

Lecture Notes in Computer Science
Edited by G. Goos, J. Hartmanis, and J. van Leeuwen

2169

Springer

Berlin

Heidelberg

New York

Barcelona

Hong Kong

London

Milan

Paris

Tokyo

Michael Jaedicke

New Concepts for Parallel Object-Relational Query Processing



Springer

Series Editors

Gerhard Goos, Karlsruhe University, Germany
Juris Hartmanis, Cornell University, NY, USA
Jan van Leeuwen, Utrecht University, The Netherlands

Author

Michael Jaedicke
Donnersbergerstr. 34
80634 München
Germany
E-mail: mjaedicke@yahoo.com

Cataloging-in-Publication Data applied for

Die Deutsche Bibliothek - CIP-Einheitsaufnahme

Jaedicke, Michael: New concepts for parallel object relational query processing /
Michael Jaedicke. - Berlin ; Heidelberg ; New York ; Barcelona ;
Hong Kong ; London ; Milan ; Paris ; Tokyo : Springer, 2001
(Lecture notes in computer science ; Vol. 2169)
Zugl.: Stuttgart, Univ., Diss., 1999
ISBN 3-540-42781-3

CR Subject Classification (1998): E.2, H.2

ISSN 0302-9743

ISBN 3-540-42781-3 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 New York
a member of BertelsmannSpringer Science+Business Media GmbH

<http://www.springer.de>

© Springer-Verlag Berlin Heidelberg 2001
Printed in Germany

Typesetting: Camera-ready by author, data conversion by Christian Grosche, Hamburg
Printed on acid-free paper SPIN: 10840389 06/3142 5 4 3 2 1 0

Preface

When I started this research, no commercial database system supported both object-relational features and parallelism. In the meantime this situation has changed dramatically. All major vendors now offer a parallel, object-relational DBMS as their high-end product. However, there is still a lot to do. Firstly, object-relational (or extensible) DBMS have yet to mature fully. Secondly, the integration of parallelism and extensibility has not yet been completed. This work is my attempt to make a contribution to both issues.

Some of the concepts and techniques developed have been implemented in a prototypical parallel database system called MIDAS. This system is the result of a team effort to which many people have contributed. My own contributions to the system are the user-defined functions and user-defined table operators, the extension of the system tables for those user-defined objects and for detailed statistics, and the new query optimizer, for which I worked with Clara Nippl. She contributed especially to the cost model, the physical operators, and implementation rules. Furthermore, I provided support for the development of concepts for the query parallelizer, the query execution control system, and the send/receive operators of the execution system.

Acknowledgments

I am very thankful to Professor Dr.-Ing. habil. Bernhard Mitschang, who gave me the opportunity to carry out this research. His continuous support, encouragement, and useful comments during many fruitful discussions and during the preparation of joint publications are especially acknowledged.

I am very grateful to Prof. Dr.-Ing. Theo Härdter who introduced me thoroughly to database systems and taught me a lot about engineering. I am also grateful to him for the analysis of my results and his valuable remarks.

I also acknowledge the help of my colleagues Giannis Bozas, Clara Nippl, Angelika Reiser, and Stephan Zimmermann with whom I worked together on the MIDAS project. Stephan provided software engineering support for the whole group and was always helpful. His work centered on PVM communication, the process architecture, buffer and transaction management, the parallelization of the execution system, the

scheduler and benchmarking, and performance. Clara worked on the parallelization of the execution system, the scheduler, the parallelizer, and the optimizer. Giannis focussed on the lock manager and the buffer management. I also thank Prof. Rudolf Bayer, who led the project together with Prof. Mitschang, for his support.

I also enjoyed giving practical courses on RDBMS together with Angelika. I thank all my colleagues in Munich and Stuttgart, especially Aiko Frank, Volker Markl, Jürgen Sellentin, and Martin Zirkel for their help and comradeship during the last years. I am grateful to my colleague Ralf Rantzau and to Prof. Härdter for improving my English.

Special thanks also go to the students that I have supervised: Pascal Frantz, Stefan Haas, Karl Hahn, Sebastian Heupel, Bernd Holzhey, Kay Krueger-Barvels, Sabine Perathoner, Ralf Schumacher, Robert Seehafer, and Susanne Stamp. It was a pleasure to work with them. Many of them made significant contributions to this work. I also thank the many other students who worked on the MIDAS project.

I also gratefully acknowledge the valuable comments of the anonymous referees of diverse papers, which have also improved this work. The feedback from C. Mohan, G. Lohman, and M. Carey on the issue of interfaces for query execution plans and the feedback from Harald Schöning on the parallel execution of user-defined functions is much appreciated.

Finally, I would like to thank the Deutsche Forschungsgemeinschaft (DFG) for funding the MIDAS project (SFB 342, B2). I also acknowledge the support of the Studienstiftung des Deutschen Volkes during my university course.

Special thanks go to Springer-Verlag, especially to Alfred Hofmann, and the series editors, for publishing my work in LNCS.

Last but not least I am grateful to my parents for their affection and continuous care.

July 2001

Michael Jaedicke

Abstract

During the last few years parallel object-relational database management systems have emerged as the leading data management technology on the market place. These systems are extensible by user-defined data types and user-defined functionality for the data. This work focuses on the efficient parallel execution of user-defined functionality. The main contributions describe techniques to support data parallelism for user-defined scalar and aggregate functions, to support intra-function parallelism for the execution of a scalar function on a large object, and a new technology to provide extensibility with regard to new set-oriented database operations that can efficiently implement user-defined functionality in parallel object-relational database management systems. Some of these techniques have been implemented in the MIDAS prototype or on top of a commercial object-relational database management system.

Table of Contents

CHAPTER 1 Introduction

| | | |
|-----|-----------------------------------|---|
| 1.1 | ORDBMS: The Next Great Wave | 1 |
| 1.2 | Extensible DBMS | 2 |
| 1.3 | Overview | 3 |

CHAPTER 2 Background on User-Defined Routines

| | | |
|-------|--|----|
| 2.1 | User-Defined Routines | 5 |
| 2.2 | Definition, Implementation, and Execution of New UDR | 6 |
| 2.2.1 | User-Defined Scalar Functions | 7 |
| 2.2.2 | User-Defined Aggregate Functions | 9 |
| 2.2.3 | User-Defined Table Functions | 10 |
| 2.2.4 | User-Defined Functions and Large Objects | 11 |
| 2.3 | Comparison with Stored Procedures | 12 |
| 2.4 | Optimization of Queries with UDF | 12 |

CHAPTER 3 Parallel Processing of User-Defined Functions

| | | |
|-------|---|----|
| 3.1 | Introduction | 14 |
| 3.2 | Limits of Current ORDBMS | 15 |
| 3.3 | Parallel Processing of UDF | 17 |
| 3.3.1 | Two Step Parallel Aggregation of UDAF | 17 |
| 3.3.2 | Partitioning Classes and Partitionable Functions | 18 |
| 3.3.3 | Parallel Sorting as a Preprocessing Step for UDAF | 21 |
| 3.3.4 | Extended Syntax for Function Registration | 22 |
| 3.4 | Example Applications | 24 |
| 3.4.1 | The UDAF Most_Frequent | 24 |
| 3.4.2 | The UDSF Running_Average | 25 |
| 3.4.3 | The UDAF Median | 25 |
| 3.4.4 | Further Applications | 26 |
| 3.5 | Plausibility Considerations Regarding Performance | 28 |
| 3.6 | Related Work | 30 |
| 3.7 | Summary | 31 |

CHAPTER 4 Intra-function Parallelism

| | | |
|-------|--|----|
| 4.1 | Introduction | 33 |
| 4.2 | Compose/Decompose Operators for Intra-function Parallelism | 34 |
| 4.2.1 | Compose/Decompose Operators | 34 |
| 4.2.2 | Extensibility of Compose Operators by Combine Functions | 36 |
| 4.2.3 | Application of Intra-function Parallelism | 37 |
| 4.2.4 | Intra-function Parallelism for Function Pipelines | 38 |
| 4.3 | Experimental Performance Study | 39 |
| 4.3.1 | Experimental Scenario and Implementation | 39 |
| 4.3.2 | Performance Results | 41 |
| 4.4 | Related Work | 43 |
| 4.5 | Summary | 44 |

CHAPTER 5 The Multi-operator Method

| | | |
|-------|---|----|
| 5.1 | Introduction | 45 |
| 5.2 | Performance Problems with Complex UDF in Current ORDBMS | 46 |
| 5.2.1 | The PBSM Algorithm as a Sophisticated UDP Implementation | 47 |
| 5.3 | The Multi-operator Method as a New Technique to Implement Complex UDF | 49 |
| 5.3.1 | The Multi-operator Method and Its Benefits | 49 |
| 5.3.2 | A Multi-operator Implementation of the PBSM Algorithm | 51 |
| 5.4 | Supporting the Multi-operator Method | 53 |
| 5.4.1 | Executing Query Execution Plans | 53 |
| 5.4.2 | Example for a Textual Specification of Query Execution Plans | 55 |
| 5.4.3 | Parallel Evaluation | 55 |
| 5.5 | Performance Evaluation | 56 |
| 5.5.1 | Experimental Scenario | 56 |
| 5.5.2 | Performance Results | 62 |
| 5.6 | Related Work | 64 |
| 5.7 | Summary | 65 |

CHAPTER 6 User-Defined Table Operators

| | | |
|-------|---|-----|
| 6.1 | Introduction | 67 |
| 6.2 | User-Defined Table Operators | 68 |
| 6.2.1 | A Generalization Relationship for Row Types | 68 |
| 6.2.2 | Defining and Implementing UDTO | 69 |
| 6.2.3 | The Different Usages of the UDTO Concept | 74 |
| 6.2.4 | Parallel Processing of Procedural UDTO | 77 |
| 6.2.5 | Extension to Multiple Output Tables | 80 |
| 6.3 | Example Applications for UDTO | 81 |
| 6.3.1 | Computing a Spatial Join | 81 |
| 6.3.2 | Different UDTO for the Same Predicate | 85 |
| 6.3.3 | Computing the Median: An Aggregation Operator | 89 |
| 6.3.4 | A UDTO for a Complex Aggregation | 90 |
| 6.3.5 | Association Rule Mining | 94 |
| 6.4 | Related Work | 101 |
| 6.5 | Summary and Conclusions | 102 |

CHAPTER 7 Implementation of UDTO

| | | |
|-------|--|-----|
| 7.1 | Introduction | 106 |
| 7.2 | The MIDAS Prototype | 106 |
| 7.2.1 | Architectural Overview | 107 |
| 7.2.2 | Query Compilation and Execution | 108 |
| 7.2.3 | The MIDAS System Tables | 111 |
| 7.2.4 | UDSF in MIDAS | 112 |
| 7.3 | Implementation of SQL Macros | 113 |
| 7.3.1 | DDL Statements | 113 |
| 7.3.2 | SQL Macro Expansion in DML Statements | 115 |
| 7.3.3 | Expanding SQL Macros in Preprocessors and Middleware | 116 |

| | | |
|-------|--|-----|
| 7.4 | Implementation of Procedural UDTO | 123 |
| 7.4.1 | Extensions to the SQL Compiler | 123 |
| 7.4.2 | Extensions to the Optimizer and the Parallelizer | 125 |
| 7.4.3 | Extensions to the Scheduler | 126 |
| 7.4.4 | Extensions to the Execution Engine | 126 |
| 7.4.5 | Extensions to Transaction Management | 128 |
| 7.4.6 | Implementation of Input and Output Tables | 131 |
| 7.5 | Optimization Issues for UDTO | 134 |
| 7.5.1 | UDTO and Implied Predicates | 134 |
| 7.5.2 | Estimating Costs and Selectivity of UDTO | 135 |
| 7.5.3 | Application of Traditional Optimization Rules | 137 |
| 7.6 | Using UDTO to Generate Alternative Execution Plans for UDF | 138 |
| 7.7 | Evaluation of the Implementation | 139 |
| 7.7.1 | Evaluation of SQL Macros | 140 |
| 7.7.2 | Evaluation of Procedural UDTO | 142 |
| 7.8 | Summary | 144 |

CHAPTER 8 Summary, Conclusions, and Future Work

| | | |
|-----|-------------------|-----|
| 8.1 | Summary | 145 |
| 8.2 | Conclusions | 146 |
| 8.3 | Future Work | 149 |

References

| | |
|------------------|-----|
| References | 151 |
|------------------|-----|

Appendix A

| | | |
|-----|--|-----|
| A.1 | The Program sequential_invert | 157 |
| A.2 | The Program parallel_invert | 158 |
| A.3 | The Query Execution Plan for the Spatial Join with SQL Macro | 159 |