



Discovering SQL

A Hands-On Guide for Beginners

Alex Kriegel

DISCOVERING SQL

INTRODUCTION	xxv
CHAPTER 1	Drowning in Data, Dying of Thirst for Knowledge1
CHAPTER 2	Breaking and Entering: Structured Information 29
CHAPTER 3	A Thing You Can Relate To — Designing a Relational Database 79
CHAPTER 4	Overcoming the Limitations of SQL 103
CHAPTER 5	Grouping and Aggregation137
CHAPTER 6	When One Is Not Enough: A Query Within a Query 155
CHAPTER 7	You Broke It; You Fix It: Combining Data Sets 173
CHAPTER 8	What Else Is There, and Why? 207
CHAPTER 9	Optimizing Performance 231
CHAPTER 10	Multiuser Environment 251
CHAPTER 11	Working with Unstructured and Semistructured Data287
CHAPTER 12	Not by SQL Alone 329
APPENDIX A	Installing the Library Database 353
APPENDIX B	Installing RDBMSs Software 375
APPENDIX C	Accessing RDBMSs 377
APPENDIX D	Accessing RDBMSs with the Squirrel Universal SQL Client. 379
INDEX	381

Discovering SQL

Discovering SQL

A HANDS-ON GUIDE FOR BEGINNERS

Alex Kriegel



Wiley Publishing, Inc.

Discovering SQL

Published by
Wiley Publishing, Inc.
10475 Crosspoint Boulevard
Indianapolis, IN 46256
www.wiley.com

Copyright © 2011 by Wiley Publishing, Inc., Indianapolis, Indiana

Published simultaneously in Canada

ISBN: 978-1-118-00267-4

ISBN: 978-1-118-09279-8 (ebk)

ISBN: 978-1-118-09277-4 (ebk)

ISBN: 978-1-118-09278-1 (ebk)

Manufactured in the United States of America

10 9 8 7 6 5 4 3 2 1

No part of this publication may be reproduced, stored in a retrieval system or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, scanning or otherwise, except as permitted under Sections 107 or 108 of the 1976 United States Copyright Act, without either the prior written permission of the Publisher, or authorization through payment of the appropriate per-copy fee to the Copyright Clearance Center, 222 Rosewood Drive, Danvers, MA 01923, (978) 750-8400, fax (978) 646-8600. Requests to the Publisher for permission should be addressed to the Permissions Department, John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, (201) 748-6011, fax (201) 748-6008, or online at <http://www.wiley.com/go/permissions>.

Limit of Liability/Disclaimer of Warranty: The publisher and the author make no representations or warranties with respect to the accuracy or completeness of the contents of this work and specifically disclaim all warranties, including without limitation warranties of fitness for a particular purpose. No warranty may be created or extended by sales or promotional materials. The advice and strategies contained herein may not be suitable for every situation. This work is sold with the understanding that the publisher is not engaged in rendering legal, accounting, or other professional services. If professional assistance is required, the services of a competent professional person should be sought. Neither the publisher nor the author shall be liable for damages arising herefrom. The fact that an organization or Web site is referred to in this work as a citation and/or a potential source of further information does not mean that the author or the publisher endorses the information the organization or Web site may provide or recommendations it may make. Further, readers should be aware that Internet Web sites listed in this work may have changed or disappeared between when this work was written and when it is read.

For general information on our other products and services please contact our Customer Care Department within the United States at (877) 762-2974, outside the United States at (317) 572-3993 or fax (317) 572-4002.

Wiley also publishes its books in a variety of electronic formats. Some content that appears in print may not be available in electronic books.

Library of Congress Control Number: 2011922790

Trademarks: Wiley, the Wiley logo, Wrox, the Wrox logo, Programmer to Programmer, and related trade dress are trademarks or registered trademarks of John Wiley & Sons, Inc. and/or its affiliates, in the United States and other countries, and may not be used without written permission. All other trademarks are the property of their respective owners. Wiley Publishing, Inc., is not associated with any product or vendor mentioned in this book.

To Liana

ABOUT THE AUTHOR



ALEX KRIEDEL is an Enterprise Systems Architect for the Oregon Health Authority. He has over 20 years of professional experience designing and developing software, implementing and administering enterprise RDBMS, as well managing software development processes. Alex graduated from National Technical University of Belarus with a Master's of Science in Physics of Metals. He also holds several industry certifications, including PMP from Project Management Institute, TOGAF 8 Certified Practitioner from the Open Architecture Group, Certified Scrum Master from Scrum Alliance, and Microsoft Certified Technology Specialist (MCTS) from Microsoft.

Alex provides online training and consulting services through the www.agilitator.com website.

Alex is author of *Microsoft SQL Server 2000 Weekend Crash Course* (Wiley, 2001) and a co-author on several other titles: *SQL Bible* (Wiley, 2003), *SQL Functions* (Wrox, 2005), *Introduction to Database Management* (Wiley, 2007) and *SQL Bible, 2nd Edition* (Wiley, 2008). His books have been translated into Chinese, Portuguese and Russian.

ABOUT THE TECHNICAL EDITOR

BORIS TRUKHNOV is a Principal Oracle Engineer for NexGen Data Systems, Inc. He has been working with relational databases (primarily Oracle) since 1994. Boris is an author of several technical books published in US and translated into Portuguese, Chinese, and Russian, including *SQL Bible* (1st and 2nd editions) and *Introduction to Database Management*.

Boris's areas of expertise include RAC, ASM, RMAN, performance tuning, database and system architecture, platform migrations, and system upgrades.

Boris is an Oracle 11g Database Administrator Certified Professional (OCP) and Oracle Real Application Clusters Administrator (OCE).

CREDITS

EXECUTIVE EDITOR

Robert Elliott

PROJECT EDITOR

Christopher J. Rivera

TECHNICAL EDITOR

Boris Trukhnov

PRODUCTION EDITOR

Rebecca Anderson

COPY EDITOR

Nancy Sixsmith

EDITORIAL DIRECTOR

Robyn B. Siesky

EDITORIAL MANAGER

Mary Beth Wakefield

FREELANCER EDITORIAL MANAGER

Rosemarie Graham

ASSOCIATE DIRECTOR OF MARKETING

David Mayhew

PRODUCTION MANAGER

Tim Tate

**VICE PRESIDENT AND EXECUTIVE GROUP
PUBLISHER**

Richard Swadley

**VICE PRESIDENT AND EXECUTIVE
PUBLISHER**

Barry Pruett

ASSOCIATE PUBLISHER

Jim Minatel

PROJECT COORDINATOR, COVER

Katie Crocker

PROOFREADER

Carrie Hunter, Word One New York

INDEXER

Johnna VanHoose Dinse

COVER DESIGNER

Ryan Sneed

COVER IMAGE

© Henry Chaplin / iStockPhoto

ACKNOWLEDGMENTS

I would like to thank Robert Elliott, executive editor at Wiley Publishing for the wonderful opportunity to work on this book, and for the patience with which he helped me to navigate the editorial process. His friendly managerial style and valuable insights helped to keep the project on track and on time.

Many thanks to the Wiley Editorial team, especially to my project editor, Christopher Rivera, for the patience and meticulousness in preparing the text for publication. His suggestions and guidance helped to make this book better.

I would like to thank my technical editor and my friend, Boris M. Trukhnov, for the thorough technical editing of the book and his illuminating insights into the subject.

I would like to thank Robert M. Manning for helping with SQuirreL Universal SQL Client introduction (Appendix D) and to the entire SQuirreL development project team for the work that went into delivering this great free open source application.

My thanks go to Dzmitry Aliaksandrau, CCNA, for preparing screenshots for the database products used in the book and help in putting together the presentations. I'd like to thank Andrey Pfliger for help with testing SQL scripts in the book and suggestions on how to make the content more accessible for the readers.

CONTENTS

INTRODUCTION

xxv

CHAPTER 1: DROWNING IN DATA, DYING OF THIRST FOR KNOWLEDGE	1
Data Deluge and Informational Overload	2
Database Management Systems (DBMSs)	2
Storage Capacity	2
Number of Users	2
Security	2
Performance	3
Scalability	3
Costs	3
Recording Data	3
Oral Records	3
Pictures	4
Written Records	4
Printed Word	4
All of the Above	4
Analog versus Digital Data	4
To Store or Not to Store?	5
Relational Database Management Systems	6
IBM DB2 LUW	6
Oracle	7
Microsoft SQL Server	7
Microsoft Access	7
PostgreSQL	8
MySQL	8
HSQLDB and OpenOffice BASE	9
What Is SQL?	9
The SQL Standard	10
Dialects of SQL	10
Not the Only Game in Town	11
Let There Be Database!	11
Creating a Table	13
Getting the Data In: INSERT Statement	14
Give Me the World: SELECT Statement	16

Good Riddance: the DELETE Statement	22
I Can Fix That: the UPDATE Statement	25
Summary	28

CHAPTER 2: BREAKING AND ENTERING: STRUCTURED INFORMATION	29
---	-----------

A Really Brief Introduction to Data Modeling	29
Conceptual Modeling	30
Logical Modeling	30
Physical Modeling	31
Why Can't Everything Be Text?	31
Character Data	32
Fixed Length and Variable Strings	32
Binary Strings	34
Character versus Special Files	35
Numeric Data	36
Exact Numbers	36
Approximate Numbers	38
Literals for the Number	39
Once Upon a Time: Date and Time Data Types	40
Binary Data	42
It's a Bird, It's a Plane, It's ... a NULL!	43
Much Ado About Nothing	43
None of the Above: More Data Types	46
BOOLEAN	46
BIT	46
XML Data Type	46
DDL, DML, and DQL: Components of SQL	47
Refactoring Database TABLE	47
DROP TABLE	48
CREATE TABLE	48
ALTER TABLE	49
Populating a Table with Different Data Types	52
Implicit and Explicit Data Conversion	53
SELECT Statement Revisited	55
Selecting Literals, Functions, and Calculated Columns	55
Setting Vertical Limits	56
Alias: What's in a Name?	56
Setting Horizontal Limits	58
DISTINCT	58

Get Organized: Marching Orders	59
ORDER BY	59
ASC and DESC	60
TOP and LIMIT	60
INSERT, UPDATE, and DELETE Revisited	61
INSERT	61
SELECT INTO	63
UPDATE	63
DELETE	65
TRUNCATE That Table!	66
SQL Operators: Agents of Change	67
Arithmetic and String Concatenation Operators	67
Comparison Operators	68
Logical Operators	69
ALL	70
ANY SOME	70
BETWEEN <EXPRESSION> AND <EXPRESSION>	70
IN	71
EXISTS	72
LIKE	72
AND	74
NOT	75
OR	75
Assignment Operator	76
Bitwise Operators	76
Operator Precedence	77
Summary	78
 CHAPTER 3: A THING YOU CAN RELATE TO — DESIGNING A RELATIONAL DATABASE	 79
Entities and Attributes Revisited	80
Keys to the Kingdom: Primary and Foreign	81
Relationship Patterns	83
Domain Integrity	87
Am I Normal? Basics of Relational Database Design	89
Specifying Constraints	92
Selecting a Flavor For Your Data Model	93
Data Warehouses and Data Marts	93
Star and Snowflake Schemas	94
What Could and Does Go Wrong	94

Working with Multiple Tables	95
JOIN Syntax	95
UNION Operator	96
Dynamic SQL	97
Ultimate Flexibility, Potential Problems	99
Summary	101
CHAPTER 4: OVERCOMING THE LIMITATIONS OF SQL	103
In Numbers, Strength	104
Building Character	107
“X” Marks the Spot: Finding the Position of a Character in a String	112
CHARINDEX	113
CHAR	113
SUBSTRING	114
LENGTH	114
TRIM, LTRIM, and RTRIM	116
Date and Time Functions	117
What Time Is It?	117
Date Arithmetic	118
A Glimpse of Aggregate Functions	121
Conversion Functions	123
Conversion Between Different Data Types	125
Conversion Between Different Character Sets	125
Miscellaneous Functions	126
Making the CASE	127
SQL Procedural Extensions	129
Happy Parsing: Stored Procedures	131
User-Defined Functions (UDFs)	132
Why Use Procedural Extensions?	134
Performance and Network Traffic	134
Database Security	134
Code Reusability	135
Summary	135
CHAPTER 5: GROUPING AND AGGREGATION	137
Aggregate SQL Functions Revisited	137
AVG()	137
COUNT()	139
MAX()	140
MIN()	141
SUM()	142

Eliminating Duplicate Data	143
GROUP BY: Where Your Data Belongs	144
GROUP BY with HAVING Clause	148
ORDER BY Clause: Sorting Query Output	149
Summary	153
 CHAPTER 6: WHEN ONE IS NOT ENOUGH: A QUERY WITHIN A QUERY	 155
 What You Don't Know Might Help You	 155
Subquery in the WHERE Clause	155
EXISTS Operator	156
ANY Operator	157
ALL Operator	157
Subquery in the SELECT List	158
Subquery in the FROM Clause	160
Subquery in the HAVING Clause	161
Subqueries with INSERT	163
Subqueries with UPDATE	165
Subqueries with DELETE	166
Correlated Query	167
How Deep the Rabbit Hole Goes: Nesting Subqueries	169
A Subquery or a JOIN?	170
Summary	171
 CHAPTER 7: YOU BROKE IT; YOU FIX IT: COMBINING DATA SETS	 173
 Joins Revisited	 173
INNER JOIN	175
N-way INNER JOIN	179
LEFT OUTER JOIN	182
RIGHT OUTER JOIN	184
FULL JOIN	185
Self JOIN: Looking Inside for an Answer	186
CROSS JOIN (aka Cartesian Product)	187
State of the UNION	189
A Point of VIEW	193
CREATE VIEW	194
ALTER VIEW	198
DROP VIEW	198
Updatable VIEW	198
WITH CHECK OPTION	200

Hierarchical Views	201
Benefits and Drawbacks	202
But Wait; There's More!	203
INTERSECT	203
EXCEPT and MINUS	204
Summary	205
CHAPTER 8: WHAT ELSE IS THERE, AND WHY?	207
An INDEX for All Seasons	207
UNIQUE Index	209
CLUSTERED Index	209
An INDEX Destroyed	211
TABLE Revisited	211
VIEW Revisited	214
By Any Other Name: Aliases and Synonyms	214
Auto-Incremented Values	216
Identity Columns	217
Microsoft SQL Server	218
IBM DB2	220
PostgreSQL	221
MySQL	221
Microsoft Access	222
OpenOffice BASE with HSQLDB	222
Who Am I: Finding One's IDENTITY	223
Sequences	224
Comparing Identity Columns and Sequences	227
Triggers	228
One Happy Family: Working in Heterogeneous Environments	229
Summary	229
CHAPTER 9: OPTIMIZING PERFORMANCE	231
Database Performance	231
Performance Benchmarks	231
Order of Optimization	233
Hardware Optimization	234
Operating System Tune-up	234
Optimizing RDBMSs	234
Optimizing Database/Schema	234
Application Optimization	236
SQL Optimization	237

RDBMS-Specific Optimization	243
Oracle 10/11g	244
IBM DB2 LUW 9.7	244
Microsoft SQL Server 2008	245
PostgreSQL	245
MySQL	246
Desktop RDBMSs	247
Microsoft Access	247
OpenOffice BASE with HSQLDB Backend	248
Your DBA Is Your Friend	249
Summary	249
 CHAPTER 10: MULTIUSER ENVIRONMENT	 251
Sessions	251
Orphaned Sessions	254
Transactions	254
Understanding Locks	262
SQL Security	264
Basic Security Mechanisms	265
Defining a Database User	266
Managing Security with Privileges	268
Operating System Security Integration	272
INFORMATION_SCHEMA and SQL System Catalogs	279
Oracle Data Dictionary	281
IBM DB2 LUW System Catalogs	282
Microsoft SQL Server 2008 System Catalog	283
Summary	285
 CHAPTER 11: WORKING WITH UNSTRUCTURED AND SEMISTRUCTURED DATA	 287
SQL and XML	287
A Brief Introduction to XML	289
Formatted XML	290
DTD and Schema	290
Document Type Definition (DTD)	291
XML Schema Definition (XSD)	291
Namespaces	292
XML as a DataSource	294
Accessing XML Documents in an Application	294
XML Path Language: XPath	294

XML Query Language: XQuery	294
Encoding XML	294
Presenting XML Documents	296
XSL and XSLT	296
XML and RDBMSs	296
Implementation Details	299
Oracle 11g XML DB	302
IBM DB 9.7 pureXML	307
Microsoft SQL Server	311
PostgreSQL 9.0	316
MySQL 5.5	317
XML for RDBMS: Best Practices	318
All Bits Considered	320
What Would Google Do?	320
Getting Binary Data In and Out of the RDBMS Table	323
Best Practices for Binary Data	325
SQL and Text Documents	326
Summary	327
CHAPTER 12: NOT BY SQL ALONE	329
The Future Is Cloudy	329
Key/Value Pair	331
What in the World Is Hadoop?	334
Google's BigTable, Base, and Fusion Tables	334
Amazon SimpleDB	336
MongoDB	337
Microsoft SQL Azure	338
SQL and Business Intelligence	339
OLAP Rules	340
ROLAP, MOLAP, and HOLAP	341
Oracle 11g	342
IBM DB2	342
Microsoft SQL Server	343
XML for Analysis (XMLA)	344
Elementary, My Dear Watson!	344
Column-Oriented DBMS	345
Object Databases	346
Object-Oriented Programming (OOP) Paradigm	346
Objects and Classes	346

Object-Relational Mapping Frameworks	349
Hibernate/NHibernate	350
Microsoft LINQ and Entity Framework	350
Summary	350
APPENDIX A: INSTALLING THE LIBRARY DATABASE	353
Oracle 10g XE	354
Installing Library Sample Database with SQL*Plus	354
Installing with Oracle Web Interface	356
IBM DB2 9.7 Express-C	360
IBM Command Editor	360
IBM Command Window	362
Microsoft SQL Server 2008 Express	363
SQL Server Management Studio Express	363
PostgreSQL 9.0	365
Installing with pgAdmin III	366
MySQL 5.1	369
Installing with the MySQL CommandA-Line Utility	370
Microsoft Access 2007/2010	371
OpenOffice BASE 3.2	372
APPENDIX B: INSTALLING RDBMSS SOFTWARE	375
APPENDIX C: ACCESSING RDBMSS	377
Oracle	377
IBM DB2	377
Microsoft SQL Server 2008	377
MySQL	378
PostgreSQL	378
Microsoft Access 2007/2010	378
Open Office BASE with HSQLDB	378
APPENDIX D: ACCESSING RDBMSS WITH THE SQUIRREL UNIVERSAL SQL CLIENT	379
INDEX	381

INTRODUCTION

THE INFORMATIONAL DELUGE shows no signs of abating. We are inundated with data from the TV, from the Internet, and from advertisements stuffed in our mail boxes, virtual and otherwise. Unfortunately, as the quantity of information increased, its quality declined dramatically: Books were replaced by journals; then magazines; then newspapers; then web pages, blogs, and finally, tweets. The information becomes ever-more voluminous and ever-less trustworthy. Even worse, in the age of the Internet data never really disappears; it keeps accumulating, tucked away in files, logs, and databases. According to Google's former CEO Eric Schmidt, we create as much data in two days as we did from the first written record until 2003 (a date as good as any); this is about five *exabytes* (that is five billion gigabytes!) of data in just two days, and the pace keeps accelerating.

When electronic data storage became a reality, it brought about its own set of rules: To make sense out of the data, one had to learn the language. Relational database theory was so far the most successful attempt to bring electronic data under control, and it brought Structured Query Language (SQL) to go along with it.

The relational databases and SQL have evolved quite a bit since the 1970s when they made their first appearance, and the concepts embedded into the database SQL might appear counterintuitive to the uninitiated. By unraveling the SQL story, the reader will understand the rationale behind it and will learn to appreciate both the power and the limitations of SQL.

WHO THIS BOOK IS FOR

This book starts at the beginning, and no prior knowledge of SQL or relational databases is assumed. Along the way, on a voyage of discovery, you will participate in the creation of the sample database, which not only incorporates all SQL concepts taught in the book but also undergoes several refactoring iterations to introduce data modeling, query tuning and optimization, and set of best practices for everything SQL.

This book is for computer programmers ready to add relational database programming to their skill sets, for the business users who want more power over the data locked away in their databases, and everybody else who might be interested in learning the powerful language, the lingua franca of the relational databases.

Readers with previous database experience might want to skim through the first couple of chapters and delve into more advanced topics, or they might decide to revisit the first principles introduced in these chapters; the choice is yours.

WHAT THIS BOOK COVERS

The book covers the current release of the SQL Standard, SQL:2008, but it mostly focuses on the practical side of the language, highlighting the differences between particular implementations. It provides examples using SQL implementations in the latest versions of the following modern database systems either available for download as free express editions, or as free open source software. The most popular desktop database packages, Microsoft Access and OpenOffice, are also covered:

- IBM UDB2 9.7
- Oracle 10g
- Microsoft SQL Server 2008/2005/2000
- MySQL 5.1/5.5
- PostgreSQL 9.0
- Microsoft Access 2007/2010
- OpenOffice 3.2 BASE (with embedded HSQLDB)

HOW THIS BOOK IS STRUCTURED

The book takes a holistic approach, introducing the reader to the concepts of the relational databases in general, and SQL in particular, by gradually building an understanding of the subject through the iterative process of refactoring the ideas, where each concept introduced at the beginning will be revisited in greater detail later on, illuminating the interconnectedness of the underlying principles.

Chapter 1 introduces the story behind SQL and the relational theory behind it. It is a whirlwind tour in which the basic concepts are introduced; all further chapters build upon it. The distinction between data and information is illuminated, and foundations are laid for further exploration. The chapter gives an overview of the relational database management systems (RDBMSs) used in this book.

We revisit these concepts again in Chapter 2 and add some more. The amorphous data becomes structured as it is being analyzed and conformed to the relational model. The “fridge magnets” paradigm becomes the “chest drawer” one, and then morphs into a bona fide relational database table.

The relational model is further explored in Chapter 3, as we step through the basics of the database design and normalization process. The SQL tools for working with normalized data are introduced. Dynamic SQL makes its appearance in this chapter.

To highlight both the power and limitations of SQL as a set-based language, some of the most popular procedural extensions (such as Oracle’s PL/SQL and Microsoft’s Transact-SQL) are discussed in Chapter 4. This chapter will also introduce SQL functions as a means of alleviating innate deficiencies of the language when dealing with a record-based logic.

Aggregate data are explored in Chapter 5, summarizing the power and limitations of the approach. The aggregate SQL functions introduced in the previous chapter are taken to the next level to show how SQL works with data stripped of its individuality.

Chapter 6 deals with subqueries when data sets are being staggered, and data discovery is based upon multilevel data filtering, one query providing selection criteria to another. The subqueries are precursors to the more SQL attuned JOIN(s), a recurring theme throughout the book.

The power of SQL comes from its ability to deal with data locked in relational tables. Chapter 7 explores the ways SQL can combine this data into a single data set.

This book introduces basic SQL concepts, opening the door for further exploration, and Chapter 8 lays out the next steps of this voyage, with concepts you might consider to explore further later on.

Chapter 9 deals with performance optimization, describing general approach and best practices in optimizing your queries and database environment.

Chapter 10 discusses how relational databases work in multiuser environments, and what mechanisms were implemented in SQL to deal with concurrent data access.

SQL is all about structure and order — it is Structured Query Language, after all! But the real data comes in every shape and size, and Chapter 11 shows how SQL accommodates semistructured (XML documents), unstructured (text files), and binary (such as pictures and sounds) data.

Chapter 12 briefly discusses the latest developments, such as columnar databases, NoSQL databases, object databases, and service oriented architecture (SOA), and how they relate to SQL.

Appendix A describes, step by step, the procedure for installing the sample Library database and populating it with an initial set of data with specific instructions for each of the seven databases discussed in this book. The SQL scripts for this are available for download from the book's support-
ing websites.

Appendix B provides step-by-step instructions for installing relational database software packages used in this book.

Appendix C describes facilities provided with each of the respective databases to access, create database objects, and manipulate data stored in the tables.

Appendix D introduces the open source project Squirrel Universal SQL client that can be used to access every database used in this book via Java Database Connectivity (JDBC) interface. It describes, step by step, the process of setting up and configuring the software.

WHAT YOU NEED TO USE THIS BOOK

To make the most out of this book, we recommend downloading and installing the relational database software used throughout the book. Most of the software is free or available on a free trial basis. You'll find step-by-step instructions in Appendix B.

CONVENTIONS

To help you get the most from the text and keep track of what's happening, we used a number of conventions throughout the book.

TRY IT OUT

Try It Out is an exercise you should work through, following the text in the book.

1. It usually consists of a set of steps.
2. Each step has a number.
3. Follow the steps through with your copy of the database.

How It Works

After each *Try It Out*, the code you typed will be explained in detail.



Boxes with a warning icon like this one hold important, not-to-be-forgotten information that is directly relevant to the surrounding text.



The pencil icon indicates notes, tips, hints, tricks, or asides to the current discussion.

As for styles in the text:

- We *highlight* new terms and important words when we introduce them.
- We show keyboard strokes like this: Ctrl+A.
- We show file names, URLs, and code within the text like so: `INSERT INTO...SELECT FROM`.

We present code in two different ways:

We use a monofont type with no highlighting for most code examples.

We use **bold** to emphasize code that is particularly important in the present context or to show changes from a previous code snippet.

SUPPORTING WEBSITES AND CODE

As you work through each chapter, we recommend that you download the SQL scripts to create and populate the database. The code is available at www.wrox.com or at www.agilitator.com. You can use the search box on the website to locate this title. After you have located this book, click the Download Code link to access the files that can be downloaded. You can download the files via HTTP or FTP. All the files are stored as ZIP files.



The ISBN for this book is 978-1-118-00267-4. You may find it easier to search by the ISBN than by the title of the book.

You can also download the code from the main WROX download page: www.wrox.com/dynamic/books/download.aspx. Click the link to the Discovering SQL: A Hands-On Guide for Beginners to access the files that can be downloaded.

ERRATA

We make every effort to ensure that there are no errors in the text or in the code. However, no one is perfect, and mistakes do occur. If you find an error in one of our books, like a spelling mistake or a faulty piece of code, we would be very grateful for your feedback. By sending in errata, you may save another reader hours of frustration, and at the same time, you will be helping us provide even higher quality information.

To find the errata page for this book, go to www.wrox.com and locate the title using the Search box or one of the title lists. Then, on the book details page, click the Book Errata link. On this page, you can view all errata that have been submitted for this book and posted by Wrox editors. You may also contact the author via e-mail at discovery@agilitator.com. A complete book list, including links to each book's errata, is also available at www.wrox.com/misc-pages/booklist.shtml.

If you don't spot "your" error on the Book Errata page, go to www.wrox.com/contact/techsupport.shtml and complete the form there to send us the error you have found. We'll check the information and, if appropriate, post a message to the book's errata page and fix the problem in subsequent editions of the book.

P2P.WROX.COM

For author and peer discussions, join the P2P forums at p2p.wrox.com. The forums are a Web-based system for you to post messages relating to Wrox books and related technologies, and interact with other readers and technology users. The forums offer a subscription feature to e-mail you

topics of interest of your choosing when new posts are made to the forums. Wrox authors, editors, other industry experts, and your fellow readers are present on these forums.

At p2p.wrox.com, you will find a number of different forums that will help you, not only as you read this book but also as you develop your own applications. To join the forums, just follow these steps:

1. Go to p2p.wrox.com and click the Register link.
2. Read the terms of use and click Agree.
3. Complete the required information to join, as well as any optional information you want to provide, and click Submit.
4. You will receive an e-mail with information describing how to verify your account and complete the joining process.



You can read messages in the forums without joining P2P, but in order to post your own messages, you must join.

Once you join, you can post new messages and respond to messages that other users post. You can read messages at any time on the Web. If you want to have new messages from a particular forum e-mailed to you, click the Subscribe to this Forum icon by the forum name in the forum listing.

For more information about how to use the Wrox P2P, be sure to read the P2P FAQs for answers to questions about how the forum software works, as well as many common questions specific to P2P and Wrox books. To read the FAQs, click the FAQ link on any P2P page.

Discovering SQL

1

Drowning in Data, Dying of Thirst for Knowledge

Information may be the most valuable commodity in the modern world. It can take many different forms: accounting and payroll information, information about customers and orders, scientific and statistical data, graphics, and multimedia, to mention just a few. We are virtually swamped with data, and we cannot (or at least we'd like to think about it this way) afford to lose it. As a society, we produce and consume ever increasing amounts of information, and database management systems were created to help us cope with informational deluge. These days we simply have too much data to keep storing it in file cabinets or cardboard boxes, and the data might come in all shapes and colors (figuratively speaking). The need to store large collections of persistent data safely, and “slice and dice” it efficiently, from different angles, by multiple users, and update it easily when necessary, is critical for every enterprise.

Besides storing the information, which is what electronic files are for, we need to be able to find it when needed and to filter out what is unnecessary and redundant. With the informational deluge brought about by Internet *findability*, the data formats have exploded, and most data comes unstructured: pictures, sounds, text, and so on. The approach that served us for decades — shredding data according to some predefined taxonomy — gave in to the greater flexibility of unstructured and semistructured data, and all this can still fit under the umbrella of a database (a broader concept than the “data banks” of the 1970s).

The databases evolved to accommodate all this, and their language, which was designed to work with characters and numbers, evolved along with it. The concept of gathering and organizing data in a database replaced with the concept of a data hub (“I might not have it, but I know where to find it”) with your data at the core, surrounded with ever less related (and less reliable) data at the rim.

When does data transform into information? When it is organized and is given a context. Raw data collection does not give you much. For example, the number 110110 could be a decimal number 54 in binary representation; November 1, 2010, the date of D. Hamilton Jackson Memorial Day commemorating establishment of the first press in the U.S. Virgin Islands; House Committee Report #110 for the 110th U.S. Congress (2007–2008), you get the idea.

To transform data into information, you can aggregate the data, add context, cross-reference with other data, and so on. This is as far as databases can take you. The next step, transforming information into knowledge, normally requires human involvement.

DATA DELUGE AND INFORMATIONAL OVERLOAD

One of the reasons behind building a database of your information is to filter the information specific to your needs, to separate the wheat from the chaff. Anybody who uses Internet search engines such as Google or Bing can attest that results brought back are far from being unambiguous because the search engine tries to find the best matches in the sea of relevant, tangentially relevant, and absolutely irrelevant information. Your database is created to serve your unique needs: to track your sales, your employees, and your book collection. In doing so, it might reach out and get some additional information (for example, getting a book's information from Amazon.com), but it will be information specific to your particular needs.

Another important aspect of the database is security. How secure do you need your data? Can anybody see it and modify? Does it need to be protected from unauthorized access due to compliance requirements and simply common sense?

Database management systems, otherwise known as DBMSs, answer all these questions, and more.

Database Management Systems (DBMSs)

What makes a database management system a system? It's a package deal: You get managed storage for your data, security, scalability, and facilities to get data in and out, and more. These are things to keep in mind when selecting a DBMS. The following sections describe a few of the factors that you should consider.

Storage Capacity

Will the selected DBMS be sufficient for current and future needs? If you intend to store your favorite recipes or manage your home library, you might decide to use a desktop database such as Microsoft Access. When you need to store terabytes of information (for example, New York Stock Exchange financial transactions for the last 50 years), you should shop for an enterprise class DBMS such as Oracle, Microsoft SQL Server, or IBM DB2.

Number of Users

If you are the only user of your database, you might not need some of the features designed to accommodate concurrent data use in your database. The current version of Microsoft Access, for instance, supports up to 255 concurrent users (in practice, actual numbers will depend on many factors, including network, bandwidths, and processing power). And with advanced clustering technologies, there is theoretically no limit on the number of users in an enterprise DBMS such as Oracle.

Security

How secure do you want your data to be? You might not be overly concerned if your favorite recipes are stolen, but you'd want your banking or health information to be as secure as possible (and there

are regulations to mandate certain levels of protection for various kinds of data collected). One of the major differentiators between enterprise class DBMSs and their desktop counterparts is a robust, finely grained security implementation. A simple file that is a Microsoft Access database is more insecure than a server-based IBM DB2 installation with multiple levels of protection.

Performance

How fast does your database need to be? Can you wait minutes for the information to come back, or must you have a subsecond response, as in a stock trading platform? The answers tie into the question about concurrent users and also scalability. Some DBMSs are inherently slower than the others, and should not be deployed in environments they cannot handle.

Scalability

As Yogi Berra used to say, “Predictions are hard, especially about the future.” Databases must be able to accommodate changing business needs. While one cannot anticipate all the changes down the road, one could make an educated guess based upon likely scenarios and industry trends. Your business will change (growths, acquisitions), and your database needs will change with it. You can bet that your data will live longer than the database it lives in. The operating system might change (mainframe, UNIX/Linux, Windows); the programming environments might change (COBOL, C/C++, Java, .Net); regulations might change, but your data must endure, and not entirely for sentimental reasons.

Any of the modern enterprise DBMSs will get a decent score on any of these factors; ultimately, your business needs will dictate the technology choice. Expert advice will be needed for large production deployment, and qualified database administrators to keep your database in the best shape possible. Once you master the language, your data could be transformed into information; it will be up to you to take it to the next level: knowledge.

Costs

Of course, it is important to consider costs associated with installing and operating a database. Vendors might charge hundreds of thousands of dollars for an enterprise class DBMS or it could be had for free as an open source DBMS. Remember: “There ain’t no such thing as a free lunch.” An open source DBMS might save you money in upfront costs, but would quickly catch up in expertise, time, tools availability, and maintenance costs later on. The total cost of ownership (TCO) must be considered for every DBMS installation.

Recording Data

As far as recorded history goes, humans kept, well, records. Some philosophers even argue that one of the major differences between humans and animals is the ability to record (and recall) past events.

Oral Records

In all probability, oral records were the first kind of persistent storage that humans mastered. The information was transmitted from generation to generation through painstaking memorization; mnemonic techniques such as melody and rhyming were developed along the way. Information transmitted orally was highly storage-dependent, and could deteriorate (as in a game of Chinese

whispers) or disappear altogether after an unfortunate encounter by the bearer with a lion, a shark, or a grizzly bear.

Pictures

Pictures such as petroglyphs or cave paintings were much sturdier and somewhat less dependent on vagaries in an individual's fate. They were recorded on a variety of media: clay, stone, bark, skin; and some have survived to the modern age. Unfortunately, much of the context for these pictures was lost, and their interpretation became a guessing game for the archeologists.

Written Records

The beginning of written records, first pictographs and then hieroglyphs, dates back to around 3000 BC, when the Sumerians invented wedge-shaped writing on clay tablets, or cuneiform. This activity gradually evolved into a number of alphabets, each with its own writing system, some related, some autochthonous. This opened the door to storing textual information in pretty much the same form that we use even now. The medium for the writing records also improved over time: clay, papyrus, calf skin, silk, and paper.

Printed Word

Recording and disseminating the information was a painstakingly manual process. Each record had to be copied by hand, which severely limited access to information. The next step was to automate the process with printing. First came woodblock printing, with the earliest surviving example in China dating back to 220 AD. This sped up the process dramatically; a single woodblock could produce hundreds of copies with relatively little effort. The invention of movable type, first by the Chinese and Koreans (1040 and 1230, respectively) and then by Johannes Gutenberg in 15th century Europe, led to dramatically increased access to information through automated duplication. Still, single storage (book) could only be used by a single user (reader) at a time, and searching was a painstaking manual process, even with invention of indexing systems (a list of keywords linked to the pages where these keywords were used).

All of the Above

Technological advances made it possible to accumulate information in a variety of media (text, pictures, and sounds). Not until electronic data storage was developed did it become possible to store them all together and cross-reference them for later automated retrieval. The data had to be digitized first.

Analog versus Digital Data

Up until the invention of the first computers, most information was created and stored in human-readable format. Various mechanical systems were invented to facilitate storage and retrieval of the information, but the information itself remained analogous: print, painting, and recorded sound. Sounds recorded on LP disks are analog, and sounds recorded on CD are digital. The most dedicated audiophiles claim that a CD is but an approximation of the real sound (and they are correct), but most people do not notice the difference. One cannot deny the convenience afforded by a digital CD (or, better yet, an audio file stored on one's computer).

The idea to represent data in binary format came independently to several people around the world, with MIT engineer Claude Shannon formulating principles of binary computation in 1938, and German scientist Konrad Zuse creating a fully functional binary computer in 1941. It turns out that a binary system is uniquely suited for the electrical signal processing; it was humans' turn to adapt to a machine.

The familiar letters and punctuation were translated into combinations of ones and zeroes, starting with the Extended Binary Coded Decimal Interchange Code (EBCDIC), developed by IBM in the early 1950s; through the American Standard Code for Information Interchange (ASCII) character-encoding scheme introduced in the early 1960s; to the advent of Unicode, which made its debut in 1991. The latter system was designed to accommodate every writing system on Earth, and can currently represent 109,000 characters covering 93 distinct scripts.

While initial efforts were focused on representing characters and numbers, the other types of data were not far behind. Pictures and then sounds became digitized and eventually made their way into databases.

To Store or Not to Store?

In 1956, IBM was selling five megabyte persistent storage drives for a whopping \$10,000 per megabyte (no wonder it had to make this agonizing decision to store dates as two digits instead of four; also known as the Y2K problem); this came down to just under \$200 per megabyte in 1981 (Morrow Designs). In August 2010, a Western Digital 1 terabyte hard drive was selling for \$70, which translates into 122 megabytes per one cent!

When storage was dear, people had to be very selective about what data they wanted to keep; with costs plummeting, we've set our sights on capturing and storing *everything*.

The Holy Grail of the DBMS for years was to structure and organize data in a format that computers could manipulate; the preferred way was to collect the data and sort it, and then store it in bits and pieces into some sort of a database (it was called a *data bank* in those days, with policies to match). You had to own all your data. With the proliferation of the Internet, this is no longer the case. Distributed data is now the norm; instead of bringing the data in, you might choose to store information about where the data could be found and leave it at that.

Of course, you may need to keep some of your data closer to the vest (financial data and personal data, for example). Storing the actual data will give you full control of how this data is accessed and modified; this is what databases do best.

With all this dizzying variety of data formats, one needs to make a decision on how this data is to be stored. Despite advances in processing unstructured data, organizing it into taxonomies (a process called *data modeling*; see Chapters 2 and 3 for more information) has distinct benefits both in speed and flexibility. Breaking your data down into the smallest bits and pieces requires a lot of upfront effort, but it gives you an ability to use it in many more ways than when stored as monolithic blocks. Compare a Lego bricks castle with a premolded plastic castle. The latter stays a castle forever, while the former could be used to build a racing car model, if needed. The tradeoffs between structured and unstructured data (and everything in-between) will be discussed in Chapter 11.

Relational Database Management Systems

This book is about SQL, the language of relational databases, or relational database management systems (RDBMSs). Since the theoretical foundations were laid down in the 1970s by Dr. Codd, quite a few implementations have come into existence, and many more are yet to come.

Many people consider DB2 to be the granddaddy of all databases, given that the very term *relational* was introduced by IBM researcher Dr. Edgar Frank Codd in 1969, when he published his paper, “Derivability, Redundancy, and Consistency of Relations Stored in Large Data Banks” in an IBM research report. This assertion is contested by others who point to Oracle’s version 2 commercial release in 1979; Multics Relational Data Store sold by Honeywell Information Systems in 1976; or the Micro DBMS experimental designs (pioneering some of the principles formulated by Dr. Codd two years later) of the University of Michigan from 1968 (the last instance of Micro DBMS in production was decommissioned in 1998). The RDBMS road is marked by a multitude of milestones (and an occasional gravestone) of other RDBMS products, including IBM PRTV (1976); IBM SQL/DS (1980); QBE(1976); Informix (1986); Sybase (1986); Teradata (1979); and Ingres, an open source project that gave inspiration to many other successful systems such as PostgreSQL (1996), Nonstop SQL (1987), and Microsoft SQL Server (1988) — to mention but a few. These systems used different dialects of primordial SQL: SEQUEL, QUEL, Informix-SQL, and so on. It was not until 1987 when the first attempt was made to standardize the language; arguably, the battle is still going on.

The current RDBMS market is split among heavyweight proprietary relational databases Oracle (48 percent), IBM (25 percent), and Microsoft (18 percent); smaller proprietary systems Teradata and Sybase, each with a distant 2 percent; and the other vendors, as well as open source databases, comprising about 10 percent of the total market.

For a sizeable enterprise, selecting a database foundation for their applications is a decision not to be taken lightly. Not only does it cost tens of thousands of dollars in upfront licensing fees for the software, and hundreds of thousands of dollars in maintenance and support fees, but it is also an important factor in determining the overall enterprise architecture that aligns all other investments in software, hardware, and human resources. Although migrating from one RDBMS to another became easier in recent years, still the mere thought of it might give your CFO nightmares.

IBM DB2 LUW

IBM is a long-term front-runner in the RDBMS arena, from the mainframe world with the MVS family of operating systems, to z/OS, and later to UNIX and Windows. The current version is IBM DB2 9.7 LUW (Linux, UNIX, and Windows).

The IBM DB2 9.7 keeps the absolute record in transaction processing speed (see Chapter 9 for more information) and comes in a variety of editions, from Advanced Server Enterprise to a free (albeit limited) DB2 Express-C edition used to run samples provided with this book.

DB2 in its version 9.7 is still only compliant with the ANSI/ISO SQL 92 Entry standard (see later in this chapter) and supports some of the more advanced features from other standards organizations such as the Open Geospatial Consortium, JDBC, X/Open XA, as well as bits and pieces of the latest SQL:2008 Standard. In addition to its own built-in procedural extension language, SQL PL, it also provides support for Oracle’s PL/SQL, Java, and even Microsoft’s .NET family languages for creating stored procedures (see Chapter 4 for more information).

Oracle

Oracle traces its roots back to the first release of Oracle version 2 in 1979, initially for older VAX/VMS systems, with UNIX support following in 1983. Over the years, it added support for most of the features specified in SQL Standard, culminating in the latest release of Oracle 11g, which claims compliance with the “many features” of the latest release of SQL:2008 Standard.

Oracle holds second place in the high-performance transaction processing benchmarking and is at the center of the company’s ecosystem. It is a secure, robust, scalable, high-performance database that has dominated the UNIX market for decades. In addition to SQL support, it comes with a built-in procedural language, PL/SQL (see Chapter 4 for more information on procedural extensions), as well as support for general programming languages such as Java.

At the time of this writing, the latest version is Oracle 11g; the free express edition is available only for Oracle 10g, which has some limits on the data storage size and number of processors (CPUs) the RDBMS is capable of utilizing. The express edition has full support for all SQL features discussed in this book.

Microsoft SQL Server

SQL Server began as partnership between Sybase, Microsoft, and Ashton-Tate, with the initial idea to adapt existing UNIX-only Sybase SQL Server to then-new IBM operating system OS/2. Ashton-Tate later dropped out of the partnership, and the IBM OS/2 operating system faded into oblivion. Microsoft and Sybase were to share the world, being careful not to step on each other’s toes. Microsoft was to develop and support SQL Server on Windows and OS/2, and Sybase was to take over UNIX platforms. The partnership formally ended in 1994, although at its core, Microsoft SQL Server still used fair chunks of Sybase technology. In 1998, beginning with the release of Microsoft SQL Server 7.0, the last traces of Sybase legacy were eliminated, and a brand spanking new RDBMS set out to conquer the world (the Windows world, that is). As of today, Microsoft holds about 20 percent of the RDBMS market, though on Windows it reigns supreme.

The latest version as of this writing is Microsoft SQL Server 2008 Release 2; a limited Express edition available for free that supports all features of SQL covered here.

Microsoft Access

Microsoft Access, known lately as Microsoft Office Access, is a desktop relational database (relatively relational, as some might quip). It purports to be an integrated solution combining elements of a relational database engine, application development infrastructure (complete with built-in programming language and programming model), and reporting platform. Unlike other RDBMSs discussed in the book, this is a file-based database and as such has inherent limitations in performance and scalability. For example, while the latest version theoretically allows for up to 255 concurrent users, in practice anything more than a dozen users slows the performance to a crawl. It also supports only a subset of SQL Standard, as well as a number of features available in its own environment only.

One of the features is linking in tables from remote databases that allow it to be used as an application front end to any ODBC/OLEDB-compliant database.

PostgreSQL

PostgreSQL evolved from a project at the University of California at Berkeley lead by Michael Stonebraker, one of the pioneers of the relational databases theory. The principles that went into the original Ingres project, and its successor PostgreSQL, also found their way into many other RDBMS products such as Sybase, Informix, EnterpriseDB, and Greenplum.

The first version of PostgreSQL (with this exact name) came in 1996; it was released in version 6.0 the next year, and remained an open source project maintained by group of dedicated developers. There are numerous commercial versions of PostgreSQL; most notable is EnterpriseDB, a private company that offers enterprise support (along with variety of proprietary management tools) for the product and has convinced many high-profile customers such as Sony and Vonage to rely on an open source RDBMS for some critical enterprise class applications.

PostgreSQL is arguably the closest in terms of support for the SQL standards in addition to a number of features found nowhere else. Unlike its peers (such as MySQL), it provided referential integrity and transactional support from the beginning. It also comes with built-in support for the PL/pgSQL procedural extension language, as well as the capability to adapt virtually any other language to the same purpose.

MySQL

MySQL was first developed by Michael Widenius and David Axmark back in 1994, with its first release in 1995. It was initially positioned as a lightweight, fast database to serve as the back end for data-driven websites. Even though it was lacking many features of the more mature RDBMS products, it was fast in serving information and “good enough” for many scenarios. (To be really fast, MySQL can bypass referential integrity constraints and ditch transactional support; see Chapters 3 and 10 for additional information.) Plus you could not beat the price; it was free. No wonder it grew up to be the most popular relational database among small- and medium-sized users. There were a number of other free database products on the market that lacked features, near-commercial polish, or both. Not one of the big guys — Oracle, IBM, Microsoft, and Sybase — offered free express versions of their respective RDBMSs back then. MySQL was acquired by Sun Microsystems in 2008, which was subsequently swallowed by Oracle.

Currently, Oracle offers a commercially supported version of MySQL as well as a Community Edition. Following this acquisition, a number of fork versions sprang up, such as MariaDB and Percona Server, committed to maintain free status under the General Public License (GPL), one of the least restrictive open source licenses.

The latest released version of MySQL is 5.5, with version 6 on the horizon. It is multiplatform (Linux/UNIX/Windows), and supports most of the features of SQL:1999; some of the features depend on the selected options (for example, a storage engine).



The storage engine option is a feature unique to MySQL, which allows handling of different table types differently. Each engine comes with unique capabilities and limitations (transactional support, index clustering, storage limits, and so on). A database table could be created with different storage engine options, with the default being MyISAM engine.

HSQldb and OpenOffice BASE

Hyper Structured Query Language Database (HSQldb), a relational database management system implemented in the Java programming language, is available as open source under the Berkley Software Distribution (BSD) license (meaning pretty much free for all).

This is a default RDBMS engine shipped with the OpenOffice.org BASE, a desktop database positioned to compete in the same market as Microsoft Access. It is a relational database, robust, versatile, and reasonably fast, and is supported on multiple platforms including Linux, various flavors of UNIX, and Microsoft Windows. It claims to be almost fully compliant with SQL:1992 Standard, which covers most of the SQL subset discussed in this book.

An adaptation of HSQldb serves as an embedded back end to the OpenOffice.org suite component BASE and became part of the suite starting with version 2.0. Like Microsoft Access, the OpenOffice BASE can connect to a variety of RDBMSs, provided that there is a suitable driver; a number of Java Database Connectivity (JDBC) and ODBC (Open Database Connectivity) drivers are available and ship with the product.



Following Oracle's acquisition of OpenOffice and its uncertain status as an open source project under Oracle's patronage, the OpenOffice.org community decided to start a new project called LibreOffice, with the intent of implementing all the functionality of OpenOffice as free software under the original BSD license.

Relational databases are not the only game in town. Some of the older technologies, seemingly forever defeated by relational database theory, came back, helped by ever faster/cheaper hardware and software innovations. The quest for better performance and ease of creating applications spawned research into columnar and object-oriented databases, frameworks that make the “all data in one bucket” approach workable, domain-specific extensions (such as geodetic data management or multimedia), and various data access mechanisms. We discuss this topic in Chapter 12.

WHAT IS SQL?

Before the advent of commercially available databases, every system in need of persistent storage had no choice but to implement its own, usually in some proprietary file format (binary or text) that only this application could read from and write to. This required every application that used these files to be intimately familiar with the structure of the file, which made switching to a different storage all but impossible. Additionally, you had to learn a vendor-specific access mechanism to be able to use it. Relational model dealt with complexities of data structures, organizing data on logical level, but it had nothing to say about the specifics of storage and retrieval except that it had to be set-based and follow relational algebra rules. Left to their own devices, the early RDBMSs implemented a number of languages, including SEQUEL, developed by Donald D. Chamberlin and Raymond F. Boyce in the early 1970s while working at IBM; and QUEL, the original language of Ingres. Eventually these efforts converged into a workable SQL, the Structured Query Language.

SQL is a RDBMS programming language designed to define relational constructs (such as schemas and tables) and provide data manipulation capabilities. Unlike many programming languages in

general use, it does not exist outside the relational model. It cannot create stand-alone programs; it can only be used inside RDBMSs. This is a declarative type of language. It instructs the database about what you want to do, and leaves details of implementation (*how* to do it) to the RDBMS itself. In Chapter 2, we will go over the elements of the language in detail.

From the very beginning there were different dialects bearing the same SQL name, some of them quite different from each other. This worked for the vendors, as it assured lock-in to specific technology, but it also defied the purpose of creating SQL in the first place.

The SQL Standard

To bring greater conformity among vendors, the American National Standards Institute (ANSI) published its first SQL Standard in 1986 and a second widely adopted standard in 1989. ANSI released updates in 1992, known as SQL92 and SQL2, and again in 1999: SQL99 and SQL3. Each time, ANSI added new features and incorporated new commands and capabilities into the language.

The ANSI standards formalized many SQL behaviors and syntax structures across a variety of products. These standards become even more important as open source database products (such as MySQL, mSQL, and PostgreSQL) grow in popularity and are developed by virtual teams rather than large corporations.

The SQL Standard is now maintained by both ANSI and the International Standards Organization (ISO) as ISO/IEC 9075 standard. The latest released standard is SQL:2008, and work is underway to release the next version of the standard to accommodate new developments in the way RDBMSs collect and disseminate data.

Dialects of SQL

Even with a standard in place, the constantly evolving nature of the SQL Standard has given rise to a number of SQL dialects among the various vendors and products. These dialects most commonly evolved because the user community of a given database vendor required capabilities in the database before the ANSI committee created a standard. Occasionally, though, a new feature is introduced by the academic or research communities due to competitive pressures from competing technologies. For example, many database vendors are augmenting their current programmatic offerings with Java (as is the case with Oracle and Sybase) or .Net (Microsoft's SQL Server Integration Services, embedded common language runtime [CLR]).

Nonetheless, each of these procedural dialects includes conditional processing (such as that controlled through IF ... THEN statements), control-of-flow functions (such as WHILE loops), variables, and error handling. Because ANSI had not yet developed a standard for these important features at the time, RDBMS developers and vendors were free to create their own commands and syntax. In fact, some of the earliest vendors from the 1980s have variances in the most fundamental language elements, such as SELECT, because their implementations predate the standards. Some popular dialects of SQL include the following:

- **PL/SQL** — Found in Oracle. PL/SQL, which stands for Procedural Language/SQL and contains many similarities to the general programming language Ada; IBM DB2 added (limited) support for Oracle's PL/SQL in version 9.5.

- **Transact-SQL** — Used by both Microsoft SQL Server and Sybase Adaptive Server. As Microsoft and Sybase have moved away from the common platform they shared early in the 1990s, their implementations of Transact-SQL have also diverged, producing two distinct dialects of Transact-SQL.
- **SQL PL** — IBM DB2's procedural extension for SQL, introduced in version 7.0, provides constructs necessary for implementing control flow logic around traditional SQL queries and operations.
- **PL/pgSQL** — The name of the SQL dialect and extensions implemented in PostgreSQL. The acronym stands for Procedural Language/postgreSQL.
- **MySQL** — MySQL has introduced a procedural language into its database in version 5, but there is no official name for it. It is conceivable that with Oracle's acquiring the RDBMS it might introduce PL/SQL as part of the MySQL.

Not the Only Game in Town

Over the years there were many efforts to improve upon SQL and extend it beyond original purpose. With the advent of object-oriented programming, there came demand to store objects in the database; proliferation of Internet and multimedia increased demand for storage, indexing and retrieval of the binary information and XML data, and so on. While SQL standards were keeping pace with these and other demands, some decided to create a better mousetrap and came up with some ingenious ideas. For instance, HTSQL is a language that allows you to query data over Internet HTTP protocol; Datalog was envisioned as a data equivalent of Prolog, an artificial intelligence language; and MUMPS (going back to the 1960s!) mixes and matches procedural and data access elements.

The latest entry came from the NoSQL family of databases that depart from conventional relational database theory and eerily reminds us of a data bucket with key/value indexed storage. We will have a brief discussion about evolution of SQL in the last chapter of this book.

LET THERE BE DATABASE!

There is a bit of groundwork to be performed before we could submit our SQL statements to RDBMSs. If you have followed the instructions in Appendix B, complemented by the presentation slides on the accompanying book sites (both at www.wrox.com and at www.agilitator.com), you should have an up-and-running one (or all) of the RDBMSs used in this book; alternatively, you should have Microsoft Access or OpenOffice BASE installed. Please refer to Appendix B for step by step installation procedures for the RDBMS, and to Appendix A for instructions on how to install the Library sample database.



The following, with minor modifications, will work in server RDBMSs: Oracle, IBM DB2, Microsoft SQL Server, PostgreSQL, and MySQL. In Microsoft Access and OpenOffice BASE/HSQLDB, you'd need to create a project.

The concept of a database, a logically confined data storage (exemplified by the now rarely used term *data bank*), managed by a program is rather intuitive. When using a desktop database such

as Microsoft Access, your database is a file that Access creates for every new project you start; the server-based RDBMSs use a similar concept, though the details of implementation are much more complex. Fortunately, the declarative nature of SQL hides this complexity. It tells what needs to be done, not how to do it.

In the beginning, there was a database. The database we will use throughout the book will contain all the books we have on the shelves; a book tracking database that stores titles, ISBN numbers, authors, price, and so on — quite helpful in figuring out what you have.

The following statement creates a database named LIBRARY in your RDBMS (as long as it is Microsoft SQL Server, IBM DB2, PostgreSQL and MySQL; things are a bit different with Oracle, which subscribes to a different notion of what is considered a database; see Appendix A for more details).

```
CREATE DATABASE library;
```

If you have sufficient privileges in the RDBMS instance, the preceding statement will create a database, a logical structure to hold your data, along with all supporting structures, files, and multitudes of other objects necessary for its operations. You need to know nothing about these; all the blanks are filled with default values. Behold the power of a declarative language!



Oracle's syntax would be similar to this:

```
CREATE USER library IDENTIFIED BY discover;
```

With USER being roughly an equivalent of the DATABASE in other RDBMS. A discussion of the similarities and differences between the two are outside scope of this book.

Of course, there is much more to creating a database that would adequately perform in a production environment; there are a myriad of options and tradeoffs to be considered, but the basic data storage will be created and made available to you with these three words.

Once created, a database can be destroyed just as easily, using SQL's DROP statement; you cannot destroy objects that do not exist (and the RDBMS will warn you about it should you attempt to):

```
DROP DATABASE library;
```

In Oracle, of course, you'd be dropping a USER.

Now the database is gone from your server; in Microsoft Access and OpenOffice BASE, this is equivalent to deleting corresponding files.



Due to certain differences in terms of usage across RDBMSs, the concept of database is different among various proprietary databases. For example, what SQL Server defines as a database is in a way similar to both the SCHEMA and USER in Oracle, but in the context of this book, these differences are not particularly important.

Creating a Table

Now that we have a database, we can use it to create objects *in* the database, such as a table. A *table* is place where all your data will be stored, and this is where common sense logic and that of RDBMS begin to diverge.

If your refrigerator is anything like ours, you will have all kind of things held to its surface by magnets, some goofy keepsakes from a trip to a zoo, a calendar sent to you by your friendly insurance agent, your kid's school menu (and school attendance phone line), a shopping list, photos of your dog, photos of your children, the pizza hotline... Think of it as your personal database. You could just stick anything to it: text, pictures, calendars, and what not. In contrast, the RDBMS is much more particular. It will ask you to sort your data according to data types. A detailed discussion of data types will take place later, in Chapter 2. Here, we just stick to the data type most intuitively understood and best dealt with by the RDBMS: the text.

Creating a table is just as easy as creating the database in the previous example, with a minor difference of specifying a name for the table column and its data type:

```
CREATE TABLE myLibrary (all_my_books VARCHAR(4000));
```

The column ALL_MY_BOOKS is defined as a character data type (see Chapter 2 for more information of data types), and it can hold as many as 4,000 characters.



As you might have guessed, there is much more to the CREATE TABLE syntax than the preceding example implies. A full syntax listing all options in any given RDBMS would span more than one page, and mastering these options requires advanced understanding of SQL, for which this book is but a first step.

As you'll see in Chapter 2, a table, once created, can be modified (altered), or dropped from the database altogether. The SQL provides you with full control over the database objects, with power to create, change, and destroy.

TRY IT OUT Creating a Database in Microsoft SQL Server 2008

Creating a database is normally a database administrator's task, especially in a production environment; there are too many options and tradeoffs to consider to leave everything set to the default. For our purposes, we can use the basic syntax, however. There are several ways to create a database in Microsoft SQL Server, and using SQL Server Management Studio Express is arguably the easiest one. Follow these steps:

1. Make sure that you have your SQL Server instance up and running (refer to Appendix B for installation instructions).
2. Start SQL Server Management Studio Express by going to the Microsoft SQL Server 2008 menu option (this exercise assumes that SQL Server is installed on your local computer so you can connect automatically with Windows Authentication).

3. The first screen you see is a prompt to connect to your server. If not already filled by default, select the server type Database Engine, the server name .SQLEXPRESS (if you followed the instructions in Appendix B; otherwise, select another name from the drop-down box; it only displays instances of SQL Server visible from your computer), and authentication set to Windows Authentication.
4. Click the Connect button.
5. SQL Server Management Studio Express will display a window with several panes; for the purposes of this tutorial, we are only interested in the New Query button located in the upper-left corner of the window, right under the File menu (shown in Figure 1-1). Click the New Query button.
6. A new query window would appear in the middle of the window; this is where you will enter your SQL commands.
7. Type in the SQL statement for creating a database:

```
CREATE DATABASE library;
```



FIGURE 1-1

8. Click the Execute button located on the upper toolbar, as shown in Figure 1-2.
9. Observe the message “Command(s) completed successfully” in the lower pane, Messages tab.
10. Your newly created database will appear on the Databases list in the pane on the left, with the title Object Explorer (see Figure 1-3). Click the plus sign next to the node Databases node.



FIGURE 1-2

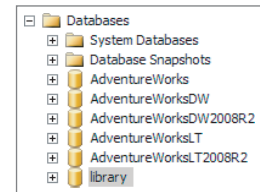


FIGURE 1-3

How It Works

Microsoft SQL Server takes out much of the complexity from creating the database process. Behind the scenes, the SQL Server created several files on the hard drive of your computer (or on an external storage device), created dozens of entries in the Windows registry and the SQL Server–specific configuration files, and created additional supporting objects for the database operations (you can take a look at these by expanding the node LIBRARY in your newly created database).

By omitting all optional configuration options, your database was created using all the default values: storage file names, locations, and initial sizes; collation orders; and so on. While this is not a recipe for creating an optimally performing database (see Chapter 9 for optimization considerations), it will be adequate for the purposes of this book.

Getting the Data In: INSERT Statement

The myLibrary table in our LIBRARY database is now ready to be populated with data, which is a task for the INSERT statement. Since the stated purpose of our database is to keep track of the books, let's insert some data using one of the books we do have on our shelf, *SQL Bible*. Here is some data.


```
SQL Bible by Alex Kriegel Boris M. Trukhnov Paperback: 888 pages
Publisher: Wiley; 2 edition (April 7, 2008) Language: English
ISBN-13: 978-0470229064
```

This is a lot of information and all in one long string of characters. The INSERT statement would look like follows:

```
INSERT INTO myLibrary VALUES ('SQL Bible by Alex Kriegel Boris M. Trukhnov
Paperback: 888 pages Publisher: Wiley; 2 edition (April 7,2008)
Language:English ISBN-13: 978-0470229064');
```

The keywords INSERT, INTO, and VALUES are the elements of the SQL language and together instruct the RDBMS to place the character data (in the parentheses, surrounded by single quotation marks) into the myLibrary table. Note that we did not indicate the column name; first because we have but a single column in which to insert, and second because RDBMS is smart enough to figure out what data goes where by matching a list of values to the implied list of columns. Both parentheses and quotation marks are absolutely necessary: the former signifies a list of data to be inserted, and the latter tells the RDBMS that it is dealing with text (character data type).

In database parlance, we have created a record in the table. There are many more books on the shelf, so how do we enter them? One way would be to add all of them on the same line, creating a huge single record. Although that is possible, within limits, it would be impractical, creating a pile of data not unlike the refrigerator model we discussed earlier: easy to add and difficult to find. Do I hear “multiple records”? Absolutely!

The previous statement could be repeated multiple times with different data until all books are entered into the table; creating a new record every time. Instead of a refrigerator model with all data all in one place, we moved onto “chest drawer model” with every book having a record of its own.

TRY IT OUT Inserting Data into a Column

Make sure you are at the step where you can enter and execute SQL commands. Repeat Steps 1 through 6 of the first Try It Out exercise and then run these statements to insert four records in your single table, single column database:

1. Type in (or download from a website) the following queries:

```
USE library;
INSERT INTO myLibrary VALUES ('SQL Bible by Alex Kriegel Boris M. Trukhnov
Paperback: 888 pages Publisher: Wiley; 2 edition (April 7,2008) Language:English
ISBN-13: 978-0470229064');
```

```
INSERT INTO myLibrary VALUES ('Microsoft SQL Server 2000 Weekend Crash Course by
Alex Kriegel Paperback: 408 pages Publisher: Wiley (October 15, 2001)
Language:English ISBN-13: 978-0764548406');
```

```
INSERT INTO myLibrary (all_my_books ) VALUES ('Letters From The Earth by Mark Twain
Paperback: 52 pages Publisher: Greenbook Publications, LLC (June 7, 2010)
Language:English ISBN-13: 978-1617430060');
```

```
INSERT INTO myLibrary (all_my_books ) VALUES ('Mindswap by Robert Sheckley
Paperback: 224 pages Publisher: Orb Books (May 30, 2006)
Language:English ISBN-13: 978-0765315601');
```

2. Click the Execute button located on the upper toolbar, as shown on Figure 1-2.
3. Observe four confirmations “(1 row(s) affected)” in the Messages tab in the lower window.

How It Works

The INSERT statement populates columns in the table by creating a record, a single row of data. The list of columns could be omitted as the list of values corresponds exactly to the list of columns (see later in this chapter for more information). If a column is specified, it has to appear in parentheses without any quotation marks; and the corresponding data goes into the list after the VALUES keyword, in parentheses, with quotation marks around the data to indicate the character nature of the value.

Give Me the World: SELECT Statement

Now that we have our data, we could query it to find out exactly what we have. The SELECT statement will help us to get the data out of the table; all we need is to tell it what table and what column.

```
SELECT all_my_books FROM myLibrary;
```

While it did produce a list of the books’ information, it is far cry from being useful. Let’s face it; it is a mess of a data, and the only advantage from being stored in a relational database is that it can be easily recalled, and possibly printed. What about search? To find out whether you have a specific book, you’d have to pull all the records and *manually* go over each and every one of them! Hardly a result you would expect from a sophisticated piece of software, which is RDBMS.

We need a way to address specific keywords in the records that we store in the table, such as the book title or ISBN number. A standard programming answer to this problem is to parse the record: chop it into pieces and scroll them in a loop looking for a specific one, repeating this process for each record in the table. The SQL cannot do any of this without vendor-specific procedural extensions. This would defy declarative nature of the language and would require intimate understanding of the data structure. Let’s take another look at the first record of data we entered:

```
SQL Bible by Alex Kriegel Boris M. Trukhnov Paperback: 888 pages  
Publisher: Wiley; 2 edition (April 7, 2008) Language: English  
ISBN-13: 978-0470229064
```

How would you go about chopping the record into chunks? What would be the markers for each, and how do you distinguish a book title from an author? Using a blank space for this purpose would put “SQL” and “Bible” into different buckets while they logically belong together. How do we know that “by” is a preposition, and not part of the author’s name? The answer comes from the structured nature of SQL, which is, after all, a *structured* query language; we need more columns. Splitting the one unwieldy string into semantically coherent data chunks would allow us to address each of them separately as each chunk becomes a column unto its own. Back to the CREATE TABLE (but let’s first drop the existing one):

```
DROP TABLE myLibrary;
```

Create a new one according to the epiphany we just had:

```
CREATE TABLE myLibrary
(
    title          VARCHAR(100)
  , author        VARCHAR(100)
  , author2       VARCHAR(100)
  , publisher      VARCHAR(100)
  , pages         INTEGER
  , publish_date   VARCHAR(100)
  , isbn          VARCHAR(100)
  , book_language  VARCHAR(100)
)
```

A single column became eight columns with an opportunity to add a ninth by splitting the authors' first and last names into separate columns (this is part of the data modeling process to be discussed in Chapter 3). For now, we've used the same data type, albeit shortened the number of characters, with a single exception: We made the PAGES column a number for reasons to be explained later in this chapter. You might also consider changing the data type of the column PUBLISH_DATE. Normally, a date behaves differently from a character, and the DBMS offers a date- and time-specific data type.

Now that we don't have to dump all data into the same bucket, we can be much more selective about data types, and use different types for different columns. It is not recommended that you mix up the data types when inserting or updating (see later in this chapter) the columns.

We will revisit data types again later in this chapter, and in more detail in Chapter 2.



You might have noticed that we have two “author” columns in our table now, to accommodate the fact that there are two authors. This raises the question of what to do when there is only one author, or when there are six of them. These questions will be explored in depth in a data modeling session in Chapters 2 and 3; here we just note that unused columns are populated automatically with default values, and if you find yourself needing to add columns to your table often, it might be the time to read about database normalization (see Chapter 3).

Now we need to populate our new table. The process is identical to the one described before, only the VALUES list will be longer as it will contain eight members instead of one. All supplied data must be in single quotes with the exception of the one going to PAGES column; quotes signify character data, absence thereof means numbers:

```
INSERT INTO myLibrary VALUES (
    'SQL Bible'
  , 'Alex Kriegel'
  , 'Boris M. Trukhnov'
  , 'Wiley'
  , 888
  , 'April 7,2008'
  , '978-0470229064'
  , 'English');
```