

Data Warehouse Performance

W.H. Inmon

Ken Rudin

Christopher K. Buss

Ryan Sousa

| | |
|--|--------------------|
| Technische Universität Darmstadt FACHBEREICH INFORMATIK | |
| B I B L I O T H E K | |
| Inventar-Nr.: | <u>1199-00268</u> |
| Sachgebiete: | <u>H. 2 / Inmo</u> |
| Standort: | <u>1993</u> |

WILEY COMPUTER PUBLISHING

Fachbereichsbibliothek Informatik
TU Darmstadt



59380117



John Wiley & Sons, Inc.

New York • Chichester • Weinheim • Brisbane • Singapore • Toronto

Contents

| | | |
|------------------------|---|-------------|
| Acknowledgments | | xi |
| Preface | | xiii |
| Chapter One | Introduction to Data Warehouse Performance | 1 |
| | Measuring Performance | 2 |
| | Productivity and Performance | 3 |
| | Data Volumes and Performance | 4 |
| | <i>Why the Growth of Data?</i> | 5 |
| | <i>Data Gets in the Way</i> | 6 |
| | <i>Finding Hidden Data</i> | 6 |
| | User Expectations and Performance | 10 |
| | Education and Performance | 12 |
| | Achieving Performance | 13 |
| | When Should Performance Be Considered? | 13 |
| | Summary | 15 |
| PART ONE | Usage, Data, and Performance | |
| Chapter Two | User Community and Performance | 19 |
| | Know Your End Users: Farmers and Explorers | 19 |
| | <i>The World of the Explorer</i> | 22 |
| | <i>Profile of Execution</i> | 26 |
| | <i>World of the Farmer</i> | 27 |
| | Typing the Farmer and the Explorer | 30 |
| | <i>Clerical versus Management</i> | 31 |

| | | |
|----------------------|---|-----------|
| | <i>Casual versus Power</i> | 32 |
| | <i>Predefined versus Ad Hoc</i> | 32 |
| | <i>Summary versus Detailed</i> | 33 |
| | <i>Simple versus Complex</i> | 34 |
| | Summary | 34 |
| Chapter Three | Farmers and Explorers | 37 |
| | Optimizing for the Explorer | 39 |
| | <i>Subsetting and Removing Data</i> | 40 |
| | <i>Creating a Living Sample</i> | 42 |
| | <i>Exploration Warehouse</i> | 44 |
| | Optimizing for The Farmer | 54 |
| | <i>Strategic Approaches to Performance</i> | 55 |
| | <i>Tactical Approaches to Performance</i> | 57 |
| | <i>Operational Approaches to Performance</i> | 64 |
| | Summary | 72 |
| Chapter Four | Data Marts | 75 |
| | What is a Data Mart? | 75 |
| | <i>Data Mart Community</i> | 77 |
| | <i>Data Mart Appeal</i> | 77 |
| | <i>Data Mart Source</i> | 79 |
| | Building the Data Mart First | 80 |
| | <i>Different Kinds of Data Marts</i> | 81 |
| | <i>Loading the Data Mart</i> | 82 |
| | <i>Metadata in the Data Mart</i> | 84 |
| | <i>Data Modeling for the Data Mart</i> | 85 |
| | <i>Purging the Data Mart</i> | 86 |
| | <i>Data Mart Contents</i> | 87 |
| | <i>Structure within the Data Mart</i> | 87 |
| | Performance in the Data Mart | 90 |
| | Monitoring the Data Mart Environment | 91 |
| | Summary | 92 |
| Chapter Five | Dormant Data | 95 |
| | Understanding Dormant Data | 96 |
| | <i>Summary Tables and Dormant Data</i> | 98 |
| | <i>Misjudgment of History and Dormant Data</i> | 98 |
| | <i>Reality of Requirements and Dormant Data</i> | 98 |
| | <i>Insistence of Detail and Dormant Data</i> | 99 |

| | | |
|--|--|------------|
| | Calculating Dormant Data | 99 |
| | Finding Dormant Data | 101 |
| | Removing Dormant Data | 102 |
| | <i>Selecting Data to be Removed</i> | 103 |
| | <i>Determining the Probability of Access</i> | 105 |
| | Summary | 107 |
| Chapter Six | Data Cleansing | 109 |
| | How Dirty Data Gets In | 109 |
| | Cleansing Dirty Data | 111 |
| | <i>Cleansing the Legacy Environment</i> | 111 |
| | <i>Cleansing at the Point of Integration</i> | 113 |
| | <i>Cleansing after Loading</i> | 114 |
| | Different Kinds of Audits | 116 |
| | Managing Required Resources | 116 |
| | Cleansing Data Over Time | 119 |
| | Bounded Referential Integrity | 120 |
| | Summary | 122 |
| Chapter Seven | Monitors | 125 |
| | Activity Monitors | 125 |
| | <i>Finding Dormant Data</i> | 127 |
| | <i>Understanding Dormant Data</i> | 129 |
| | <i>Removing Dormant Data</i> | 130 |
| | <i>Capturing Activity Information</i> | 131 |
| | <i>Reviewing the Output</i> | 136 |
| | Resource Governors Versus Query Blocking | 138 |
| | <i>Resource Governors—Nothing New</i> | 138 |
| | <i>Why Are Resource Governors Inadequate?</i> | 138 |
| | <i>Query Blocking</i> | 142 |
| | Monitoring Data Content | 143 |
| | Data Warehouse Alarm Clock | 147 |
| | Summary | 150 |
| PART TWO Platform and Performance | | |
| Chapter Eight | Components of the High-Performance Platform | 155 |
| | Performance Chain | 156 |
| | Scalability Requirement | 156 |

| | |
|--|-----|
| Parallelism and Its Relationship to Performance | 159 |
| <i>What Is Parallelism?</i> | 159 |
| <i>Types of Parallelism</i> | 162 |
| High-Performance Hardware | 166 |
| <i>Symmetric Multiprocessors (SMPs)</i> | 167 |
| <i>Clusters</i> | 170 |
| <i>Massively Parallel Processors (MPPs)</i> | 174 |
| <i>Nonuniform Memory Access (NUMA)</i> | 178 |
| High-Performance Databases | 182 |
| <i>Parallel Queries</i> | 182 |
| <i>Shared-Disk and Shared-Nothing Database Architectures</i> | 186 |
| Other Parts of the Performance Chain | 190 |
| <i>The Extract/Transform/Load Component</i> | 190 |
| <i>End-User Access Tools</i> | 193 |
| <i>Scalable Application Frameworks</i> | 194 |
| Summary | 196 |

| | | | |
|---------------------|----------|---|------------|
| Chapter Nine | X | Building a High-Performance Platform | 197 |
| | | System Architecture | 198 |
| | | <i>Three-Tiered Architectures</i> | 199 |
| | | <i>Two-Tiered Architectures</i> | 204 |
| | | <i>The Solution: Scalable Data Marts</i> | 206 |
| | | Building a Balanced Hardware System | 210 |
| | | <i>Estimating the Business Requirements</i> | 211 |
| | | <i>Determining the Technical Configuration</i> | 213 |
| | | <i>Iterate</i> | 216 |
| | | Designing the Physical Database for Performance | 218 |
| | | <i>Denormalizing the Database</i> | 219 |
| | | <i>Indexing Your Data</i> | 226 |
| | | <i>Designing Your Disk Layout</i> | 232 |
| | | Taking Advantage of I/O Parallelism | 234 |
| | | <i>Striping Techniques</i> | 234 |
| | | <i>Table Partitioning Techniques</i> | 242 |
| | | Optimizing the Queries | 248 |
| | | <i>Execution of an SQL Query</i> | 249 |
| | | <i>Effects of Parallelism</i> | 251 |
| | | <i>CPU Utilization</i> | 253 |
| | | <i>Query-Optimization Questions</i> | 254 |
| | | Summary | 256 |

| | | |
|--------------------|--|------------|
| Chapter Ten | Advanced Platform Topics | 259 |
| | Building a Performance Assurance Environment | 259 |
| | <i>Defining the Performance Assurance Metrics</i> | 260 |
| | <i>Building Performance Assurance Tests</i> | 263 |
| | Very Large Database (VLDB) Issues with Data Warehousing | 270 |
| | <i>Custom Code Requirements</i> | 271 |
| | <i>One-In-A-Million Odds Occur Frequently</i> | 272 |
| | <i>Statistical Effects</i> | 273 |
| | <i>Algorithm Changes</i> | 274 |
| | <i>Exceeding Batch Windows</i> | 277 |
| | × Data Warehouses and the Web | 278 |
| | <i>Web Access Means More Users</i> | 280 |
| | <i>Web Access Means More Data</i> | 282 |
| | Data Warehouses and Data Mining | 283 |
| | <i>Data Mining Requires Scalability</i> | 284 |
| | <i>The Basics of Scalable Data Mining</i> | 286 |
| | × Data Warehouses and Object-Relational Databases | 287 |
| | <i>Scalable Performance for Complex Data Types</i> | 288 |
| | <i>Scalable Functionality</i> | 290 |
| | <i>Scalability Issues Regarding Object-Relational Technology</i> | 290 |
| | Summary | 294 |

PART THREE Service Management and Performance

| | | |
|-----------------------|--|------------|
| Chapter Eleven | Service Management and the Service Management Contract | 299 |
| | Service Management Defined | 300 |
| | Business Need for Service Management | 302 |
| | The Service Management Contract | 304 |
| | Creating the Service Management Contract | 305 |
| | <i>Step 1: Establish Relevant Data Warehouse Resources</i> | |
| | <i>Inventory or Services Catalog</i> | 307 |
| | <i>Step 2: Characterize Usage of Data Warehouse Resources</i> | 308 |
| | <i>Step 3: Determine Current or Projected Service Levels</i> | 308 |
| | <i>Step 4: Understand the Customer's Requirements</i> | 309 |
| | <i>Step 5: Determine Cost and Feasibility of Customer's Requirements</i> | 310 |

| | | |
|-------------------------|--|------------|
| | <i>Step 6: Create the Service Management Contract</i> | 311 |
| | <i>Step 7: Track Compliance to the Service Management Contract</i> | 319 |
| | Summary | 320 |
| Chapter Twelve | Putting the Service Management Contract in Motion | 323 |
| | Putting the SMC in Context | 324 |
| | <i>Organization Layer</i> | 325 |
| | <i>Service Dimensions Layer</i> | 325 |
| | <i>Service Reporting Layer</i> | 327 |
| | Data Warehouse Administration (DWA) | |
| | Organization Layer | 327 |
| | Service Administrator | 329 |
| | <i>Measuring the Right Metrics</i> | 329 |
| | <i>DWA May Not Be In IT</i> | 330 |
| | <i>Challenges Existing DW Support Structure</i> | 330 |
| | Service Reporting Layer | 333 |
| | Service Dimensions Layer | 335 |
| | <i>Query Response Time Dimension</i> | 336 |
| | <i>User Concurrency Dimension</i> | 343 |
| | <i>Data Storage Dimension</i> | 348 |
| | <i>System Availability Dimension</i> | 351 |
| | <i>Data Currency Dimension</i> | 356 |
| | <i>Data Quality Dimension</i> | 364 |
| | Summary | 369 |
| | | |
| PART FOUR | Piecing Together the Elements | |
| | | |
| Chapter Thirteen | Delivering a High-Performance Data Warehouse Environment | 373 |
| | Case Study Overview | 375 |
| | <i>Company Background</i> | 376 |
| | <i>Changing Business Landscape</i> | 376 |
| | <i>Marketing Challenges</i> | 377 |
| | <i>Justifying the Data Warehouse Environment</i> | 378 |
| | Building the Team | 378 |
| | Scoping the Three- to Six-Month Deliverable | 380 |
| | Delivering the Data Warehouse Environment | 383 |
| | <i>Assess the Capabilities</i> | 384 |

| | |
|---|------------|
| <i>Aligning the Users, Workload, and Capabilities</i> | 388 |
| <i>Design the Databases</i> | 390 |
| <i>Data Mart and Exploration Warehouse Designs</i> | 393 |
| <i>Data Mart Design</i> | 394 |
| <i>Exploration Warehouse Design</i> | 400 |
| <i>Data Warehouse Design</i> | 404 |
| <i>Configure the Hardware</i> | 415 |
| Servicing the Data Warehouse Environment | 418 |
| Summary | 420 |
| Recommended Reading | 421 |
| Articles | 421 |
| Books | 429 |
| Book Reviews | 430 |
| Index | 431 |