

Lecture Notes in Computer Science

1182

Edited by G. Goos, J. Hartmanis and J. van Leeuwen

Advisory Board: W. Brauer D. Gries J. Stoer

Waqar Hasan

Optimization of SQL Queries for Parallel Machines



Springer

Series Editors

Gerhard Goos, Karlsruhe University, Germany

Juris Hartmanis, Cornell University, NY, USA

Jan van Leeuwen, Utrecht University, The Netherlands

Author

Waqar Hasan

Stanford University, Department of Computer Science

Stanford, CA 94305, USA

E-mail: hasan@db.stanford.edu

Cataloging-in-Publication data applied for

Die Deutsche Bibliothek - CIP-Einheitsaufnahme

Hasan, Waqar:

Optimization of SQL queries for parallel machines / Waqar

Hasan. - Berlin ; Heidelberg ; New York ; Barcelona ;

Budapest ; Hong Kong ; London ; Milan ; Paris ; Santa Clara ;

Singapore ; Tokyo : Springer, 1996

(Lecture notes in computer science ; 1182)

Zugl.: Stanford, CA, Univ., Diss.

ISBN 3-540-62065-6

NE: GT

CR Subject Classification (1991): H.2, H.3, E.5

ISSN 0302-9743

ISBN 3-540-62065-6 Springer-Verlag Berlin Heidelberg New York

This work is subject to copyright. All rights are reserved, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, re-use of illustrations, recitation, broadcasting, reproduction on microfilms or in any other way, and storage in data banks. Duplication of this publication or parts thereof is permitted only under the provisions of the German Copyright Law of September 9, 1965, in its current version, and permission for use must always be obtained from Springer - Verlag. Violations are liable for prosecution under the German Copyright Law.

© Springer-Verlag Berlin Heidelberg 1996

Printed in Germany

Typesetting: Camera-ready by author

SPIN 10550447 06/3142 - 5 4 3 2 1 0 Printed on acid-free paper

To my father

Dr. Amir Hasan

for showing me the paths that I follow.

Foreword

Performance in computing, and particularly in data access, is crucial as our dependence on computing becomes pervasive. Continued innovation is essential for increases in performance to keep up with our requirements: at any time in the past there has been a majority of tasks which can be performed in time, while there has been a remainder, the tail of the distribution, that provides a challenge. In addition, we are always facing some problems that are in the infeasible range.

As examples for the three classes of problems, namely satisfactory, challenging, and infeasible, we see in the arena of information processing, respectively, routine business processing, delivering information for decision-making from distributed large databases, and prediction of future events from models and historical data.

Performance in a computing system is improved in two dimensions: using higher speed in the individual modules and increasing the number of modules that can operate in parallel. These two dimensions are not independent, since the coordination and communication required in parallel execution increases the amount of work to be done by each module. Crucial in this balance is hence the granularity, the size of the modules. At the very fine grain the tradeoffs are well understood. Computer architectures, in various generations, have moved from say 8-bit, to 16-bit, 32-bit, and even wider paths. There is, however, a diminishing return here: if the natural data elements are small, and interact to inhibit parallel operation, then wider data paths do not provide an advantage.

Databases deal with mixed granularity. The primitive elemental values being stored are often small, but are aggregated into records of hundreds of elements, and then stored in files or tables containing thousands or millions of records. Relevant data for even a modest problem may be found on dozens of computers distributed anywhere over world-wide networks.

This monograph provides both insights and algorithms pertaining to parallel operation at a practical granularity relevant to database system operations. Tables or object classes provide sizable objects that can be treated in parallel while permitting serial, pipelined overlap. Modules can also be cloned by partitioning and replicating data objects. These approaches interact with each other as well, making the selection of effective processing schedules yet

VIII Foreword

more complex. Fortunately, the algebras over these objects are well-behaved, so that their operations can be completely and precisely defined.

This work provides an approach that balances the advantages and costs of parallel execution. Module granularity is determined by the actual operations being scheduled while respecting intrinsic limits on available parallelism such as timing and data-placement constraints and accounting for the trade-off between using parallel execution and incurring communication costs. The result is applicable to modern system configurations, where computation is performed on pipelining-capable workstations operating in parallel. Further research will have to focus on dynamic aspects of parallel computation, letting the scheduling itself overlap with the computation, since this work seems to be able to exploit all the information likely to be available prior to execution in practical systems.

Stanford, California, USA
September 1996

Gio Wiederhold

Preface

This book is about optimization techniques to determine the best way of exploiting parallel execution for SQL queries against large databases. It is the published version of my PhD dissertation at Stanford University. The techniques in this book are useful in the construction of SQL compilers that can exploit parallel machines effectively.

SQL permits questions to be posed declaratively. Users are insulated from the physical hardware and the layout of the data and thus are able to avoid the complex procedural details of programming a parallel machine. A Data Base Management System (DBMS) answers a SQL query by first finding a procedural plan to execute the query and subsequently executing the plan to produce the query result. This book provides techniques for the problem of *parallel query optimization*: Given a SQL query, find the parallel plan that delivers the query result in minimal time.

I express my gratitude to the people and organizations that made my thesis possible. Gio Wiederhold was a constant source of intellectual support. He encouraged me to learn and use a variety of techniques from different areas of Computer Science. Rajeev Motwani helped enhance my understanding of theory and contributed significantly to the ideas in my thesis. Jeff Ullman was a source of useful discussions and I thank him for his helpful and incisive comments. Ravi Krishnamurthy served as a mentor and a source of interesting ideas and challenging questions. Hector Garcia-Molina provided helpful advice. Jim Gray helped me understand the realities of parallel query processing.

My thesis topic grew out of work at Hewlett-Packard Laboratories and was supported by a fellowship from Hewlett-Packard. I express my gratitude to Hewlett-Packard Company and thank my managers Umesh Dayal, Dan Fishman, Peter Lyngbaek, and Marie-Anne Neimat for management and intellectual and moral support.

I thank Tandem Computers for providing access to a parallel machine, the NonStop SQL/MP parallel DBMS, and for permitting publication of experimental results. I am grateful to Susanne Englert, Ray Glasstone, and Shyam Johari for making this possible and for helping me understand Tandem systems.

The following friends and colleagues were a source of invaluable discussions and diversions: Sang Cha, Surajit Chaudhuri, Philippe DeSmedt, Mike Heytens, Curt Kolovson, Stephanie Leichner, Sheralyn Listgarten, Arif Merchant, Inderpal Mumick, Pandu Nayak, Peter Rathmann, Donovan Schneider, Arun Swami, Kevin Wilkinson, Xiaolei Qian.

My thesis would not have been possible without the support and understanding of my family. I thank my father, Dr. Amir Hasan, for providing the inspiration to pursue a PhD. I thank my mother, Fatima Hasan, my brothers Safdar, Javed, and Zulfiquar, and sister Seemin for their love and encouragement. I owe a debt to my wife Shirin and son Arif for putting up with the long hours that made this work possible.

Stanford, California, USA
September 1996

Waqar Hasan

Abstract

Parallel execution offers a solution to the problem of reducing the response time of SQL queries against large databases. As a declarative language, SQL allows users to avoid the complex procedural details of programming a parallel machine. A DBMS answers a SQL query by first finding a procedural plan to execute the query and subsequently executing the plan to produce the query result. We address the problem of *parallel query optimization*: Given a SQL query, find the parallel plan that delivers the query result in minimal time.

We develop optimization algorithms using models that incorporate the sources of parallelism as well as obstacles to achieving speedup. One obstacle is inherent limits on available parallelism due to parallel and precedence constraints between operators and due to data placement constraints that essentially pre-allocate some subset of operators. Another obstacle is that the overhead of exploiting parallelism may increase total work thus reducing or even offsetting the benefit of parallel execution. Our experiments with Non-Stop SQL, a commercial parallel DBMS, show communication of data across processors to be a significant source of increase in work.

We adopt a two-phase approach to parallel query optimization: *join ordering and query rewrite* (JOQR), followed by *parallelization*. The JOQR phase minimizes the total work to compute a query. The parallelization phase extracts parallelism and schedules resources to minimize response time. We make contributions to both phases. Our work is applicable to queries that include operations such as grouping, aggregation, foreign functions, and set intersection and difference, and joins.

We develop algorithms for the JOQR phase that minimize total cost while accounting for the communication cost of repartitioning data. Using a model that abstracts physical characteristics of data, such as partitioning, as colors, we devise tree coloring algorithms that are efficient and guarantee optimality.

We model the parallelization phase as scheduling a tree of inter-dependent operators with computation and communication costs represented as node and edge weights. Scheduling a weighted operator tree on a parallel machine poses a class of novel multi-processor scheduling problems that differ from the classical in several ways.

We develop and compare several efficient algorithms for the problem of scheduling a pipelined operator tree in which all operators run in parallel

XII Abstract

using inter-operator parallelism. Given the NP-hardness of the problem, we assess the quality of our algorithms by measuring their performance ratio which is the ratio of the response time of the generated schedule to that of the optimal. We prove worst-case bounds on the performance ratios of our algorithms and measure the average cases using simulation.

We address the problem of scheduling a pipelined operator tree using both pipelined and partitioned parallelism. We characterize optimal schedules and investigate two classes of schedules that we term symmetric and balanced.

The results in this thesis enable the construction of SQL compilers that can exploit parallel machines effectively.

Table of Contents

1. Introduction	1
1.1 Minimizing Response Time: Sources and Deterrents	1
1.1.1 Sources of Speedup	2
1.1.2 Deterrents to Speedup	3
1.2 Model for Parallel Query Optimization	4
1.2.1 Annotated Query Trees	5
1.2.2 Operator Trees	5
1.2.3 Parallel Machine Model	7
1.3 Organization of Thesis	8
1.4 Related Work	9
1.4.1 Query Optimization for Centralized Databases	9
1.4.2 Query Optimization for Distributed Databases	9
1.4.3 Query Optimization for Parallel Databases	10
2. Price of Parallelism	13
2.1 Introduction	13
2.2 Tandem Architecture: An Overview	14
2.2.1 Parallel and Fault-Tolerant Hardware	14
2.2.2 Message Based Software	16
2.2.3 Performance Characteristics	16
2.3 Parallelism in NonStop SQL/MP	17
2.3.1 Use of Intra-operator Parallelism	17
2.3.2 Process Structure	18
2.4 Startup Costs	20
2.5 Costs of Operators and Communication	20
2.5.1 Experimental Setup	22
2.5.2 Costs of Scans, Predicates and Aggregation	23
2.5.3 Costs of Local and Remote Communication	24
2.5.4 Cost of Repartitioned Communication	26
2.5.5 Costs of Join Operators	27
2.5.6 Costs of Grouping Operators	30
2.6 Parallel Versus Sequential Execution	31
2.6.1 Parallelism Can Reduce Work	31
2.6.2 Parallelism Can Increase Response Time	33

2.7	Summary of Findings	33
3.	JOQR Optimizations	35
3.1	A Model for Minimizing Communication	36
3.1.1	Partitioning	36
3.1.2	Repartitioning Cost	38
3.1.3	Optimization Problem	38
3.2	Algorithms for Query Tree Coloring	39
3.2.1	Problem Simplification	40
3.2.2	A Greedy Algorithm for Distinct Pre-Colorings	42
3.2.3	Algorithm for Repeated Colors	43
3.2.4	Extensions: Using Sets of Colors	46
3.3	Model for Methods and Physical Properties	48
3.3.1	Annotated Query Trees and Their Cost	50
3.4	Extension of ColorSplit for Methods and Physical Properties ..	52
3.5	Model with Join Ordering	53
3.5.1	Join Ordering Without Physical Properties	54
3.5.2	Join Ordering with Physical Properties	55
3.6	Usage of Algorithms	56
4.	Scheduling Pipelined Parallelism	59
4.1	Problem Definition	59
4.2	Identifying Worthless Parallelism	62
4.2.1	Worthless Edges and Monotone Trees	63
4.2.2	The GreedyChase Algorithm	65
4.2.3	Lower Bounds	65
4.3	The Modified LPT Algorithm	66
4.4	Connected Schedules	68
4.4.1	Connected Schedules When Communication is Free ...	68
4.4.2	BalancedCuts with Communication Costs	73
4.5	Connected Schedules as an Approximation	73
4.6	Heuristics for POT Scheduling	77
4.6.1	A Hybrid Algorithm	78
4.6.2	The Greedy Pairing Algorithm	78
4.7	Approximation Algorithms	79
4.7.1	A Two-Stage Approach	80
4.7.2	The LocalCuts Algorithm	82
4.7.3	The BoundedCuts Algorithm	84
4.8	Experimental Comparison	89
4.8.1	Experimental Setup	90
4.8.2	Experimental Comparison	90
4.8.3	Performance of Hybrid	91
4.8.4	Comparison of Hybrid, LocalCuts and BoundedCuts ..	91
4.8.5	Behavior of Lower Bound	92
4.9	Discussion	94

5. Scheduling Mixed Parallelism	95
5.1 Problem Definition	95
5.2 Balanced Schedules	99
5.3 Symmetric Schedules	102
5.4 Scheduling Trees with Two Nodes	111
5.5 Discussion	112
6. Summary and Future Work	115
6.1 Summary of Contributions	115
6.2 Future Work	118
References	121
Index	131

List of Figures

1.1	Query Processing Architecture	2
1.2	Phases and Sub-phases of Parallel Query Optimization	4
1.3	(A) Annotated Query Tree (B) Corresponding Operator Tree	6
2.1	(A) Tandem Architecture (B) Abstraction as Shared-Nothing	15
2.2	Process Startup: With (Solid) and Without (Dotted) Process Reuse.	21
2.3	Local, Remote and Repartitioned Communication	21
2.4	Scan with 1 Predicate (Dotted), 2 Predicates (Solid), Aggregation (Dashed)	24
2.5	Scan and Aggregation	25
2.6	Process Structure: (A) No Communication (B) Local (C) Remote.	26
2.7	Local and Repartitioned Execution	28
2.8	Local (Dotted) and Repartitioned (Solid) Communication	29
2.9	Query Using Simple-Hash (Dashed), Sort-Merge (Solid) and Nested Join (Dotted)	29
2.10	Hash (Solid) and Sort (Dotted) Grouping Costs	30
2.11	Process Structure: Sequential and Parallel Execution	32
3.1	Query Trees: Hatched Edges Show Repartitioning	37
3.2	(i) Query Tree; (ii) Coloring of Cost 7; (iii) Minimal Coloring of Cost 6	40
3.3	(i) Split Colored Interior Node (ii) Collapse Uncolored Leaves	41
3.4	(i) Query Tree (ii) Suboptimal DLC Coloring (cost=9) (iii) Opti- mal Coloring (cost=8)	43
3.5	Problem Decomposition After Coloring Node i	44
3.6	Opt and Optc Tables for Tree of Figure 3.4	45
3.7	Interaction of Repartitioning with Join Predicates	48
3.8	Annotated Query Trees	49
3.9	Interaction of Repartitioning with Order of Joins	54
3.10	Decomposition of a Complex Query	57
4.1	A Pipelined Schedule and Its Execution	61
4.2	(A) Trace of <i>GreedyChase</i> (Worthless Edges Hatched) (B) Modi- fied LPT Schedule (C) Naive LPT Schedule	66
4.3	Example with Performance Ratio = n/p for <i>Modified LPT</i>	68

XVIII List of Figures

4.4	Connected Schedule as Cutting and Collapsing Edges	69
4.5	Fragments Formed by <i>BpSchedule</i> Before the Last Stage of <i>BalancedCuts</i>	73
4.6	Examples with $\frac{L_C}{L_{opt}} = 2 - \frac{1}{\lceil \frac{x+1}{2} \rceil}$	75
4.7	Performance Ratio=3 for Star of 10 Nodes Scheduled on 5 Processors	77
4.8	Subtrees $T_m, T_{m'}, T_{m''}$ for Nodes m, m', m''	86
4.9	C_{opt}^m	86
4.10	Performance of Hybrid (Solid), BalancedFragments (Dotted) and Modified LPT (Dashed) on Wide Trees	91
4.11	Performance of Hybrid (Solid), BalancedFragments (Dotted) and Modified LPT (Dashed) on Narrow Trees	92
4.12	Comparison of Hybrid (Solid), LocalCuts (Dashed) and Bound-edCuts (Dotted) on Narrow Trees	92
4.13	Comparison of Hybrid (Solid), LocalCuts (Dashed) and Bound-edCuts (Dotted) on Wide Trees	93
4.14	Performance of Optimal (Dotted) and Hybrid (Solid)	93
5.1	Execution with Mixed Parallelism	97
5.2	Structure of (Strongly) Minimal Schedule	102
5.3	Matrices for $p = 3$	107
5.4	Counter-Example: Tree for Which Symmetric Schedule is a Saddle Point	110
5.5	Plot of $z = a_{11} + a_{21} - 2a_{11}a_{21}$ with a_{11} on x-Axis, a_{21} on y-Axis	111
5.6	One Sided Schedule	113
5.7	Balanced Schedule for $n=2$ (Some Communication Arcs omitted)	113
6.1	Phases and Sub-phases of Parallel Query Optimization	116

List of Tables

2.1	Parallelization Strategies and Join Methods	19
2.2	CPU Costs of Transfer and Computational Operations. (1K Tuples Occupy 1 Mbyte)	22
3.1	Examples of Input-Output Constraints	51