FROM clause shift entire sequences so that different alignments can be used in the joins.

The PROJECT clause is similar to the SQL’s SELECT. It implicitly assumes that elements of the sequences are ordered. Herein lies both an advantage and a problem. The advantage is that enforcing order early in a query makes order visible to all subsequent clauses. For instance, had we wanted to find the months whose delta sales were greater than a given threshold, we would have put the

expression 't1.sales - t2.sales > threshold' in the WHERE clause. In SQL:1999, where order is supported only in the SELECT clause, this query would have required nesting. The disadvantage of implicit order is that order is not declarative. By looking at the query alone one does not know in which order the Sales sequence is. Should the order of Sales be changed for some reason, the result of this query would have been affected. Had the query specified in which order it expected to handle data, such problem would not have happened.

An interesting order feature is that of the mixing of SEQUIN (sequences) and SQL (sets) in a query. In the following example suppose Ticks' is a table that holds two columns: an ID one, and a sequence of quote-price elements in timestamp order called priceSeq. To query the minimum price of every security within a given interval one would write in SEQUIN the following query.

```
[SEQUIN] – Mixed SQL and SEQUIN Query
SELECT ID, SEQUIN ( "PROJECT min(price)
                    FROM      $1
                    WHERE     date BETWEEN
                               '01/01/2003' AND '03/31/2003" ,
                    T'.priceSeq)
FROM   Ticks' AS T'
```

The keyword 'SEQUIN' in the query could be seen as a function call to SEQUIN's query processor. Its first argument is the query to be executed over sequence data. The second argument is the sequence data itself corresponding to ticks of the security's row that is being processed. The \$1 in the query is replaced with that parameter.

This is quite a useful feature for applications that manipulate sets of sequences, e.g. Finances and Biology. Nevertheless, nesting an entire query in the SELECT list of another is arguably hard to read.

### 3.3.2 SRQL

SRQL is a SQL dialect that also manipulates sequences and relations. It does so by representing both by the same underlying structure. A sequence is a sorted relation, in which a list of attributes exists that defines the order of its records. A relation is simply a degenerate case in which such a list is empty. That and other simplifications made the writing of order manipulations in SRQL easier.

A SRQL query structure resembles that of a SEQUIN query structure. The main improvement was to allow a query to declare in which order it requires data to be processed rather than to assume that such order comes automatically from the underlying data structures. The chosen order is enforced early in the query and, as in SEQUIN, all remaining clauses can count on it. SRQL has operators to shift sequences much as SEQUIN does, but they are not confined to the FROM clause. To illustrate the use of these mechanisms we show the SRQL rendition of the delta sales query below

```
[SRQL] – Delta Sales Query
```

```

SELECT S.month, (S.sales - SHIFT(S,-1).sales)
FROM   Sales
      SEQUENCE BY month AS S

```

In the query, the Sales relation is “sequenced” according to month order and the result is bound, a row at a time, to the tuple variable ‘S’. The expression ‘SHIFT(S,-1).sales’ reads “the previous tuple to the one pointed by the tuple variable *S*.” SRQL is the first language we present that is able to express this query without explicitly resorting to a join. From an optimization point of view such a query rendition makes the order idioms clear – it exposed both the order in which rows are to be processed and what order manipulations are involved in the query.

SRQL also supports the concept of sliding windows. The window definition, much simpler and thus less powerful than SQL’s, has to adopt the same order specified in the SEQUENCE BY clause of the query. The following SRQL rendition of the 3-month sales moving average shows the use of windows.

```

[SRQL] – 3-Month Sales Moving Average Query
SELECT month, AVG(sales) OVER -2 TO 0
FROM   Sales
      SEQUENCE BY month

```

The OVER clause defines a window based on relative positions. In this case, avg() will be called once per row and will be passed a window that ranges from two rows before the current (-2) until the current row (0).

We don’t know the personalities involved, but it seems likely that SQL:1999 borrowed the OVER constructs from SRQL. The concept here is clearer, though, because the sort order is valid throughout the entire query as opposed to just within the window. It remains that, for the same reason as SQL:1999, queries such as the best profit one would require nesting to be expressed in SQL.

SRQL has borrowed several mechanisms from SEQUIN, but it gave up an important one: the ability to manipulate order (sequences) within a field.

## 3.4 Array-Based Querying Systems

Arrays are ubiquitous in application domains such as scientific computing and finance. Attempts were made to develop data models and query languages that revolve around arrays that would support such applications. Order is intrinsically involved in array manipulation and therefore we investigate such effort.

### 3.4.1 AQL

Array Query Language (AQL) is a language that manipulates multidimensional arrays [22]. Underneath AQL there is a data model based on the nested relational calculus of [7] with the addition of primitives to handle arrays.

AQL syntax is based on set comprehensions [6]. A comprehension has the form { *f* | *q*<sub>1</sub>, *q*<sub>2</sub>, ···, *q*<sub>*n*</sub> }, where qualifiers *q*<sub>*i*</sub> can be either predicates or generators.

A generator is an expression of the form  $x \leftarrow A$  (reads  $x$  draws from  $A$ ) that sequentially binds elements of the collection  $A$  to variable  $x$ . The qualifiers are evaluated from left to right. A binding is propagated until a predicate evaluates to false under this binding. The function  $f$  in the head of the comprehension is evaluated under all bindings that survived the predicates. The AQL rendition of the delta sales query shown below illustrates the use of comprehensions.

```
[AQL] – Delta Sales Query
{ ( month, sales[i] - sales[i-1] ) |
  [ \i : (\month, \sales) ] ← Sales,
  i > 0 }
```

The first qualifier is a special generator that binds array elements to variables. We assume here that `Sales` is an array of records (`month`, `sales`). The generator binds at once the position of a record to the variable `'i'`, and the record's components to variables `'month'` and `'sales'`. A backslash preceding a variable signals that the latter is being bound at that qualifier. The second qualifier is a predicate that filters out the very first record. Array indexes in AQL start at 0. This will avoid an out-of-bounds error in the head of the comprehension. It calculates the delta for a given month and puts the result in a record format.

Note that AQL generators use implicit order. The language could easily be extended with a sort user-defined function that would explicitly enforce the desired order for a query.

AQL is the first language described so far that supports an indexing operation – that is, a positional access to an element through a subscript. The previously introduced languages do support some notion of row numbering but would not allow it to be used to directly access an element. As the query demonstrates, such a feature comes with the responsibility of avoiding out-of-bounds errors.

### 3.4.2 KSQL

KSQL is the SQL dialect of the database management system KDB [21]. KDB has a fully vertically partitioned implementation of tables in which each column is a 1-dimensional array. For that, KDB tables have been called *arrables* (array-tables) [38].<sup>2</sup> An arrable's rows are intrinsically ordered.

What makes KSQL very effective for queries that other languages can express only through nesting is its column-oriented semantics. It allows a single block rendition of the best profit query.

```
[KSQL] – Best Profit Query
SELECT max(price - mins price) BY ID
FROM Ticks
```

The variables in the query refer to entire columns rather than to single values. The construct `'BY'` has the same effect of SQL's `GROUP BY`. Therefore `'price'` in the `SELECT` clause is bound to a vector of prices partitioned by `ID`. The function

---

<sup>2</sup>This term was actually first coined in [33].

mins() is being implicitly called once for each price partition and its result is subtracted from that very same price partition. The function max() then returns the highest difference of price for that partition.

KSQL can express a large set of order-dependent queries in a very concise way. But the language does not support declarative order; a change in the underlying order of Ticks would affect the result of this query. In addition, KSQL’s syntax is somewhat far from the classical structure SQL practitioners are used to seeing.

### 3.5 Discussion

The investigation we have presented here is far from exhaustive but brings a significant cross-cut of order-manipulation mechanisms in existing query languages. We briefly cite other pertinent work. Array manipulation, in particular image pixel arrays, motivated at least two other works. AML [26] is a framework for generic function application over multidimensional arrays. It is based on a small, flexible set of algebraic operators in which several image manipulations can be encoded. RasQL [2] also manipulates image pixel arrays but using a SQL dialect. Both languages do not support declarative order. Ordered relations were also used elsewhere. In [28] the relational model was extended by providing the facility of user-defined orderings over data domains. The resulting model was called Ordered Relational Model. It supports a variation of SQL called Ordered SQL that could express the queries seen here, although with a peculiar syntax. Finally, sequences were used in an extension of Datalog called Sequence Datalog [4]. The latter kept the same syntax as Datalog but added two types of enriched terms, one capable of extracting sub-sequences from a given sequence (indexed sequenced terms) and the other capable of concatenating sequences (constructive sequence terms). Order wise, these works did not add any new mechanism that wasn’t mentioned here.

The table 3.1 presents a qualitative comparison of the languages described previously. For each language we state which underlying data model and data structure it is based upon and what querying style and semantics it adopts. We also determine for each of them what order manipulation mechanism they incorporate. By “declarative order” we refer to the ability of explicitly stating the order in which data is supposed to be processed by a query. By “running aggregates” we mean the support for cumulative aggregates (e.g., running minimum) and windowed aggregates (e.g., moving average). By “inter-row” we refer to the ability of using values of two or more distinct rows in a single expression. Finally, “best profit query” shows whether a language’s rendition of it requires a nested structure.

The table shows that there is no absolute best language. That answers our final question regarding whether any one given language would incorporate all necessary order manipulation mechanisms. AQL is arguably the most expressive of the languages – comprehensions have shown to be as expressible as languages such as OQL [14]. The missing linguistic features – declarative order and running aggregates – could be implemented through the addition of user-defined functions.

We show in Chapter 4 that such user additions may be less than satisfactorily handled by the optimization process. We also have a bias for pragmatic reasons for SQL and its dialects, but reasonable people can differ on this point. KSQL has proven to be quite appropriate for order-dependent queries and was the only SQL dialect capable of expressing the best profit query in a single block.

Linguistic differences notwithstanding, arrays were the structure that provided the greater flexibility. Any operation that can be done on an ordered set, on a sequence, or on a table can also be done on an array representation of those structures. The inverse doesn't hold. Moreover, array indexing (positional access) can be used as a primitive to express every order-manipulation mechanism studied here. It was shown that indexing comes with the risk of having run-time (out of bounds) error, though.

In conclusion, we claim that arrays are the best suited data structure to support order-manipulations. We seek a language that can take advantage of the flexibility of arrays, and can express order-dependent queries in a natural, clear way.

	Underlying Data Model	Main Data Structure	Query-Language Style	Semantics
SQL:1999	Relational	table	SQL	row
SEQUIN	Object-Relational	sequence	SQL dialect	row
SRQL	Ordered Relations	ordered sets	SQL dialect	row
AQL	Nested Relational	arrays	comprehensions	row
KSQL	Relational	arrable	SQL dialect	column

	Declarative Order	Order in all clauses	running aggregates	inter-row	best profit query
SQL:1999	yes	no	yes	window	nested
SEQUIN	no	yes	yes	join	nested
SRQL	yes	yes	yes	yes	nested
AQL	UDF	yes	UDF	yes	single block
KSQL	no	yes	yes	yes	single block

Table 3.1: Comparative table of languages with order constructs

# Chapter 4

## AQuery Syntax and Semantics

### 4.1 An Array-Based Data Model

AQuery is a SQL dialect that is based upon an ordered data structure called an *arrable*, for array-table. Informally, an arrable is a collection of named arrays that, in their simplest form, are vectors of elements of a base type. In this form, an arrable is essentially a table organized by columns. An arrable's arrays may assume more complex shapes, though. They may contain array-valued fields themselves, but nesting beyond this point is not allowed.

Ticks	ID	price	date	ts
	ACME	12.02	05/11/03	1
	WXYZ	43.23	05/11/03	2
	ACME	12.04	05/11/03	5
	ACME	12.05	05/11/03	9
	WXYZ	43.22	05/11/03	13

(a)

Series	ID	price	date	ts
	ACME	[12.02 12.04 12.05]	05/11/03	[1 5 9]
	WXYZ	[43.23 43.22]	05/11/03	[2 13]

(b)

Figure 4.1: Example of two well-formed arrables

Figure 4.1 shows examples of two well-formed arrables. Both store stock quotes, but each in a particular shape. The column ID refers to the security identifier of a quote, price is the quote itself, and date and timestamp record the moment the quote occurred. Formally, arrables can be described as follows.

**Definition 4.1 (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 of elements from a basic type. Let  $A$  be a finite array of elements of a type  $t \in \mathcal{T}$ . The cardinality of  $A$  is the number of elements in  $A$ 's first dimension. 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 an array of type  $t_i \in \mathcal{T}$ .  $\square$

**Definition 4.2 (Arrable Indexing)** – The  $k$ -th *row* 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  $r[k] = \langle A_1[k], \dots, A_n[k] \rangle$ .  $\square$

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

Because an arrable consists of arrays and arrays are ordered, an arrable's rows are ordered.

**Definition 4.3 (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  $B_1[k_1] = B_1[k_2], \dots, B_{i-1}[k_1] = B_{i-1}[k_2]$ . Or, put informally, the tuples in the ordered projection of  $r$  onto  $B_1, \dots, B_m$  are lexicographically ordered. The definitions are symmetric for descending orders, but for the purpose of exposition, we will consider order to be ascending throughout this chapter.  $\square$

For instance, the arrable `Ticks` shown in Figure 3.1(a) could be defined as `Ticks(ID, price, date, timestamp) ORDERED BY timestamp`.<sup>1</sup>

The AQuery language borrows its syntax from SQL but it permits a query to specify the order in which it wishes to process rows. It does so through a new clause called `ASSUMING ORDER` introduced between the `FROM` and the `WHERE` clauses. The `ASSUMING ORDER` clause's semantic effect is to sort data immediately after the Cartesian product indicated by the `FROM` clause.

**Definition 4.4 (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 `ORDERED BY`  $B_1, \dots, B_m$ .  $\square$

For instance, to find ACME's quotes ordered by timestamp, one would write

```
[AQuery]
SELECT price
FROM   Ticks
       ASSUMING ORDER timestamp
WHERE  ID = 'ACME'
```

---

<sup>1</sup>We are omitting the typing information here for convenience. A complete definition would include NULLs and referential integrity information also.

The order declared in the query is independent of that of the arrable to which it refers. A smart optimizer should be able to detect the coincidence and take advantage of it. (The topic of AQuery optimization will be further explored in Chapter 4.)

By enforcing order in the FROM clause of a query, all of the subsequent clauses may take advantage of it. We believe though that order throughout all clauses is a necessary, but not a sufficient condition for having clear, concise order-dependent queries.

## 4.2 Column-Oriented Semantics

One problem in expressing order-dependent queries is that often each resulting row is a combination of values of more than one input row. For example, consider a query to find the difference between each price and its previous value, assuming a time order. It needs to access two prices at once that are in distinct rows to calculate each pair's difference. *Row-oriented languages* such as SQL-92 can only iterate over only one row at a time, though. Thus they need to resort to either a self-join or 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 semantics* in that *variables are bound to entire arrays at a time*. 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 array such that  $\text{prev}_A[i] = A[i - 1]$  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 *order-dependent* 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 4.5 (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  $A_i^{\text{perm}}$ , then  $e(A_1, \dots, A_m)$  and  $e(A_1^{\text{perm}}, \dots, A_m^{\text{perm}})$  represent the same multiset. 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.  $\square$

Other functions in the order-dependent, size-preserving category are the running aggregates. A running minimum over an array  $A$ , mins( $A$ ), is  $\text{mins}_A[i] = \min(A[i], \text{mins}_A[i - 1])$  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  $\text{sums}_A[i] = A[i] + \text{sums}_A[i - 1]$  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  $\text{avgs}(w, A)$  and is defined as  $\text{avgs}_{w,A}[i] = \text{sum}(A[i - (w - 1)]..A[i])/w$ , for  $w - 1 \leq i < |A|$  or  $\text{sum}(A[0]..A[i])/i$  for  $0 \leq i < w - 1$ .<sup>2</sup>

Another category of vector-to-vector functions are those that are order-dependent but not size-preserving. They reduce an array's cardinality and as such are called edge functions (i.e., they keep either the beginning or the end of an array). For instance, the first  $n$  positions of an array  $A$ , denoted  $\text{first}(n, A)$ , is  $\text{first}_{A,n} = A[0..n - 1]$ . Similarly,  $\text{last}_{A,n} = A[|A| - n..|A| - 1]$ .

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

**Example 4.1** – The combination of column-oriented semantics and array-typed expressions make it easier to write the best profit query (chapter 3). Recall that this query uses the Ticks arrable of Figure 4.1 to find what would be the best profit one could make by buying and selling a given stock in a given date. This query required a nested formulation when written in row-oriented languages, even order-aware ones. In AQuery, it can be written in a single block.

```
[AQuery] – Best Profit Query
SELECT max(price - mins(price))
FROM   Ticks
       ASSUMING ORDER timestamp
WHERE  ID = "ACME" AND
       date = "05/11/2003"
```

The query should be read with a column-oriented mind-set. The FROM clause accesses the arrable Ticks and sorts it by timestamp order. In the WHERE clause, 'ID' and 'date' are both vectors; comparing each of them to scalars – ACME and 05/11/2003, respectively – are valid array-typed expression that result in two booleans arrays. After they are combined, the resulting array maps each position of Ticks to true or false. Processing the WHERE clause means eliminating the false positions.

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. Subtracting a vector (mins(price)) from another (price) with the same cardinality is a standard array expression as is taking the max() of the resulting vector.  $\square$

## 4.3 Relational Manipulation of Arrables

The AQuery algebra supports the operators of the relational algebra. But here each operator takes array-typed expressions as arguments. If an expression is

---

<sup>2</sup>Such a definition is commonplace in financial applications. Other domains may require  $\text{avgs}()$  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.

order-dependent, then the operator behaves in an order-preserving way. Otherwise the operator behaves in an order-cavalier way. We use the following order equivalence between arrays to define such behavior.

**Definition 4.6 (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  $r \equiv_{\{\}} s$ .  $\square$

The order-cavalier variation of an operator is simply one that is multiset equivalent to its order-preserving variation. In the remaining of the section we define the order-preserving variations of the relational algebra operators.

### 4.3.1 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.<sup>3</sup>

```
projection(e,r)
1. s:= empty arrable having the same schema as e
2. for i = 0 to |r|-1
3.     append <e1[i], ..., em[i]> 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  $\pi_e(r)$ .

### 4.3.2 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
3.     if p[i] is true
4.         append r[i] to s
5.     end if
6. end for
7. output s
```

---

<sup>3</sup>Note: After considering different formalization notations including set comprehensions and lambda calculus we have decided to use a simple minded but (to us) clear loop formulation. In fact comprehensions would have worked well for many of the operators, but introduced problems for some such as certain variants of order-preserving joins

As for a projection, a selection can be order-dependent, and either order-preserving or order-cavalier. This may be interesting if we have a hash index for example.

**Example 4.2** – Let  $e = \max(\text{price} - \min(\text{price}))$  and  $p = (\text{ID} = \text{'ACME'}) \wedge (\text{date} = \text{'05/11/2003'})$ . The Best-Profit query can be translated to the AQuery algebra as follows,

$$\pi_e^{op}(\sigma_p^{op}(\text{sort}_{timestamp}(\text{Ticks})))$$

□

### 4.3.3 Group By

Grouping in AQuery uses an arrable's facility to store array valued fields. Intuitively, grouping in AQuery partitions the operand arrable into disjoint sub-arrables 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 4.1 shows the effect of grouping the arrable Ticks in the same figure by ID and date.

Formally, 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 array A in  $r$
7.             if A is not a grouped-by column
8.                 concat  $r[i].A$  to  $s[j].A$
9.             end if
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 are. Group-by can also have an order-cavalier variation in which the assembled arrays in fields may not be in the same order as in the original arrable.

Grouping in AQuery is independent of aggregation. To apply a function to each array-valued element of a column, AQuery provides an operator modifier

called *each*. Formally, let the array  $A$  be a parameter (array) 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 can be naturally extended for cases where  $F$  takes more than one argument.

**Example 4.3** – Consider the schema Packets(pID, src, dest, length, timestamp), where pID identifies a packet exchanged between a source (src) and a destination (dest) host. Length refers to the size of the packet and timestamp to the moment this packet was exchanged. A “flow” from a source  $s$  to a destination  $d$  ends whenever there is a 2-minute gap between consecutive packets from  $s$  to  $d$  [9]. Suppose a network administrator wants to know the count of packets and their average length within each flow. This query would need to group packets of each flow, and compute the count and average needed. Finding the flows is very hard to express, though, because it involves order.

In AQuery, such a grouping expression corresponds to the arrable 'src, dest, sums(deltas(timestamp)>120)'. 'deltas(col)' is the abbreviation of 'col - prev(col)'. Figure 4.2(a) shows how this expression is computed, supposing the Packets arrable is sorted over src, dest, and timestamp. The expression 'delta(timestamp)>120' finds for each packet whether it starts a new flow. Assuming that the boolean TRUE carries a value of 1, and FALSE of 0, the expression 'sums( deltas(timestamp) > 120)' generates a unique flow identifier, when concatenated with src and dest.

The arrable we see in Figure 4.2(b) is the grouped one. Note that the columns of Packets that are not columns of  $g$  (the grouping expression) have arrays within fields. Because fields may be arrays (though not arrables), aggregate functions may apply over an entire column or over each field. In figure 4.2(c) we see that avg() was applied to each of the array-values of the column length.

The AQuery rendition is given below.

```

[AQuery] – Network Management Query
SELECT src, dest, avg(length), count(timestamp)
FROM   Packets
       ASSUMING ORDER src, dest, timestamp
GROUP BY src, dest, sums(deltas(timestamp) > 120)

```

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

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

□

Packets	src	dest	length	ts	deltas(ts)>120	sums(deltas(ts)>120)	
	s1	s2	250	1	F	0	] g1
	s1	s2	270	20	F	0	
	s1	s2	235	141	T	1	□ g2
	s2	s1	330	47	F	1	] g3
	s2	s1	280	150	F	1	
	s2	s1	305	155	F	1	

(a)

Packets'	src	dest	length	ts
	s1	s2	[[250, 270]]	[[1, 20]]
	s1	s2	[[235]]	[[141]]
	s2	s1	[[330, 280, 305]]	[[47, 150, 155]]

(b)

Packets''	src	dest	avg(length)	count(ts)
	s1	s2	260	2
	s1	s2	235	1
	s2	s1	305	3

(c)

Figure 4.2: Intermediate arrables in the Network Management Query

### 4.3.4 Flatten

Flatten generates a first normal form equivalent of an arrable that contains array-fields. It requires every row of the input arrable to be made of scalars, or of scalars and array-valued fields where the latter are of the same cardinality.

To define the flatten operation formally, assume that  $\text{card}()$  is a function that, given a row, returns the maximum cardinality of any of the row's elements. Further let  $r$  be an arrable made of arrays  $A_1, \dots, A_n$  and let there be a  $m$  such that  $A_1, \dots, A_m$ , are vectors (contain only scalar elements) and  $A_{m+1}, \dots, A_n$  contain array fields as described earlier. Flatten over  $r$  could be defined as follows

`flatten(r)`

1.  $s :=$  empty arrable having the same schema as  $r$
2. for  $i = 0$  to  $|r|-1$
3.     for  $j = 0$  to  $\text{card}(r[i])$
4.         append  $\langle r.A_1[i], \dots, r.A_m[i], r.A_{m+1}[i][j], \dots, r.A_n[i][j] \rangle$
5.     end for
6. end for
7. output  $s$

**Example 4.4** – Financial analysts often observe stock tendencies before making purchase decisions. Moving averages are capable of smoothing the volatile stock

price curves and exposing underlying optimistic and pessimistic sentiment. For instance, whenever a short-term trend curve (a 5-day moving average) crosses above a longer-term one (21-day moving average) technical analysts would suspect the stock will move up soon.

The following query would be involved in this analysis.

```
[AQuery] – Crossing Averages Query
WITH
  averages (ID, date, a21, a5) AS
  (SELECT ID, date,
         avgs(21, price) as a21,
         avgs(5, price) as a5
   FROM   Ticks
   ASSUMING ORDER ID, timestamp
   GROUP BY ID)
SELECT ID, date
FROM   FLATTEN(averages)
       ASSUMING ORDER ID, timestamp
WHERE  a21 > a5 AND
       prev(a21) <= prev(a5) AND
       prev(ID) = ID
```

This query finds the dates where the 21-day and the 5-day moving average for a given set of stocks cross. The WITH construct from AQuery was borrowed from SQL:1999. It defines a “local view” that can be referenced only in the FROM clauses of subsequent WITH queries or in the main query. Note that the view returns the averages as array fields for each ID and date (non-1NF arrable). The next step is simply to check crossings.

Let  $e = ID, date, avgs(21, price), avgs(5, price)$  and  $p = a21 < a5 \wedge prev(a21) > prev(a5) \wedge prev(ID) = ID$ . The crossing averages query is translated to the AQuery algebra as follows

$$r \leftarrow \pi_e^{op}(gby_{ID}^{op}(\text{sort}_{ID, timestamp}(\text{Ticks})))$$

$$\pi_{ID, date}^{op}(\sigma_p^{op}(\text{flatten}(r)))$$

□

### 4.3.5 Cross Product and Join

Cross-product ( $\times$ ) in AQuery is order-cavalier and hence has the same definition as in the relational algebra. By contrast, joins may have several variations depending on whether and how the order of the input arrables is preserved. Let  $r(A_1, \dots, A_n)$  and  $s(B_1, \dots, B_m)$  be arrables. A *left-right order-preserving* join of arrables  $r$  and  $s$  on join predicate  $p$ , denoted  $r \bowtie_p^{lrop} s$ , is defined in the following way.

```

join(p, r, s)
1. o:= empty arrable with schema  $\langle A_1, \dots, A_n, B_1, \dots, B_m \rangle$ 
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  $\langle A_1[i], \dots, A_n[i], B_1[j], \dots, B_m[j] \rangle$  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. Suppose that  $r(A_1, \dots, A_n)$  is the arrable for which order should be preserved. A *left order-preserving* join,  $r \bowtie_p^{lop} s$ , is one that is order-equivalent with respect only to  $A_1, \dots, A_n$  to a left-right order-preserving join of the same two arrables.

**Example 4.5** – The arrable Portfolio(ID, tradedSince) ORDERED BY ID, stores information about the stocks that makes one analyst's portfolio. It is a subset of the stocks that appear in Ticks. If this analyst wanted to extract the ten last quotes for each stock that he or she traded, then the following query could be issued.

```

[AQuery] – Non-1NF Result Query
SELECT t.ID, last(10, price)
FROM   Ticks t, Portfolio p
      ASSUMING ORDER timestamp
WHERE  t.ID= p.ID
GROUP BY t.ID

```

Semantically, the query first performs a cross-product ( $\times$ ) between Trades and Portfolio. As mentioned before, Cross-product in AQuery is order-cavalier. Next, the ASSUMING clause imposes the desired sort order and the join predicate is applied. 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 price 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 follows. (Note that last() takes two arguments and therefore that is reflected on the syntax of the 'each' call.)

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

□

## 4.4 Positional Manipulation of Arrables

AQuery exploits the ordered nature of arrables so as to support the referencing of rows by their position. We describe two order-manipulation mechanisms that use this facility.

### 4.4.1 Querying with Arrable Indexing

Because expressions in AQuery are array-typed, it is only natural to allow indexing (i.e., access to an array's element given its position) in the language. Yet an access to a non-existent position would cause a run-time error. AQuery avoids such a possibility by introducing the notion of safe index sequence generators (SISG). A SISG is an expression that is always evaluated to a valid sequence of indexes. For instance, the expression 'price[ODD]' uses the SISG 'ODD' to return all prices from position 1 until the last odd position in price in that context. Other SISG is EVEN, which works similarly to ODD but starts at 0 and uses only even positions. We use an example to introduce the SISG EVERY n.

**Example 4.6** – Suppose a financial analyst wants to know the standard deviation of prices for ACME's stocks and whether it would be accurate to work with samples of its prices instead. The following query would return the standard deviation for the price column of Ticks, for samples at every 10th price, and at every 100th price. The function stddev() is a built-in one and calculates the standard deviation for a vector.

```
[AQuery] – Array Indexing Query
SELECT stddev(price), stddev(price[EVERY 10]), stddev(price[EVERY 100])
FROM   Ticks
      ASSUMING ORDER timestamp
WHERE  ID = 'ACME'
```

□

### 4.4.2 Querying with Row Direct Addressing

After a query's FROM clause is evaluated, the resulting arrable implicitly gains an additional column (vector) called ROWID. This synthetic column can be referred to anywhere a regular column can.

**Example 4.7** – Good candidate stocks for day-trading may be among the most early-traded stocks. To discover which stocks were quoted within the first thousand quotes of a given day and how many times, the following query may be issued.

```
[AQuery] – Row Direct Addressing Query
WITH
  OneDay AS
  (SELECT ID, price, timestamp
   FROM   Ticks
        ASSUMING ORDER timestamp
   WHERE  date = '05/11/2003')
SELECT ID, count(*)
FROM   OneDay
      ASSUMING ORDER timestamp
```

```
WHERE ROWID < 1000
GROUP BY ID
```

The WITH query filters out ticks that did not occur in the desired date. The main query's WHERE clause will eliminate all the rows having ROWIDs 1000 or greater, according to timestamp order. Note that this query requires two steps so as to make sure the date filter is executed before the ROWID one. □

## 4.5 Comparing AQuery to Other Order-Aware Languages

AQuery's renditions of order-dependent queries are usually more concise than those of row-oriented languages. To illustrate this point we present the SQL:1999 rendition of the Network Management query presented earlier. Recall that this query's goal is to break sequences of packets (sessions) between pairs of hosts down into "flows" and to calculate statistics of the latter. A flow between a pair of hosts ends – and a new one starts – whenever they stop communicating for a period of 120 seconds or more.

```
[SQL:1999] – Network Management Query
WITH
  Prec (src, dest, length, timestamp, ptime) AS
  (SELECT src, dest, length, timestamp,
         min(ts) OVER
           (PARTITION BY src,dest
            ORDER BY timestamp
            ROWS BETWEEN 1 PRECEDING
            AND 1 PRECEDING)
   FROM   Connections),
  Flow (src, dest, length, timestamp, flag) AS
  (SELECT src, dest, length, timestamp,
         CASE WHEN timestamp-ptime > 120 THEN 1
              ELSE 0
         END
   FROM   Prec),
  FlowID (src, dest, length, timestamp, fID) AS
  (SELECT src, dest, length, timestamp,
         sum(flag) OVER
           (ORDER BY src, dest, timestamp
            ROWS UNBOUNDED PRECEDING)
   FROM   Flow)
SELECT src, dest, avg(length), count(timestamp)
FROM   FlowID
GROUP BY src, dest, fID
```

Separating the flows requires checking the intervals between time-consequent packets for a given pair of hosts. Expressing this calculation in SQL:1999 is not

entirely straightforward. 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 turned true (1) at each packet whose difference to the preceding one exceeds two minutes; otherwise the flag is turned to false (0). Next, the *FlowID* sub-query sums these flags cumulatively, creating an auxiliary flow ID, *fID*. The main query uses these results.

By contrast, the combination of AQuery's column-orientation, underlying data model, and built-in support for order makes it easier to write the same query. For convenience, we repeat the query here.

```
[AQuery] – Network Management Query
SELECT src, dest, avg(length), count(timestamp)
FROM Packets
      ASSUMING ORDER src, dest, timestamp
GROUP BY src, dest, sums(deltas(timestamp) > 120)
```

For all queries presented in this chapter we found the same sort of structural discrepancies between AQuery and SQL:1999 renditions that the Network Management query presents. Ultimately, if a calculation depends on several row values at once, a row-oriented language needs an auxiliary construct to align those values in a row. Nevertheless, several row-oriented languages in the literature provided inspiring insights.

AQuery borrowed from SRQL the early introduction in a query of an order defining clause. AQuery differs from SRQL in that the latter has a row-oriented semantics. Therefore several expressions that are valid in AQuery are not so in SRQL. SRQL cannot handle the table equivalent of non-1NF arrables. We have shown that this feature was useful in order-dependent queries.

In contrast, SEQUIN can handle tables that have sequence-valued fields. But whenever a query involves fields of both the table and the sequence, SQL is used to deal with the former and SEQUIN with the latter. That can lead to somewhat difficult-to-read queries.

AQuery takes its main inspiration from KSQL: namely its arrable notion, and its column-oriented semantics. AQuery differs from KSQL by trying to preserve the SQL flavor to a much greater extent than KSQL, by the introduction of the `ASSUMING ORDER` clause to make the use of order declarative, and (though this is independent of the semantics) by using cost-based optimization.

## 4.6 Conclusion

AQuery was designed to support order-dependent queries without compromising on backwards compatibility to SQL-92 (modulo nested capabilities). AQuery's clauses are the same as SQL's – `ASSUMING ORDER` is an optional clause – and support all expressions that SQL clauses do, even though the former is column-oriented and the latter is row-oriented.

AQuery meets the criteria defined earlier for order-aware languages. It has declarative order, semantically the operations preserve order and the order idioms may claim to be intuitive.

One characteristic makes AQuery particularly amenable to query optimization. It is rather simple to identify which expression in a query and thus which clauses are order-dependent. A simple type analysis can check whether the use of an order-dependent function makes an expression order-dependent. We will see in the next chapter how an optimizer can take advantage of that knowledge.

# Chapter 4

## AQuery Optimization

### 4.1 Introduction

A SQL query specifies how its result must look but it does not establish how to compute it. For instance, if a query's WHERE clause is a conjunction of two predicates ( $p_1$  AND  $p_2$ ), either predicate may be chosen to execute first or execution orders may be mixed. An optimizer decides so in a cost-conscious fashion, usually executing the most selective predicate first or the one having a useful index.

This optimization approach presupposes a query to have alternative query execution plans (QEPs). A first QEP comes naturally from translating the query's text into its algebraic equivalent using the order of operations given by the query. Other QEPs may be obtained by the application of *query transformations*, rearrangements of the operators in a QEP that do not alter the query's semantics. In our previous example, the optimizer could first use a transformation that broke a composite selection into two simple ones. It could then decide on any order, for there is a transformation that could commute this pair of selections [11].

AQuery optimization can be done in this transformational fashion. If a query does not contain any mention of order then it can be optimized exactly as SQL would be. If early ordering is used (ASSUMING ORDER clause) then transformations involving sort can be applied. In the literature, sort transformations were addressed in two distinct but complementary ways.

The transformations in [34] avoid redundant sorting work by either eliminating the sort altogether or by reducing the number of columns over which sort is done. To eliminate a sort over a table  $r$  with respect to column  $A_1$  ( $sort_{A_1}(r)$ ),  $A_1$  must be found to be a prefix of the existing order of  $r$ 's records. This is the case whenever an index (to be precise, an index that orders its key information such as a B-tree) clustered by  $A_1$  is used to scan  $r$ , or whenever a predicate such as ' $A_1 = value$ ' has been previously evaluated over  $r$ . To reduce a sort over  $r$  in respect to columns  $A_1$  and  $A_2$  to  $sort_{A_1}(r)$ , either  $A_1$  must be a key of  $r$  or  $A_1$  must functionally imply  $A_2$ .

**Example 4.1** – Suppose Connections(host, port, client, timestamp) ORDERED BY timestamp is an arrable that stores the clients' addresses that accessed a

network’s services (port, host) and when did they do so. The following query fetches the clients that connected to host ‘atlas’ in timestamp order. (The use of ASSUMING ORDER here is merely illustrative; more realistic queries will follow this introduction.)

```
[AQuery]
SELECT client
FROM   Connections
      ASSUMING ORDER timestamp
WHERE  host = 'atlas'
```

The QEP derived directly from the query text is shown in Figure 4.1(a). We show plans in the usual diagrammatic way but 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 corresponding transformation number. The formal descriptions of the transformations are given in Table 4.1.

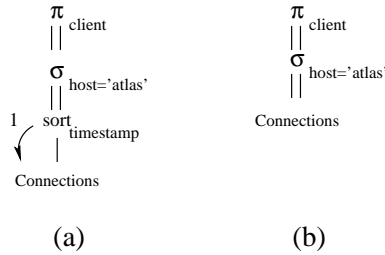


Figure 4.1: An initial QEP and the application of a sort elimination transformation

Note that operators after the sort are connected by double-arcs. This implies that they maintain the order the sort imposed. The sort can be eliminated because it matches the order defined for the arrable Connections, according to transformation 1 in Table 4.1. The resulting plan is shown in Figure 4.1(b).  $\square$

The other way order transformations were addressed in the literature was in a context in which sort and other operations interacted [36]. If relations were lists rather than sets of records and relational operations were carried in an order-preserving fashion, then sort and projection would be commutative. Similarly, sort and selection would be so, and sort could be pushed-down over a join – the complete list appears in [35]. These transformations all apply to AQuery as well.<sup>1</sup>

**Example 4.2** – Take the same query and arrable as in the previous example, but this time let Connections be ORDERED BY host and timestamp. For convenience, the initial syntax-driven QEP is repeated in Figure 4.2(a). Note that

<sup>1</sup>However, [36] and [35] consider that every expression is an order-independent one. The implications of this are discussed in 4.3.

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

Table 4.1: Equivalences between sort and remaining algebra operators

the selection (host = 'atlas') can now benefit from the existing order (host, timestamp). Note also that by evaluating the selection first, fewer records will need to be sorted. The transformation 3 commutes a selection with a sort and when applied here it generates the QEP shown in Figure 4.2(b).

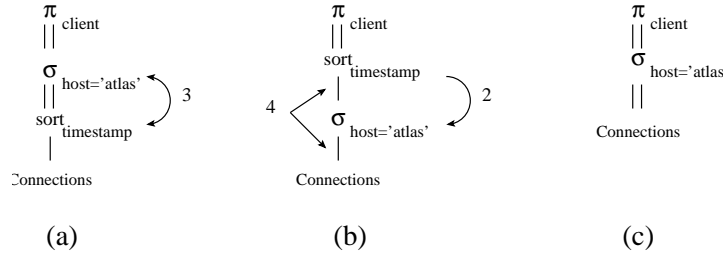


Figure 4.2: An initial QEP and the application of a selection push-down transformation

There is a way to further save the work involved in the sort. At this point, the order of records between the Connections scan and the selection is irrelevant (single arcs), as is the order between this latter and the sort. Changing the order of the records in a portion of the plan in which order is irrelevant does not impact the semantics. Thus the selection can be converted into an order-preserving one by transformation 4 so as to propagate the host, timestamp order. Since the selection output contains only records from host 'atlas' it can be said to be ordered by timestamp. This in turn makes the sort redundant. The latter manipulation is captured by transformation 2 in Table 4.1 and the final plan is

depicted by Figure 4.2(c). □

## 4.2 Optimization of Edge Selections

Optimization of AQuery goes beyond sort elimination and move-around. AQuery’s order idioms are often built around the edge-functions (e.g., `first()`, `last()`), which allow rather aggressive optimizations as well. We introduce these new techniques through examples and evaluate their relative performance improvement.

### 4.2.1 Implicit Selections and Sort-Edge

Let us use the arrable `Connections` once again, this time `ORDERED BY host`. In an intrusion detection scenario an administrator may wish to find the last client that connected to a given server. In AQuery this query would look like the following.

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

The above query’s initial plan is depicted in Figure 4.3(a). A regular selection such as  $\sigma_{host='atlas'}$  can be pushed down over a sort [36]. Transformation 3 in Table 4.1 is a slight variation of that transformation but here order-preservation or lack thereof is made explicit. The advantages of this transformation are the same as example 4.2’s: sort work reduction. But once more, the gains can go further.

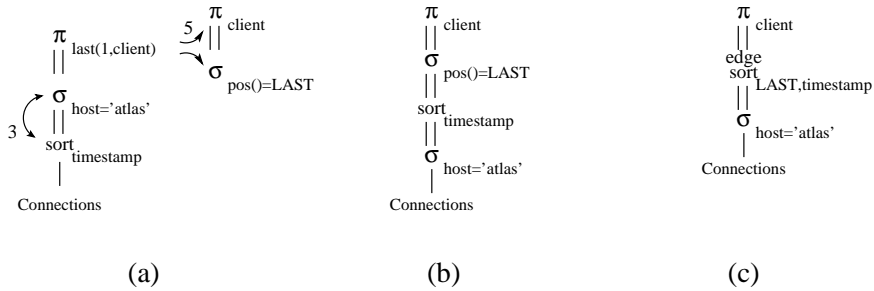


Figure 4.3: Implicit selection and sort-edge optimization

The projection  $\pi_{last(1,client)}$  includes an implicit selection, i.e., it is only interested in one client. This is a particularity of AQuery’s column-oriented semantics – a projection over a function that itself performs a selection. The transformation 5 in Table 4.1 is a new transformation that replaces a projection with indexing by a pure projection plus a selection of the desired positions. 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. If such positions are on an end of the

operand array, we call this selection an *edge selection*. The result of applying this transformation is seen in Figure 4.3(b).<sup>2</sup>

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.

AQuery implements the logical pattern  $\sigma_{edge-condition}^{op}(sort(r))$  through a physical operator called sort-edge. It uses a modified heap-sort to keep the top (or bottom)  $n$  elements, as appropriate. This is similar to the approach used in [8] except that we modify the heap-sort to make it stable.<sup>3</sup>

The gains in performance provided by sort-edge are considerable. In Figure 4.4(a) we compare a sort-edge with a regular sort operation. The graph shows that for small slabs – 1 and 10 elements – sort-edge takes a small fraction of the time needed to sort the entire set.

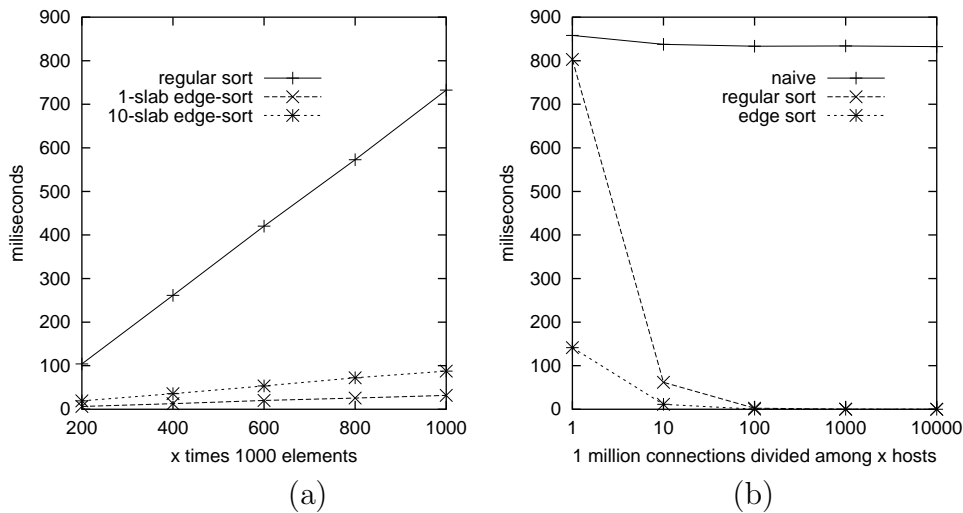


Figure 4.4: Efficiency of sort-edge technique

The benefits of such work reductions reflect in the performance improvement of our example query. Figure 4.4(b) compares the performance of the three plans of shown in Figure 4.3 for a Connections arrable of 1 million rows divided across a varying number of hosts. The “naive” plan performs the entire sort and only then filters for host 'atlas'. The other plans start by executing the latter selection, which can benefit from existing order. Note that in the 1- and 10-distinct hosts scenarios a considerable number of rows have to be sorted and sort-edge delivers a much better performance than regular sort. For the remaining scenarios, very few rows survive the selection and both optimized plans perform equally well.

<sup>2</sup>A word of caution: Had we separated the implicit edge selection before moving the regular selection, they would be adjacent. Regular and order-dependent selections do not commute.

<sup>3</sup>A stable sort is one that does not change the original order of records having identical value on the sorted key. Heap-sort is not naturally stable. It becomes stable if one concatenates a tuple ID to the key.

## 4.2.2 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. The sort splitting technique applies to this kind of scenario. Consider again the arrable `Connections ORDERED BY host`. The following query finds all the clients that connected to the last host to be accessed.

```
[AQuery]
SELECT client
FROM   Connections
      ASSUMING ORDER timestamp
WHERE  host = last(1,host)
```

An initial plan for this query appears in Figure 4.5(a). `Timestamp` is not a prefix of `order(Connections)`, thus the sort over `timestamp` may be required. However, `host` is a prefix of `order(Connections)`, and therefore the selection  $\sigma_{host=last(1,host)}$  may take advantage of it.

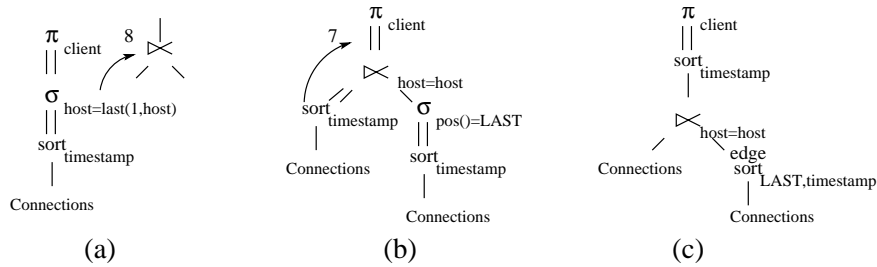


Figure 4.5: Sort-splitting optimization

The sort-splitting technique says that 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 8 in table 4.1. The benefit of the semi-join is that we can now manipulate order on each of the semi-join's arguments independently.

Figure 4.5(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 7 in table 4.1 commutes a semi-join and a sort. It states that under certain conditions sorting a semi-join is equivalent to sorting its left stream and then performing an order-preserving semi-join. The conditions hold here.

This transformation's impact here is two-fold. First, the evaluation of the semi-join predicate is facilitated by an existing order. Second, sorting over `timestamp` has to be done just over records generated by the semi-join. This is much cheaper than the original semi-join over the whole arrable. The resulting plan appears in Figure 4.5(c).

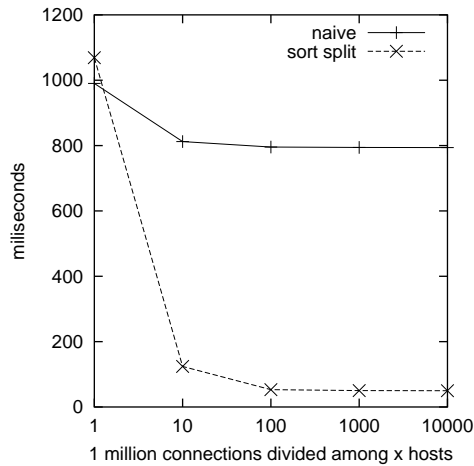


Figure 4.6: Efficiency of sort-splitting technique

Figure 4.6 shows the performance gains of applying the sort splitting technique to the example query. The efficiency of the optimized plan stems from delaying the enforcing 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.

### 4.2.3 Early Edge Selection and Edgeby

Edge selections can often be performed very early in a query. This requires transformations that push edge selections all the way through operations such as joins.

Consider the arrable Ticks(ID, date, price, timestamp) ORDERED BY timestamp, which stores stock quotes, and the arrable Portfolio(ID, name, tradedSince) ORDERED BY ID that stores the subset of securities with which an analyst deals. Name is a unique identifier of securities in Portfolio, and so is ID. An analyst may want to retrieve the last price of a security by its name through the following query.

```
[AQuery]
SELECT last(1, price)
FROM   Ticks, Portfolio
      ASSUMING ORDER timestamp
WHERE  Ticks.ID=Portfolio.ID
      AND name = 'ACME'
```

An initial plan for this query is depicted in Figure 4.7(a). Note that in this plan the selection is carried after the join and the sort. It would be more advantageous to perform it earlier. A regular selection can be commuted with the sort by the application of transformation 3, as was done before. Because the selection would

then be in a portion of the plan that is order-independent, it could then be pushed down over the join using the classic join-selection commutativity [11]. The result is seen in Figure 4.7(b).

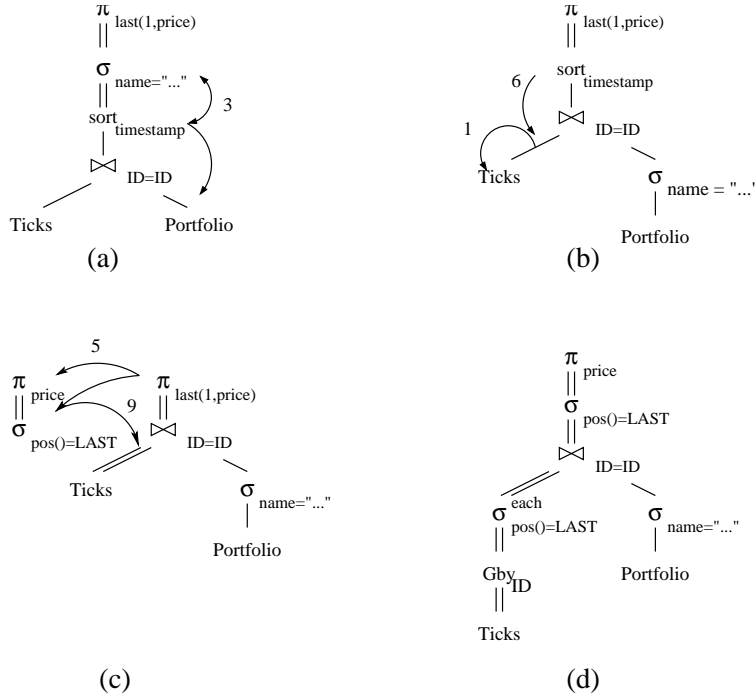


Figure 4.7: Early edgeby optimization

Upon realization that the order of Ticks matches the ASSUMING ORDER of the query, the optimizer would try to eliminate the sorting completely. Transformation 6 in table 4.1 commutes a join with a sort while still keeping track of order. That is a slight variation of a transformation in [36] in which order-preservation is made explicit. As Trades is already ORDERED BY timestamp, that sort may be eliminated, as transformation 1 in Table 4.1 makes possible. The result is seen in Figure 4.7(c).

This query also contains a projection-with-selection, and again they can be broken apart. The consequent presence of the edge selection after the join suggests that it may be unnecessary to perform the join in its entirety. Portfolio.ID is a key and therefore it guarantees that each record in Ticks will match at most one record in Portfolio. (Foreign key joins are among the most frequent of equijoins.) Under these conditions we could push down this edge selection in the following way: For each ID in Trades, find its last record by grouping Trades by ID and selecting each last record. By applying the edge selection earlier, the query's join examines far fewer rows than before. The final selection would then pick the desired price. This is what transformation 9 in table 4.1 does. The final plan is shown in Figure 4.7(d).

Replacing an edge-selection by a grouping operation and the very same edge-selection may look more expensive, but it isn't. 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  $\sigma_{edge-condition}^{each}(Gby(r))$ .

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 violate the edge condition. Depending on this condition, edgeby can scan arrables backwards or forwards. In fact, edgeby can be parameterized to perform grouping followed by any possible edge selection.

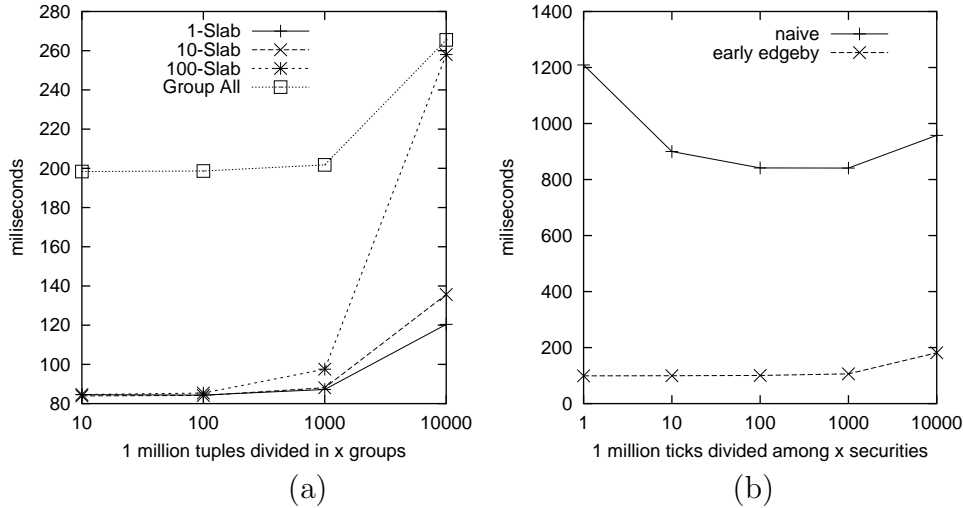


Figure 4.8: Efficiency of the early edgeby technique

In most cases, an edgeby requires a small fraction of the time required to perform the associated group-by, if done in its entirety as shown in Figure 4.8(a). We used the arrable Ticks 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 give better performance. A degenerate case is seen where a 100-slab is taken from groups that are themselves 100 records wide. Edgeby doesn't improve performance here – but doesn't hurt either.

The result of applying the early edgeby technique is shown in Figure 4.8(b). The naive and the optimized plans for the example query can be seen. By applying a 1-slab 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.

#### 4.2.4 Sort Embedding

If a GROUP BY operation follows a sort it may be advantageous to invert their order. Performing a single sort over the entire data is often more expensive than performing several sorts, each embedded within a group.

Consider again the arrable Ticks, this time with no determined ORDERED

BY. (Often ticks arrive in a “near-timestamp” order.) The following query fetches the ten most recent prices for each security ID

```
[AQuery]
SELECT ID, last(10,price)
FROM Ticks
      ASSUMING ORDER ID, timestamp
GROUP BY ID
```

An initial plan for this query is shown in Figure 4.9(a). We can separate the implicit selection from the projection as we did before. The resulting plan appears in Figure 4.9(b).

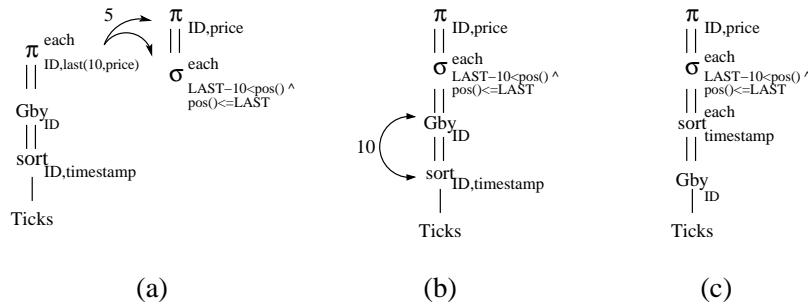


Figure 4.9: Sort-embedding optimization

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. Moreover, for this particular query the smaller sorts would then be followed by edge selections – sort-edge would apply. The transformation 10 in table 4.1 allows commuting a sort with a group-by. Note that (a) group by must deliver its results in the same order it is grouping by over (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 4.9(c). Note how a double-arc connects group-by and sort-each, because this instance of group by is order generating.

Figure 4.10(a) 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. When only one group exists, there’s no point in applying the technique – but, again, there’s no penalty in doing so. Replacing one big sort by several smaller ones starts to payoff whenever more than 10 groups exist.

The efficiency of sort-embedding reflects in the performance of the example query. Figure 4.10(b) shows the comparative performance of plans for the naive and optimized cases. 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.

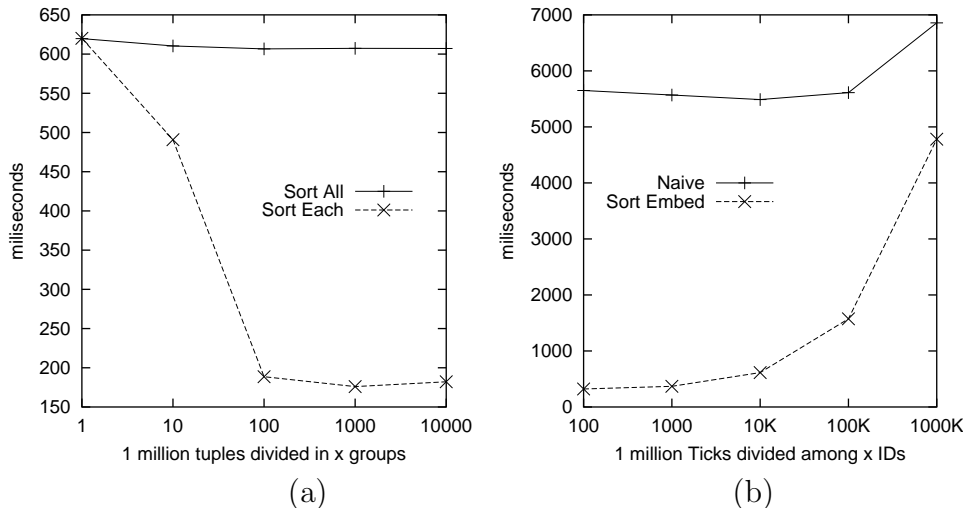


Figure 4.10: Efficiency of the sort embedding technique

### 4.3 Related Work

Order management is not normally fully integrated in a query’s optimization process. Sort has traditionally been seen as a physical property to be included in a plan (if not specified by SQL’s ORDER BY clause) only to support an efficient algorithm such as merge join. Mechanisms such as Starburst’s “glue” [23] or Volcano’s “enforcer” [13] made sure a sort step was added whenever an efficient algorithm required it. This approach to order management was shown to miss several optimization opportunities [34].

The authors in [34] suggested and implemented an order-management step in their optimization process. It greatly improved QEPs for queries that have order requirements due to the clauses ORDER BY, GROUP BY, or the DISTINCT modifier. The transformations 1 and 2 in table 4.1 come from their work. Yet this order management step did not consider interactions between sort and other operations – an important condition to perform order-management as an integrated aspect of query optimization.

To the best of our knowledge the authors in [36] were the first to promote sort to a logical operator and to suggest that order optimization could (i) be transformational and (ii) be considered at the same time as other transformations. AQuery follows this same idea.

While using several of the transformations in [36] – namely transformations 3, 6, and 7 in table 4.1 – AQuery has contributed a few of its own – transformations 5, 8, 9, and 10. The transformations suggested here go a step further by taking into account the fact that expressions in a query can be order-dependent. (In [36] expressions are all order-independent.) This allowed, for instance, to identify that two selections over order-dependent expressions do not commute, but they would had the expressions been otherwise. AQuery also differs from that work by treating edge selections as a case of order-dependent selections rather than as a new operator (called operator *top-k* in [36]).

In fact, edge selections were first described in [8]. A clause STOP AFTER

was suggested that was capable of limiting the final cardinality of queries whose results were ordered (by an ORDER BY). The sort-edge operator presented here was motivated by a similar operator described there. There are two differences that distinguishes AQuery from that work, though. First, in [8] the optimization process for STOP AFTER was considered in a phase prior to – rather than integrated with – the phase in which plan enumeration takes place. Therefore interaction between sort and other operators are not fully considered. Second, AQuery language allows edge selections to occur at any point of the query and not only at its end. We have shown, for instance, queries that applied edge selections within groups (e.g., last 10 quotes of each security in a Portfolio). We have also shown that new techniques such as sort embedding can be applied to optimize these queries.

A particularly inspiring integrated optimization technique appeared in [22]. The AQL optimizer can manipulate operators (or newly added functions) on the calculus level, i.e., by application of variations of  $\lambda$ -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. A complete fusion of these ideas requires more exploration.

# Chapter 5

## System Design and Implementation

### 5.1 A Column-Oriented Execution Model

The AQuery system is a database management system that implements the arrable-based data model and that supports AQuery as the query interface. In developing a new database system, we wanted to investigate how the use of arrays as the basic data structure could work for performance. Our studies started with how to execute a query, once a query execution plan is obtained.

To carry out an execution plan a system has to run each of the plan's operators. There are some alternatives to do so. The most common way is to equip each operator with a *get-next-row* interface (`getNext()`) that returns one row of the operator's result at a time. This execution model is called *iterator-based* [12]. To obtain a plan's first resulting row the system calls the root operator's `getNext()` interface. In turn, the root operator calls the `getNext()` of the operator from which it consumes rows. This cascading effect propagates down the plan until a leaf operator reads an input's row. The row is passed back to the caller and it gets processed on its way up. Assuming it belongs to the query's result, it eventually arrives back at the root operator. The system starts the process over again until an end-of-rows is signaled.

The iterator model has several advantages, but its efficiency is poor especially because of poor cache. Cache inefficiency is due to the loading of unnecessary columns' data to perform an operation [1]. For instance, if a selection's predicate involved only one column, it would not need an entire row of values to evaluate the predicate. The unused fields occupy precious cache space.

In addition, the iterator model has an inherent overhead of a function call per row per operator. In a simple-predicate selection, a call's cost (parameter stacking and context saving) can be of the same order as the evaluation of the predicate itself.

A contrasting execution model is the *column-oriented* one found in the Monet database system [3]. Instead of handling a row at a time, the model assumes data is vertically partitioned in columns and manipulates these columns as units. For instance, a selection would look at the column involved in the predicate, and

only at this, and would return the row IDs and that column's data that satisfy the predicate. This model presents better performance in modern architecture hardware – hierarchical memory and super-scalar CPU – than the iterator one [25].

The AQuery system uses a column-oriented model as well, but slightly more generally than Monet. The model maps easily to arrables which are naturally partitioned in arrays (columns). To perform the above selection, the AQuery system might materialize its results as Monet does. But this would involve performing memory copies of data. Alternatively, the AQuery system passes on the next operator the reference to the data it operated upon along with a collection of indexes to the still relevant rows. This avoids unnecessary materialization. We use the term *effective index array* to denote this collection of indexes resulting from one operator that may be passed along to subsequent operators. A running example of a plan demonstrates how this concept works.

Consider the arrable Trades, depicted in Figure 5.1, which stores prices at which stocks exchanged hands, and at which quantities (volume) and timestamps these transactions took place. The arrable Base, shown in the same figure, classifies stocks according to the type of business they conduct, using a Standard Industry Code (SIC). For visualization purposes we show the indexes of each arrable's rows in the left-hand column (small fonts under the arrables name). Consider further the following query that fetches the last quote of each stock in the COMPUTER industry.

```
[AQuery]
SELECT t.ID, last(price)
FROM   Trades t, Base b
      ASSUMING ORDER t.ID, ts
WHERE  t.ID = b.ID AND
      SIC = "COMP"
GROUP BY t.ID
```

The first step of a plan for this query is to initialize an effective index array. For clarity, let us assume temporarily that the arrables fit entirely in memory. The initial index array would contain all the existing indexes for both arrables, that is, 0..6 for arrable Trades and 0..2 for Base. This structure is depicted in Figure 5.2(a).

Following the query's text, the next operation in the plan would be the join of the two arrables. Because the predicate involves the columns Trades.ID and Base.ID, the join operation needs to access only these columns. The effective index array depicted in Figure 5.2(b) shows how this operation recorded which pair of positions of Trades and Base satisfied the join predicate.

Although irrelevant at this point of the query, the index array imposes an ordering for the rows of the intermediate result so far. Coincidentally, the first row of the intermediate result is a join of the first rows of the arrables, Trades[0] and Base[0]. The following row is the join of the third row of Trades, Trades[2], and the first row of Base, Base[0] – and so on. Note that the indexes used here are positions relative to the arrables. Data need not be stored in ascending index

Trades	ID	price	ts
0	ACME	12.05	1
1	XYWZ	42.35	2
2	ACME	12.04	3
3	EMCA	17.19	4
4	EMCA	17.20	5
5	ACME	12.02	6
6	XYWZ	42.37	7

Base	ID	SIC
0	ACME	COMP
1	XYWZ	AUTO
2	EMCA	COMP

Figure 5.1: Two example arrables and their rows indexes

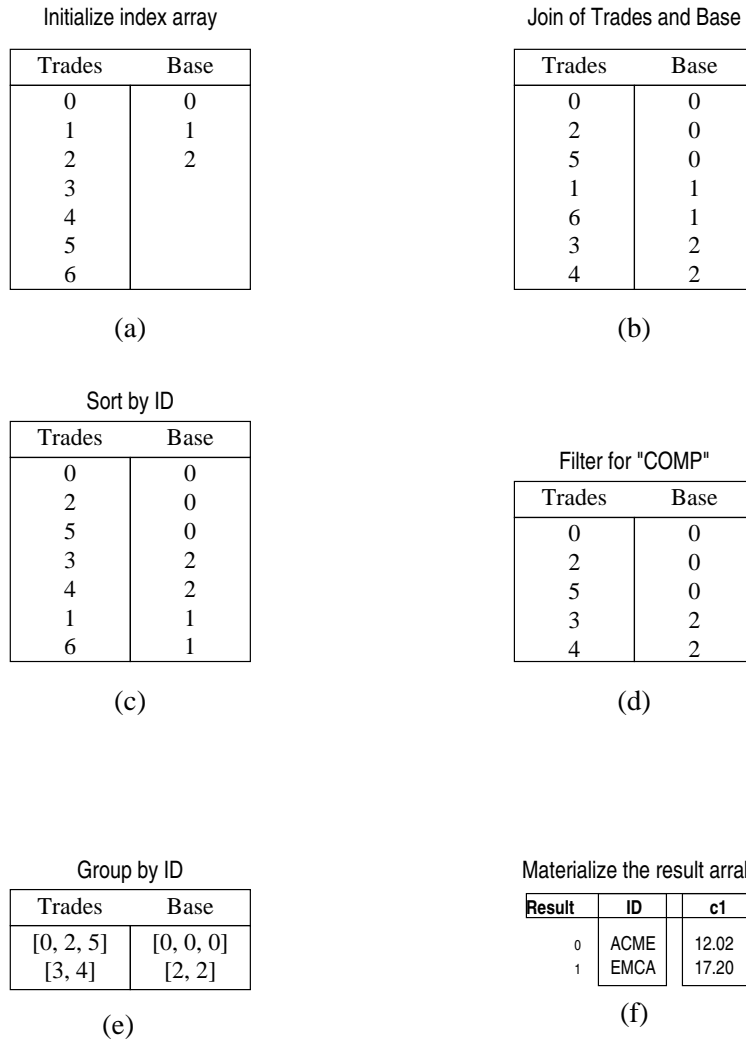


Figure 5.2: Effective index array during a query execution

order (storage physical independence). What is important here is the ability to fetch a row by its position, a natural approach for arrays.

The ordering of rows in the current index array is not the one specified by the query's ASSUMING ORDER clause and thus a next step in a plan could be a sort. Recall that the ASSUMING ORDER is over Trades.ID and Trades.timestamp. Here again only those columns need to be accessed to compute the new index array. The sort would rearrange the indexes in a way to impose an ID and timestamp order. The Figure 5.2(c) shows the result of the operation.

Note that the sort took as input the result of the join operation. The permutation occurred on pairs of indexes of both arrables.

A selection may follow that eliminates the indexes that do not correspond to computer industry's stocks. The only column needed here is Base.SIC. The selection loads it and dereferences using the effective indexes that correspond to the Base arrable. Note that the Base.SIC column has three values in the original arrable. But here, the resulting data vector has six positions (0, 0, 0, 2, 2, 1), for a join and a sort have already been applied. The selection thus generates the effective index array seen in Figure 5.2(d).

Note that the selection eliminated the indexes that contained 1 in the Base portion of the index array.

The next operation is the GROUP BY. The first step here is to index the column Trades.ID by the portion of the index array corresponding to the arrable Trades. Then, the operation captures the notion of groups by putting the indexes of rows that belong together in a same array-valued field of the effective index array. Figure 5.2(e) shows the effect of GROUP BY.

The effective index array shows that the result of calculating the query up until the GROUP by has two rows. Moreover, each of the columns that was not grouped-by will have array-values fields. For instance, if one wished to fetch Trades.prices for the first group, one would load that column and use values at positions 0, 2, and 5. For the second group, positions 3 and 4 would be used instead.

Finally, the projection in the SELECT clause may be applied. A final projection means materialization of the results. Here it takes into consideration that ID was grouped-by. Only the first index of each group is used in the result (i.e., ID[0] and ID[3]). By contrast, the function last() is called once for each group and is passed the corresponding price arrays for each group (price[0 2 5] for group 0, and price[3 4] for 1). It returns the values of price[5] and price[4] respectively. The resulting arrays are then assembled into an arrable, show in Figure 5.2(f).

The column-oriented execution model does not prevent a plan from using a constrained amount of memory. To understand why, we introduce the notion of *pipelining* among operators. A *non-blocking* operator (pipeline-friendly) is one that can decide on an answer without checking all the input rows. It therefore does not stop the flow of rows until it reads all the input. For instance, selection is a non-blocking operator. Sort is blocking.

To limit the memory a query plan uses, it suffices to adapt the execution of each operator type to the limitation. A non-blocking operator would chop the input columns to appropriately sized slices and would process one slice at a time.

A blocking operator can use implementations based on external algorithms (e.g., sort on disk).

In summary, the execution model stores tuples of indexes instead of tuples of values. It dereferences those indexes on demand.

## 5.2 Implementing the Execution Model

The AQuery system implements the column-oriented execution model using an interpreted array-oriented language called K [20]. K is a modern descendant of APL [18], the reference language on array manipulation. In these languages, any data structure can be represented by arbitrarily shaped arrays, which can be manipulated by the composition of very simple but highly efficient array primitives. These primitives operate on the array level, i.e., they are mappings from an array (or a list thereof) to an array. Because AQuery was built on top of K, it was natural to represent AQuery's arrables and operators.

A K implementation of an arrable follows the latter's definition: a collection of named arrays. On disk, an arrable is mapped to its own sub-directory. Each component array is stored in a file named after the column it represents. In memory, an arrable is represented by a dictionary structure. Entries are named after the column they represent and each entry stores a column's array. Whenever a column needs to be accessed, the system loads the corresponding array-file from disk to virtual memory and adds an entry to the corresponding arrable's dictionary. When the array (column) is no longer needed the system eliminates the dictionary entry, freeing memory.

To be executed, a plan's operators are translated into a sequence of K instructions that follow a three-step template. In step one, the operator receives both the column handles it needs (e.g., the columns involved in a selection's predicate) and the current effective index array of the query. In step two, the operator accesses data (uses the index array to retrieve the still relevant columns' values) and performs its calculation (evaluates the predicate over the values). In step three, it returns the updated index array (eliminates the rows that mapped to false).

In this context, the system translates a plan into a K function that puts together a sequence of such operators. The function takes handles to the arrables it manipulates and outputs an arrable. The Figure 5.3 shows the code generated for the following query.<sup>1</sup>

```
[AQuery]
SELECT last(price)
FROM   Trades
       ASSUMING ORDER ID,ts
WHERE  ID= 'XYZW'
```

A query plan is an anonymous function, defined between the curly braces of lines 1 and 8. Immediately after its definition, the system calls the function passing

---

<sup>1</sup>Operations on K are denoted by symbols which makes the code extremely compact. Critics and even the language designer agree that K code looks like line noise to the uninitiated.

a list of arrables (line 8). In this case, the function takes only a singleton list, a handle to the Trades arrable. The query plan starts by initializing its effective index array in line 2, with as many index vectors as there are input arrables. The code here is generic, in that it would have worked correctly if any other number of arrables were passed. The next step in that plan is to load data. Line 3 defines which column handles are going to be used; the variable 'cols' is thus initialized with handles for three columns. In line 4 the data for these columns are fetched from disk, resolving their address. A dictionary named 'r' is created to hold the Trades arrable's data in memory. Next, in line 5, the columns 'r.ID' and 'r.ts' are sorted in ascending order. That is what the K operators '<+' (flip and sort) are doing in that line. Note that all access to these variables is now indexed by the corresponding portion of the effective index array, eia[0]. As a result of the sort the index array is rearranged, as in 'eia:eia@....' Line 6 computes which positions of the contain "XYWZ" in their r.ID column. Again, the effective index array is rearranged accordingly. Finally, line 7 calls the in-line function last 'x[-1+#x]' and passes to it the price column at position that are still pertinent. Thus, only the prices of stock 'XYWZ' ordered by ID and timestamp are passed. The line 7 then creates the resulting arrable which contains a column called 'price' and the result of the last function.

```

1.  {[t]                                     / declare parameter
2.    eia:(#t)#_n                           / initialize eia
3.    cols:'ID'price'ts                     / list of needed columns
4.    r:@[_n;cols;::{1:($t[0]),"/", $x}'cols / load data into dic r
5.    eia:eia@\:<+( r.ID[eia[0]];r.ts[eia[0]] ) / sort by ID,ts
6.    eia:eia@\:&r.ID[eia[0]]~'XYWZ'         / selection over ID
7.    :.,('price; {x[-1+#x]}[r.price[eia[0]]]) / compute last
8.  }[, 'trades]
```

Figure 5.3: Example of a K plan

## 5.3 From Text to Execution: the entire flow

Having discussed the execution engine's strategy, we are now ready to discuss the entire process from text to execution.

### 5.3.1 Parsing

The text first undergoes lexical and syntactical analysis, commonly referred to as "parsing." These phases make sure a query has valid syntax and generates an abstract syntax tree. The AQuery system parses a query using a LL(1) grammar and a recursive decent parser and generates tree that represents the query's text. This format is more amenable to the semantics step.

### 5.3.2 Semantics Step

Semantics checks in AQuery do more than simply verify that all the objects in a query – columns, arrables, functions – exist. AQuery semantics requires expressions to obey shape constraints, as discussed in Chapter 4. For instance, a search condition in the WHERE clause is a boolean expression whose cardinality is equal to that of the arrable resulting from the FROM clause. Therefore, each use of a non-size-preserving function in that clause must pass a semantics check.

To illustrate how shape verification is done, consider the expression 'col + first(10,col).' Addition is valid between two equally sized arrays, between two scalars, or between an array and a scalar. Therefore the above expression is invalid.

The AQuery system captures these shape constraints in tables such as Table 5.1 for the addition operation. Initial cardinality refers to the number of elements a column first had in the context being considered. For instance, in the FROM clause initial cardinality has as many elements as the table, in the WHERE as many element as the Cartesian product of the FROM clause generated, and so on. This notion is important because some clauses require initial cardinality (WHERE, GROUP BY, HAVING) while others don't (SELECT).

+	n	m	1	initial
n	n	error	n	error
m	error	m	m	error
1	n	m	1	initial
initial	error	error	initial	initial

Table 5.1: Cardinality of the addition operation

### 5.3.3 Relational Optimization Support

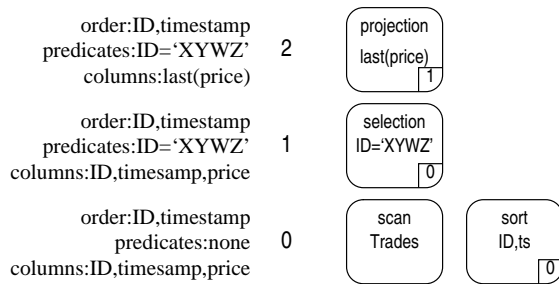
Once semantics correctness is verified, a query is translated into its algebraic format. The result is a tree in which nodes are operators from the arrable algebra described in chapter 4. This tree is a possible plan for the query, but only rarely a good one. Algebra equivalences (query transformations) may generate an equivalent plan with improved performance. These transformations are described in the relational literature [11] and new ones regarding order equivalence were introduced in chapter 4. Applying a transformation to a plan entails generating a new plan, which may share several nodes with the original one. Given the number of equivalent plans that even a small set of transformations can generate, support for plan manipulation is needed.

The AQuery system uses a data structure called MEMO [13] to store a forest of query plans in an efficient way. The main idea behind the MEMO is to group tree nodes into equivalence classes. Two nodes will belong to the same class if the result of processing the queries to which they belong up and including that node is semantically equivalent. A group is characterized by properties of the data their nodes generate. Examples of properties are: tables and columns accessed thus far,

predicates applied, etc. In the AQuery system data ordering is a discriminating property for groups.

Figure 5.4 shows two representations of the MEMO structure for the query that fetches the last quote of the ‘XYWZ’ stock. Operations are represented as squares. An operation that takes the result of another one as input has the group of the latter indicated in the lower right corner. Groups of operations are numbered and a subset of the properties of each group is shown on the left-hand side of their MEMO. Recall that the query in question performs a scan on Trades to retrieve the ID, price, and timestamp columns; a sort over ID and timestamp; a filtering (selection) over ID; and, finally, a materialization (projection) of the application of last() over the resulting price vector. The configuration of the MEMO will depend on the arrable’s order.

If the Trades arrable were ORDERED BY ID and timestamp, then the MEMO in figure 5.4(a) would result. The scan of Trades appears as the first operation in group 0. (We refer to it as operation 0.0.) The query also performs a sort, but because the arrable is already in a convenient order, the sort is redundant. Therefore, the sort over the result of the scan would not change the data’s properties and thus sort is placed in the same group as the scan. (It becomes operation 0.1.) Note that the nodes that point to group 0 may now choose between two alternate sub-plans. The subsequent operations change properties of the data and as such they are inserted into groups of their own. This MEMO represents two possible plans:  $\langle 2.0 \leftarrow 1.0 \leftarrow 0.0 \rangle$  or  $\langle 2.0 \leftarrow 1.0 \leftarrow 0.1 \leftarrow 0.0 \rangle$ .

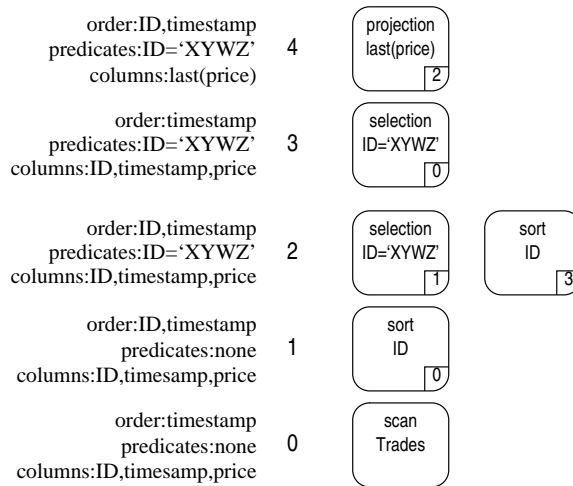


(a)

If Trades were ORDERED BY timestamp alone, The MEMO depicted in Figure 5.4(b) would result. Sorting on ID and timestamp now generates a distinct result than that of the scanning of Trades. Therefore, the scan and that sort belong to different groups (0 and 1, respectively).

This query may benefit from the transformation that commutes a sort with a selection. That is, instead of sorting the outcome of group 0 and then filtering, it does the reverse. A new filter node that takes group 0 as input is inserted in the MEMO in a group 3, for no other group in the MEMO shares its properties. The node that sorts the output of the new filter is inserted in group 2 for it shares its properties. That MEMO now encodes two possible plans:  $\langle 4.0 \leftarrow 2.0 \leftarrow 1.0 \leftarrow 0.0 \rangle$  and  $\langle 4.0 \leftarrow 2.1 \leftarrow 3.0 \leftarrow 0.0 \rangle$ .

At the time of this writing the system was still not applying query transformation automatically. An accurate cost-model is still missing. The system either



(b)

Figure 5.4: Two possible MEMO configurations for a given query

generates a syntax-directed plan automatically or it generates plans in an interactive fashion where the system finds possible transformations to be applied and prompts a user to pick one until the user is satisfied with the plan.

### 5.3.4 ‘K’-Code Generation

After a plan is chosen, it is translated into K language instructions. The translation is made one node at a time. Each type of node requires a different translation strategy. For instance, the code generated for a selection should evaluate its expression and filter the current index array accordingly. By contrast, the code generated for a sort should trigger a reordering of the index array. Moreover, nodes are sensitive to the context. The code generated for a projection changes depending on whether there is grouping or not. The latter’s expressions have to be “eached” whereas the former’s need not. The code generation module handles all these nuances.

The code generation module performs optimizations of its own. The most salient one is a limited but quite effective form of common sub-expression elimination in which column references appear more than once in a same expression. For instance, take the expression ‘price - prev(price).’ Accessing the still active positions of price twice is unnecessary. Upon detection of repeated column-references, it creates temporaries to store that result.

## 5.4 Conclusion

Performance is a major goal of the design of the AQuery system. Each technique in our framework plays a role in performance:

- Processing a query by maintaining its effective index array prevents unnecessary copies of data to be passed among operators. A selection, for instance, doesn’t

need to copy the actual selected tuples, but only to forward their positions to the next operation. A sort does not need to produce a new ordered copy of the data but only a permutation of positions. Ultimately, processing indexes as opposed to copies of data allows the system to defer materialization of results up until the end of the query.

- Column-wise representation of data avoids unnecessary disk accesses. It ensures that in every transfer only data that is actually going to be used – as opposed to entire rows with column data unused by a query – is fetched. This also benefits cache utilization, because irrelevant fields don't clutter cache lines.
- Evaluating expressions in a vector-oriented fashion promotes high reference locality. Vectors are contiguous in memory and most operations involve a sequential scan of a vector, often a contiguous one. Vector-orientation also allows calls to processing functions to be made only once for the entire vector as opposed to once for each element of the vector.

The AQuery system is an ongoing effort. We are currently investigating a query enumerating algorithm [15] that along with a suitable cost model would be able to apply the transformations automatically. We are also considering more sophisticated disk sub-systems such as the one described in [39].

# Chapter 6

## Performance Analysis

### 6.1 Introduction

This chapter addresses the efficiency of AQuery in performing order-dependent queries. It does so by comparing the performance of the AQuery system against that of a commercial SQL:1999 one.<sup>1</sup> SQL:1999 [17] was chosen because its order features are implemented in at least two major DBMS products, and its syntax and semantics are thoroughly documented.

We have resorted to somewhat advanced SQL:1999 features to write the order-dependent queries we test here. To make the chapter more accessible to unfamiliar readers, we present a quick introduction to the “window” feature in the remainder of this section. The more advanced reader may want skip to Section 6.2.

One way to introduce a SQL:1999’s window is as a function that maps each table’s rows to a subset of this table. We will use the table T1 in figure 6.1 to show two running examples.

T1	c1	c2	c3
	a	2	10
	a	1	20
	b	2	30
	a	3	40
	b	1	50

Figure 6.1: A table to be used on a window definition

In our first example, let’s assume one wanted to identify for each given row all the rows that preceded it. Let us assume an ordering based on column c2. To make the example more interesting, we impose a partitioning on column c1, that is, column c2 is ordered within partitions defined by column c1. To specify such a window in SQL:1999 one would write `WINDOW w PARTITION BY c1 ORDER BY c2 ROWS UNBOUNDED PRECEDING`. The last part of the window declaration says that the width of a window spans from the first row (of the partition) up until the row being considered. Figure 6.2 shows the calculation of such a window over

---

<sup>1</sup>Due to license agreement restrictions, we omit the name of the product.

table T1. In the figure, the window instance that corresponds to row  $i$ , according to the ordering given, is denoted  $w_i$ .

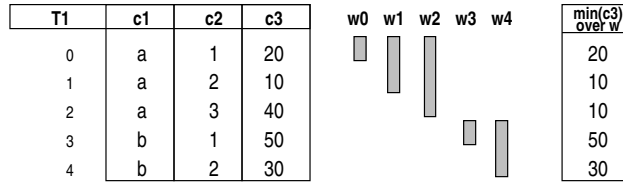


Figure 6.2: A running minimum window

Our example window is useful, for instance, to compute a running minimum. It suffices to apply the aggregate function  $\min()$  to each of the window instances, as the Figure 6.2 shows. In a query, this calculation could be used as follows.<sup>2</sup>

```
[SQL:1999]
SELECT *, min(c3) OVER ( PARTITION BY c1
                        ORDER BY c2
                        ROWS UNBOUNDED PRECEDING)
FROM T1
```

A window may be used in the SELECT clause provided that the query applies an aggregate function over its instances.

In our second example, suppose one wanted to fetch just the previous row for each of T1's row rather than all the preceding rows. Let us use the same partitioning and ordering criteria as before. The difference is thus on the window width. To use a single-row wide window that contains a previous row, one would say in SQL:1999 `WINDOW z PARTITION BY c1 ORDER BY c2 ROWS BETWEEN 1 PRECEDING AND 1 PRECEDING`. The Figure 6.3 shows the calculation of such a window over the table T1.

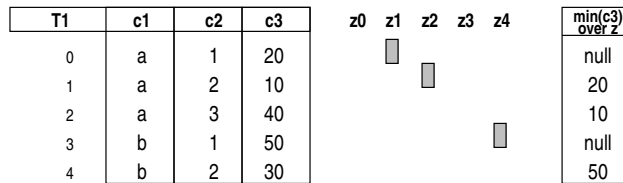


Figure 6.3: A previous row window

This window is useful, for instance, to obtain a column's previous value, the equivalent of  $\text{prev}()$  in AQuery. Note that even though the window instances contain only one row, a query must still specify an aggregating function to use the window. In this case  $\min()$ ,  $\max()$ , and  $\text{avg}()$  would all have the same effect. Suppose one wants the previous element of each  $c3$ 's value, the implementing SQL:1999 query is shown below.

<sup>2</sup>We are using the inline definition of the window here. In the full definition case, the WINDOW clause would be added at the end of the query and the expression would refer to the window name as in ' $\min(c3)$  OVER  $w$ .' As of this writing only the inline use of windows was supported, though.

[SQL:1999]

```
SELECT *, min(c3) OVER ( PARTITION BY c1
                        ORDER BY c2
                        ROWS BETWEEN 1 PRECEDING
                        AND 1 PRECEDING)
FROM T1
```

## 6.2 The Best Profit Query

The Best Profit query finds, for a given stock and a given date, the best profit one could obtain by buying it and then selling it later that same day. This query has been discussed previously in section 3.2 but for convenience, we again present the relevant aspects.

The query uses the arrable/table Ticks(ID, date, price, timestamp), where ID is the ticker of a traded security, timestamp identifies the date and time of a particular trade, date is a human-readable form of the day portion, and price is the price of the security. The table version of Ticks had indexes on timestamp (clustering), on ID (non-clustering), and on date (non-clustering). The arrable version was ORDERED BY timestamp.

The AQuery and SQL:1999 renditions of this query are the following.

[AQuery]

```
SELECT max(price - mins(price))
FROM Ticks
      ASSUMING ORDER timestamp
WHERE ID = 'ACME' AND date = '05/11/2003'
```

[SQL:1999]

```
SELECT max(running_diff)
FROM (SELECT ID, date,
            price - min(price) OVER ( PARTITION BY ID,date
                                      ORDER BY timestamp
                                      ROWS UNBOUNDED PRECEDING)
            AS running_diff,
      FROM Ticks ) AS t1
WHERE ID = 'ACME' AND date = '05/11/2003'
```

The structure of the two renditions differs because in SQL:1999 an aggregate function (max()) cannot take a running aggregate (min(price) OVER ...) as an argument. Nevertheless, both renditions sort first (ASSUMING/ORDER BY timestamp) and then filter (ID='ACME' AND date='05/11/2003'). Note that an optimal plan would do the inverse – sorting is more expensive than filtering.

The plans generated for the two queries are shown in Figurefig:6plan1. The AQuery plan exploits the possibility of performing an early selection. Sort and selection are adjacent and can be commuted with a simple transformation (Table 4.1 transformation 3). Since there were indexes available to evaluate each of the selection's conjuncts, AQuery could use them both. Sort was done only on the reduced set of rows.

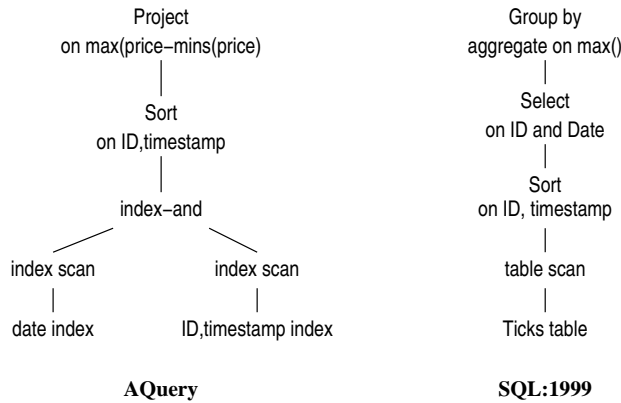


Figure 6.4: Plans for the best-profit query

By contrast, in the SQL:1999 query the sort and the selection are separated by a window operation. The SQL:1999 optimizer doesn't push the selection down. It could in this case since the selection doesn't impact the way groups are formed in the OVER clause; it just eliminates entire groups.<sup>3</sup> It remains though that pushing a selection over a projection that contains windowed functions (SELECT ... price - min(price) OVER ...) requires considerable analysis. The SQL:1999 system we have analyzed chose none of the available indexes.

The difference in plans is noticed in the response times. The chart in Figure 6.5 shows the relative improvement of the AQuery's plan over the SQL:1999 optimizer's. We used Ticks arrables/tables with varying number of securities from 200 to 1000, and using 1000/ticks per security. The chart shows that AQuery results were between eight and twenty one times faster for the best-profit query.

This and the remaining experiments were conducted on a Pentium III-M 1.13Mhz with 1Gb of memory running Linux. The timings reported correspond to wall clock timings. We were careful to allocate the same amount of memory for both optimizers and made sure execution used cold buffers (empty).

The difference grows with the size of the input because the more securities used, the more the SQL:1999 plan sorts rows that will eventually be discarded.

### 6.3 Network Management Query

The Network Management query's goal is to break sequences of packets (sessions) between pairs of hosts down into "flows" and to compute statistics of the latter. A flow between a pair of hosts ends (and a new one starts) whenever they stop communicating for a period of 120 seconds or more. The schema involved in the query is Packets(pID, src, dest, length, timestamp), where pID identifies a packet exchanged between a source (src) and a destination (dest) host, length refers to the size of the packet, and timestamp to the moment this packet was exchanged. The table version of Packets had indexes over timestamp (clustered) and a composed

<sup>3</sup>If vendors use this thesis to improve their SQL:1999 systems, we will consider that a mark of success.

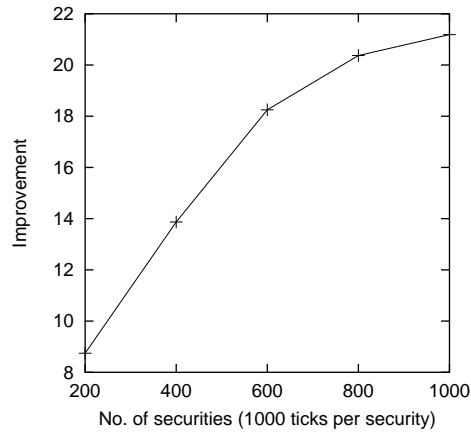


Figure 6.5: Best profit query relative improvement

one over source, destination, and timestamp (non-clustered). The arrable version was ORDERED BY timestamp. The AQuery rendition of this query was discussed in Example 3.3, but we repeat it here, along with its SQL:1999 counterpart.

**[AQuery]**

```
SELECT src, dest, avg(length), count(timestamp)
FROM Packets
    ASSUMING ORDER src, dest, timestamp
GROUP BY src, dest, sums(deltas(timestamp) > 120)
```

**[SQL:1999]**

```
WITH
  Prec (src, dest, length, timestamp, ptime) AS
  (SELECT src, dest, length, timestamp,
    min(timestamp) OVER
      (PARTITION BY src, dest
        ORDER BY timestamp
        ROWS BETWEEN 1 PRECEDING
        AND 1 PRECEDING)
  FROM Packets),
  Flow (src, dest, length, timestamp, flag) AS
  (SELECT src, dest, length, timestamp,
    CASE WHEN timestamp-pptime > 120 THEN 1
    ELSE 0
  END
  FROM Prec),
  FlowID (src, dest, length, timestamp, fID) AS
  (SELECT src, dest, length, timestamp,
    sum(flag) OVER
      (ORDER BY src, dest, timestamp
        ROWS UNBOUNDED PRECEDING)
  FROM Flow)
SELECT src, dest, avg(length), count(timestamp)
FROM FlowID
```

GROUP BY src, dest, fID

Expressing this calculation in SQL:1999 was not as straightforward as in AQuery, although the execution flow of both queries are quite similar – at least semantically. In SQL:1999, the first sub-query, Prec, creates a new column, ptime, containing the previous packet’s timestamp within each source and destination partition. It uses a window that partitions by source and destination, sorts by timestamp, and uses a one-row width window instance. In AQuery, there is no partitioning but sort is done over source, destination, and timestamp. The Prec sub-query has the same effect of a prev(timestamp) in AQuery. (Recall that deltas(col) is equivalent to col - prev(col).)

Next, the SQL:1999 Flow sub-query adds a flag column that is turned true (1) at each packet whose difference to the preceding one exceeds two minutes; otherwise the flag is turned to false (0). In the AQuery rendition, this is what the expression 'deltas(timestamp) > 120' does.

The SQL:1999 query continues by calculating FlowID, which sums the flags cumulatively, creating an auxiliary flow ID, fID. Note that a new window is required here that does not fully agree with the preceding window definition. In AQuery, the sums() function does a similar task without resorting to any additional sort. The main query in SQL:1999 uses these results in exactly the same way as the SELECT clause of the AQuery rendition.

Once more the query structure impacted the quality of the plans found, which are shown in Figure 6.6. The main difference between the two plans is an extra sort on the SQL:1999 one. Its optimizer added a sort by the same columns it is grouping by. By contrast, AQuery’s group by is dependent on – and thus benefits from – the order enforced by the ASSUMING clause. The SQL:1999 optimizer did consider both windows to have the same ordering requirements, though, by sorting only once by source, destination, and timestamp. It is unclear from the SQL:1999 plan documentation how the ptime attribute is calculated.

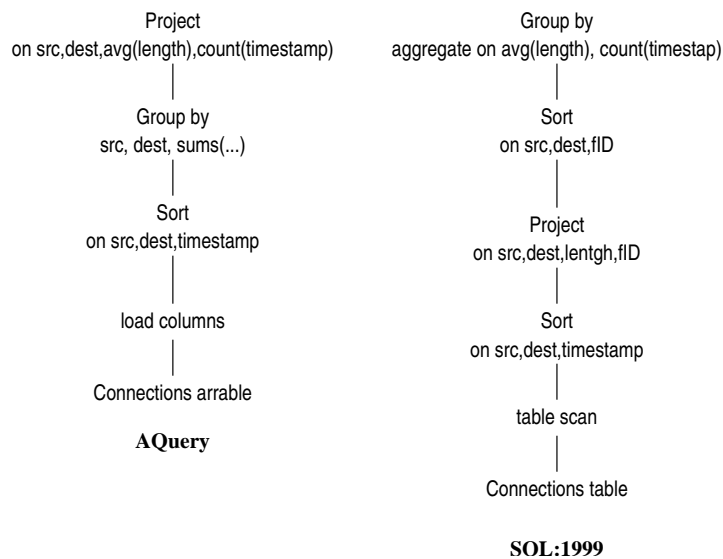


Figure 6.6: Plans for the network management query

AQuery took advantage of the existing order and used a stable algorithm to eliminate that column from the sort.

Here again the difference in the plans brought discrepant response times. The chart in Figure 6.7 shows the relative improvement of the AQuery plan over SQL:1999's. We used a Packets arrable/table with 100 sessions and varied the number of packets for each session from 2K to 10K. AQuery results were from a little less than 2.4 to 3 times faster.

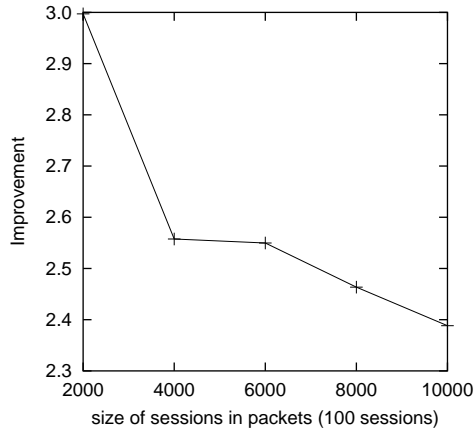


Figure 6.7: Network management query relative improvement

A clear component of the worse response time is the extra sort but that alone cannot account for the entire time difference. We believe that the vertical-partitioning, array-processing approach of AQuery is playing an important role here. For one, processing a column at a time as opposed to a row at a time demands far less overhead. While in the latter case, each operator is called once every time an operator needs a row (iterator model), in AQuery there is only one call per operator (full columns are returned). For another, vertical partitioning and array processing combined highly favored locality of reference.

## 6.4 Conclusion

These experiments suggest basic flaws in the straightforward implementation of SQL:1999 for order queries. In the best-profit query, ordering was part of a more complex operation, a window, and therefore the optimizer missed the opportunity of pushing down a highly selective filtering operation through it. The moral here is that *if order implies complex structures (windows) then even trivial transformations (selection push-down) may require serious analysis.*

In the network management query, the ordering by source, destination, and flow ID is equivalent to that by source, destination, and timestamp. Flow ID grows monotonically with the latter. The optimizer missed that and inserted an extra sort step in the query. The moral here is that *the more complex the syntax, the harder it is to find order idioms.*

The experiments showed that AQuery's structural simplicity helped find much

better plans. This translated into performance improvements often greater than an order of magnitude.

# Chapter 7

## Conclusion

### 7.1 Summary

A query in AQuery may determine the order through which it wishes to manipulate data. The query's clauses can exploit this order by the use of array-typed expressions and vector-to-vector functions. The use of order does not force a query to use a more complex structure than it normally would. As a consequence, queries that in other languages are hard to write become as natural in AQuery as SQL is for unordered queries.

The data ordering in a query is independent of the order in which data is stored. If these orderings match, then the query may execute more efficiently, but there is no semantic effect of the order. If these orderings don't match, then the sort work involved can often be diminished. AQuery's conciseness fosters recognition of common order idioms for which we have provided some effective optimization techniques. The techniques complement the vast body of knowledge in query transformations which mostly apply to AQuery as well.

The AQuery system is the validation of these ideas. It executes AQuery queries over arrables by breaking them down into a sequence of array primitives. We have used the system to successfully write several queries that occur in the finance and network management domains.

### 7.2 Ongoing Work

The motivational study we conducted in several application areas showed that ad-hoc order-dependent querying is virtually unavailable commercially. (SQL:1999 systems are known not to be widely used in finances, Biology, or network management.) We are responding to that need with a major effort for making AQuery public. We are building a database of regression tests that cover the entire language. Our ultimate goal is to have a solid system by the time of its first release.

As a parallel effort, we continue to investigate other query transformation for order dependent queries. Nevertheless, we are still experimenting on how to make the AQuery system apply these transformations automatically.

## 7.3 Future Work

Our current implementation carries several performance-improving operations such as edgeby or sort-edge. However, we have not yet touched other possibilities for gaining performance through parallelism or modern-architecture hardware exploitation (super-scalar CPUs and hierarchical memory). Our extensive use of arrays makes the latter seem particularly promising. Current super-scalar CPUs can use Single Instruction Multiple Data (SIMD) parallelism; we therefore can convert several of our array primitives to use this facility.

Finance is a domain where often order-dependent querying involves streaming data. AQuery showed a natural facility to express streams-based queries. We have however identified important queries in which the analysis of very recent data may be triggered by events on the head of the stream. This quasi-streaming approach in which one would need to backtrack to recent elements of a stream is an avenue that interests us.

Biological sequence databases for DNA or proteins make extensive use of order-dependent operations. The order idioms and the optimization techniques found here would still be valid, but the query language interface may need to be rethought. This is another avenue that we hope to pursue.

# Bibliography

- [1] Anastassia Ailamaki, David J. DeWitt, Mark D. Hill, and Marios Skounakis. Weaving Relations for Cache Performance. In *Proceedings of the International Conference on Very Large Data Bases (VLDB)*, pages 169–180, 2001.
- [2] Peter Bauman, Andreas Dehmel, Paula Furtado, Roland Ritsch, and Norbert Widmann. The Multidimensional Database System RasDaMan. In *Proceedings of the 1998 ACM SIGMOD International Conference on Management of Data (SIGMOD)*, pages 575–577, 1998.
- [3] Peter A. Boncz. *Monet: A Next-Generation DBMS Kernel for Query-Intensive Applications*. PhD thesis, Universiteit van Amsterdam, The Netherlands, 2002.
- [4] Anthony J. Bonner and Giansalvatore Mecca. Sequences, Datalog, and Transducers. *Journal of Computer and System Sciences (JCSS)*, 57(3):234–259, 1998.
- [5] Timothy Budd. *An APL Compiler*. Springer-Verlag, 1988.
- [6] Peter Buneman, Leonid Libkin, Dan Suciu, Val Tannen, and Limsoon Wong. Comprehension Syntax. *SIGMOD Record*, 23(1):87–96, 1994.
- [7] Peter Buneman, Shamim Navqi, Val Tannen, and Limsoon Wong. Principles of Programming with Complex Objects and Collection Types. *Theoretical Compute Science*, 149(1):3–48, 1995.
- [8] Michael J. Carey and Donald Kossmann. On Saying ‘Enough Already!’ in SQL. In *Proceedings of the 1997 ACM SIGMOD International Conference on Management of Data (SIGMOD)*, pages 219–230, 1997.
- [9] Chuck Cranor, Yuan Gao, Theodore Johnson, Vlaidslav Shkapenyuk, and Oliver Spatscheck. Gigascope: High Performance Network Monitoring with an SQL Interface. In *Proceedings of the 2002 ACM SIGMOD International Conference on Management of Data (SIGMOD)*, pages 623–623, 2002.
- [10] Marco de Vivo, Eddy Carrasco, Germinal Isern, and Gabriela O. de Vivo. A Review of Port Scanning Techniques. *ACM Computer Communications Review*, 29(2):41–48, 1999.

- [11] Hector Garcia-Molina, Jeffrey Ullman, and Jennifer Widom. *Database System Implementation*. Prentice-Hall, 1999.
- [12] Goetz Graefe. Query Evaluation Techniques for Large Databases. *ACM Computing Surveys*, 25(2):73–170, 1993.
- [13] Goetz Graefe and William McKenna. The Volcano Optimizer Generator: Extensibility and Efficient Search. In *Proceedings of the International Conference on Data Engineering (ICDE)*, pages 209–218, 1993.
- [14] Torsten Grust. *Comprehending Queries*. PhD thesis, University of Konstanz, 1999.
- [15] Yannis E. Ioannidis and Eugene Wong. Query Optimization by Simulated Annealing. In *Proceedings of the 1987 ACM SIGMOD International Conference on Management of Data (SIGMOD)*, pages 9–22, 1987.
- [16] ISO/IEC 9075. *Information Technology - Database Languages - SQL*, 1999.
- [17] ISO/IEC 9075. *Amendment 1:2001 - Information Technology - Database Languages - SQL (SQL/OLAP)*, 2001.
- [18] Kenneth E. Iverson. *A Programming Language*. Wiley, 1962.
- [19] Kaippallimalil J. Jacob and Dennis Shasha. FinTime - A Financial Time Series Benchmark. *SIGMOD Record*, 28(4):42–48, 1999.
- [20] KX Systems. *K Reference Manual*.
- [21] KX Systems. *KSQL Reference Manual*.
- [22] Leonid Libkin, Rona Machlin, and Limsoon Wong. A Query Language for Multidimensional Arrays: Design, Implementation, and Optimization Techniques. In *Proceedings of the 1996 ACM SIGMOD International Conference on Management of Data (SIGMOD)*, pages 228–239, 1996.
- [23] Guy M. Lohman. Grammar-like Functional Rules for Representing Query Optimization Alternatives. In *Proceeding of the 1988 ACM SIGMOD International Conference on Management of Data (SIGMOD)*, pages 18–27, 1988.
- [24] David Maier and Bennet Vance. A Call to Order. In *Proceedings of the Twelfth ACM Symposium on Principles of Database Systems (PODS)*, pages 1–16, 1993.
- [25] Stefan Manegold, Peter A. Boncz, and Martin L. Kersten. Optimizing Database Architecture for the New Bottleneck: Memory Access. *The VLDB Journal*, 9(3):231–246, 2000.
- [26] Arunprasad P. Marathe and Kenneth Salem. Query Processing Techniques for Arrays. *The VLDB Journal*, 11(1):68–91, 2001.

- [27] Jim Melton. *Advanced SQL:1999 – Understanding Object-Relational and Other Advanced Features*. Morgan Kaufmann Publishers, 2002.
- [28] Wilfred Ng. An Extension of the Relational Data Model to Incorporate Ordered Domains. *ACM Transactions on Database Systems (TODS)*, 26(3):344–383, 2001.
- [29] Oracle. *Analytic SQL Features in Oracle 9i*, december 2001.
- [30] Raghu Ramakrishnan, Donko Donjerkovic, Arvind Ranganathan, Kevin S. Beyer, and Muralidhar Krishnaprasad. SRQL: Sorted Relational Query Language. In *Proceedings 10th International Conference on Scientific and Statistical Database Management (SSDBM)*, pages 84–95, 1998.
- [31] Praveen Seshadri, Miron Livny, and Raghu Ramakrishnan. SEQ: A Model for Sequence Databases. In *Proceedings of the 11th Intenational Conference on Data Engineering (ICDE)*, pages 232–239, 1995.
- [32] Praveen Seshadri, Miron Livny, and Raghu Ramakrishnan. The Design and Implementation of a Sequence Database System. In *Proceedings of the International Conference on Very Large Data Bases (VLDB)*, pages 99–110, 1996.
- [33] Dennis Shasha. Time Series in Finances. Summer School in Extending Database Technologies in La Baule, France, 1999.
- [34] David E. Simmen, Eugene J. Shekita, and Timothy Malkemus. Fundamental Techniques for Order Optimization. In *Proceedings of the 1996 ACM SIGMOD International Conference on Management of Data (SIGMOD)*, pages 57–67, 1996.
- [35] Giedrius Slivinskas. *A Middleware Approach to Temporal Query Processing*. PhD thesis, Aalborg University, Denmark, 2001.
- [36] Giedrius Slivinskas, Christian S. Jensen, and Richard T. Snodgrass. Bringing Order to Query Optimization. *SIGMOD Record*, 31(2):5–14, 2002.
- [37] Igor Tatarinov, Stratis Viglas, Kevin S. Beyer, Jayavel Shanmugasundaram, Eugene J. Shekita, and Chun Zhang. Storing and querying ordered XML using a relational database system. In *Proceedings of the 2002 ACM SIGMOD International Conference on Management of Data (SIGMOD)*, pages 204–215, 2002.
- [38] Arthur Whitney and Dennis Shasha. Lots o’ Ticks: Real-Time High Performance Time Series Queries on Billions of Trades and Quotes. In *Proceedings of the 2001 ACM SIGMOD International Conference on Management of Data (SIGMOD)*, 2001.
- [39] Jingren Zhou and Kenneth A. Ross. A Multi-Resolution Block Storage Model for Database Design. In *To appear in the Proceedings of the 2003 International Database Engineering and Application Symposium (IDEAS)*, 2003.

- [40] Yunyue Zhu and Dennis Shasha. StatStream: Statistical Monitoring of Thousands of Data Streams in Real Time. In *Proceedings of the International Conference on Very Large Data Bases (VLDB)*, pages 358–369, 2002.