

# WELCOME TO ISSUE 2, 2011



## Already a subscriber?

Turn the page or click here for instructions.

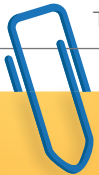
## Like what you're reading?

Share it with a colleague now!

## Not a subscriber?

Click here to sign up now and enjoy these exclusive digital-edition benefits:

- › Search through the entire issue
- › Access back-issue archives
- › Share articles with others
- › Download a PDF copy



This Digital Edition Sponsored By:

**IBM will  
tell your  
story.**

# GET NOTICED.

## Get started today.

Contact the IBM client references team at [cusref@us.ibm.com](mailto:cusref@us.ibm.com), and find out how we can help your company with:

- Case studies
- Speaking engagements
- Press releases
- Video testimonials
- Analyst interviews
- Advertising

Find more details at [www.ibm.com/ibm/clientreference/us/en/](http://www.ibm.com/ibm/clientreference/us/en/).

**We'll help tell your story.  
You'll reap the rewards.**

"Being a reference company for IBM gives Synopsis many opportunities to gain exposure with influential audiences, including reporters, IT analysts and potential customers."

—Ricardo Palma,  
General Manager,  
Synopsis

When Synopsis helped Peru's largest bank migrate from Oracle to DB2, the IBM client references team was there to tell the story to a global audience of decision makers.

It's a story about cutting data management costs in half—and it's a story about Synopsis. How well are you telling *your* story?



# data management

IBM.COM/DMMAGAZINE

KNOWLEDGE. PERFORMANCE. RESULTS.



## TUNING SQL FOR THE PEOPLE

Data pros get it done in the Senate

## WATSON'S HELPERS

The road to *Jeopardy!*

## SEED THE CLOUDS

How to deliver databases as a service

# TAMING BIG DATA

Get a handle on the next (big) thing

It's a simple concept:  
When you get answers faster,  
you can ask more questions.



### IBM Netezza: The Simple Appliance Built for Serious Analytics

IBM Netezza high-performance data warehouse appliances are purpose-built to make advanced analytics on big data simpler, faster and more accessible. By bringing simplicity to advanced analytics, we help companies maximize the full potential of their data in ways that were previously unthinkable—fueling entirely new, game-changing ways of working.

To learn more and find out how you can test drive an IBM Netezza data warehouse appliance on your site, with your data, visit [www.netezza.com](http://www.netezza.com).



IBM  
**data**  
**management**

KNOWLEDGE. PERFORMANCE. RESULTS.

# Taming Big Data

12



What happens when data gets big, fast, and complicated?

20

## Tuning SQL at the Senate

Meet the data pros who make the Senate go

29

## The Man to See About Certification

He knows what you need to know. Do you?

33

## Get Your Head in the Clouds

How to prepare for databases as a service

## Departments

- 2 **Editor's Note**  
By Cameron Crotty
- 4 **NewsBytes**
- 8 **IIUG User View**  
By Stuart Litel
- 10 **IDUG User View**  
By David Beulke
- 38 **Data Architect**  
Securing DB2 Data  
By Robert Catterall
- 40 **Distributed DBA**  
The DB2 Problem Determination Tool  
By Roger E. Sanders
- 42 **Programmers Only**  
New *ORDER BY* Information: Part 2  
By Bonnie Baker
- 45 **Informix DBA**  
Fastest Informix DBA Contest III  
By Lester Knutsen
- 48 **Smarter is...**  
Making Watson Smarter...Faster  
By Howard Baldwin

## Sponsor Index

|  |        |  |   |       |  |
|--|--------|--|---|-------|--|
| Applied Analytix, Inc. ....                | 19     | www.applied-analytix.com/spreadsheet   | IIUG Informix Conference .....            | 9     | www.iug.org/conf                                   |
| DBI .....                                  | C3, 19 | www.dbisoftware.com/smarter            | Melissa Data .....                        | 5, 19 | www.melissadata.com/myibm                          |
| Fourth Millennium Technologies .....       | 47     | www.fmtusa.com                         | Netezza .....                             | C2    | www.netezza.com                                    |
| IBM .....                                  | C4     | www.ibm.com/facts                      | Niteo Partners .....                      | 19    | www.niteo.com/SmartPredictHC                       |
| IBM Client Reference Program .....         | 43     | www.ibm.com/ibm/clientreference/us/en/ | Quest Software .....                      | 7, 19 | www.quest.com/ToadDB2forDBAs                       |
| IBM Information On Demand .....            | 19, 37 | www.ibm.com/events/InformationOnDemand | Relational Architects International ..... | 3     | www.relarc.com/form/beyond_restart_and_concurrency |
| International DB2 Users Group (IDUG) ..... | 11     | www.idug.org                           | Safari Books Online .....                 | 15    | www.safaribooksonline.com/ibmmagQ2                 |

## Sponsor Article

IBM, Intel Accelerate Terabyte-Class XML Database Processing ..... **25-28**  
[www.ibm.com](http://www.ibm.com), [www.intel.com](http://www.intel.com)

# IBM data management

KNOWLEDGE. PERFORMANCE. RESULTS.

#### EDITOR

Cameron Crotty  
editor@tdagroup.com

#### MANAGING EDITOR

Stephanie S. McLoughlin

#### SENIOR EDITOR

Lisa K. Stapleton

#### ART DIRECTOR

Iva Frank

#### DESIGNER

Lalaine Gagni

#### CONTRIBUTING WRITERS

Bonnie Baker, Howard Baldwin, David Beulke, Ives Brant, Robert Catterall, Lester Knutsen, Stuart Litel, Roger E. Sanders, Jin Zhang

#### EDITORIAL BOARD OF DIRECTORS

Jeannie Cramer, Jeff Jones, Steve Miller, Jennifer Reese, Bernie Spang

#### AD SALES

Jill Thiry  
advertise@tdagroup.com

#### AD COORDINATORS

Katherine Hartlove, Kim Johnsen

#### SPECIAL THANKS TO

Lea Anne Bantsari

#### SUBSCRIPTION SERVICES AND REPRINTS

To subscribe to the print or digital version of *IBM Data Management* magazine, change your address, or make other updates to your information, please go to [ibm.com/dmmagazine](http://ibm.com/dmmagazine). For instant access to the *IBM Data Management* magazine digital edition, visit [www.ibm.dmmagazinedigital.com](http://www.ibm.dmmagazinedigital.com). For information about reprints, please send an e-mail to [customerservice@tdagroup.com](mailto:customerservice@tdagroup.com).

IBM and the IBM logo are registered trademarks of the International Business Machines Corporation and are used by TDA Group under license.

Material published in *IBM Data Management* magazine copyright © 2011, International Business Machines. Reproduction of material appearing in *IBM Data Management* magazine is forbidden without prior written permission from the editor.



#### PRESIDENT

Paul Gustafson

#### VICE PRESIDENT, STRATEGY AND PROGRAMS

Nicole Sommerfeld

#### VICE PRESIDENT, EDITORIAL DIRECTOR

Debra McDonald

#### VICE PRESIDENT, CONTENT SERVICES

Paul Carlstrom

Printed in the U.S.A.

# W

hen the lights went up on a famous television studio stage in February, it felt like the start of a new era. Even five years ago, who could have envisioned a computer capable of playing *Jeopardy!*, let alone winning?

With this issue of the magazine, we mark both the achievement that Watson represents and the critical contributions that data professionals made to the project. And in our cover story on big data, we suggest that—amazing as it is—Watson is the second act of an ongoing production.

In the first act of *The Machine Show*, we asked computers to find the answers to questions faster than ever before. Deep Blue played the part of *The Machine*, successfully answering the question, “What’s the best move?” better than its human co-star, Gary Kasparov.

In the second act, the challenge was even greater: before it could search for the answer, Watson had to tease apart the subtleties of language and layers of meaning to decide what the question was. The Watson team needed to solve big data: storing and searching a nearly infinite variety of data at previously unimaginable volumes and speeds. And as they and others grapple with that challenge, the outlines of the next stage are beginning to come into view.

With the right tools and approach, big data becomes a place that you explore without preconceived notions about what you’ll find—without predefined structures or queries. In the third act of *The Machine Show*, we will ask computers to go beyond finding the answers to our questions, or even understanding what the questions are—we’ll ask them to decide what the questions should be.

I usually close by asking you to send your thoughts, comments, suggestions, and article ideas to us at [editor@tdagroup.com](mailto:editor@tdagroup.com). This issue, I also want to encourage you to check out our electronic edition at [www.ibm.dmmagazinedigital.com](http://www.ibm.dmmagazinedigital.com). It has everything that you love about the print edition in a convenient, paper-free format. You can easily share articles with friends and colleagues, click directly through to article and partner resources, and access a full archive available to subscribers. You’ll also find us on IBM developerWorks at [ibm.com/developerworks](http://ibm.com/developerworks). See for yourself today!

Thanks for reading,

Cameron Crotty

# ABEND

## *Restart made simple™*

Smart/RESTART lets your applications restart from near the point of failure — after abends, recompiles, even system IPLs. Your applications can run restartably, often without source changes.

Smart/RESTART guarantees that your program's sequential file and cursor position, working storage and VSAM updates stay in sync with changes to DB2, MQ, IMS and other RRS compliant resources. So you can restart fast with assured integrity.

Smart/RESTART is a robust, reliable and proven solution used by Global 2000 organizations worldwide to run their mission-critical z/OS batch applications. It's the standard for z/OS batch restart.

DB2, WebSphere MQ, z/OS and IMS are registered trademarks of IBM Corp.

*For a free trial visit [www.relarc.com](http://www.relarc.com), or call +1 201.420.0400*

**Download  
our  
White Paper**

“Beyond Restart  
and Concurrency:  
z/OS System Extensions  
for Restartable  
Batch Applications”

**fai Relational  
Architects  
International**

# NEWSBYTES

## INFORMATION GOVERNANCE **COMMUNITY GROWING WORLDWIDE**

IGC draws  
1,500 members,  
solicits best  
practices ideas



# t

he wisdom of crowds isn't an oxymoron. Large groups have an uncanny ability to get the right answers—just see James Surowiecki's famous book on the subject for proof. So it's not unreasonable to say that the Information Governance Community (IGC) is getting smarter with each passing month. The group, chaired by IBM's Steven Adler, recently hit a major milestone when it attracted its 1,500<sup>th</sup> new member, a signal that the desire for stronger information governance capabilities is growing worldwide.

Increasing concern about the size and complexity of information that companies need for compliance and competitiveness—and how to manage it well—is driving the rapid growth of the organization. The IGC offers the opportunity to network with peers and thought leaders. Its aim: to develop standards and best practices to deal with the challenges of the emerging discipline.

"We started it as a way to listen to our customers, mostly in the United States and Canada," says Adler. "Now we've got a really diverse membership, including more than 100 members in South Africa and more than 80 in Belgium, which proves that information governance is a big challenge internationally," he says.

Make no mistake, organized information governance is still in its infancy. For example, there is still plenty of discussion about exactly which topics should be included in the realm of information governance, which is sometimes defined as "multidisciplinary structures, policies, procedures, processes, and controls implemented to manage information on all media in such a way that it supports the organization's immediate and future regulatory, legal, risk, environmental, and operational requirements," according to Wikipedia.

That's a mouthful—but is it useful or even meaningful? The IGC is working to define what information governance should mean in daily practice. It is also refining a maturity model that organizations can use to evaluate and improve their capabilities. (A maturity model explains different stages, techniques, and milestones in companies' paths toward effective information management.)

"We needed an organic way to look at maturity, so we open-sourced the task," says Adler. The goal, he says, is to share ideas and best practices on everything related to handling information well, so that people can see how their enterprises are doing in critical areas.

Members can also perform an online audit of their organizations' information

governance practices by answering a series of 180 questions. Based on the evaluation results, the IGC provides suggestions and networking opportunities for members to discover new ways to improve their organizations' governance. IGC also uses the survey results to understand the state of information governance worldwide.

At the upcoming IBM Information On Demand Global Conference in Las Vegas, Nevada, the group plans to give awards for outstanding examples of information governance. For more information, visit [www.infogovcommunity.com](http://www.infogovcommunity.com). To join the IGC, see [www.infogovcommunity.com/login\\_form](http://www.infogovcommunity.com/login_form).

In a related story, IBM is accepting registrations for its Information Integration and Governance Forum 2011, which consists of in-person seminars in major cities in North America, Europe, and the Asia-Pacific region. Knowledgeable professionals will discuss new ways to achieve data-driven insights and make better business decisions. The forum will also demonstrate why investing in governance programs is critical to achieving better business results.

To register for the IBM Information Integration and Governance Forum 2011, please visit [ibm.com/software/data/information-integration-governance/forum](http://ibm.com/software/data/information-integration-governance/forum).

## Do You Know SUSE Like IBM Knows SUSE?

### IBM, Novell improve tools for DB2-on-SUSE Linux development

IBM and Novell have partnered to add new free SUSE templates for IBM DB2 for Linux, UNIX, and Windows (LUW) to Novell's SUSE Studio, a service that helps developers build software appliances for the SUSE platform.

Software appliances are preconfigured disk images containing an application, middleware, a database, and a minimal operating system. Now, SUSE Studio users can build appliances that incorporate DB2 applications, thanks to the new integrated DB2 Express-C template; DB2 is also now available as a software appliance in the SUSE Gallery of Linux Appliances.

"We're lowering the time and cost required to get up and running on DB2 on the SUSE platform," says Boris Bialek, director of Information Management Technology Ecosystem and the leader of the project. "ISVs and other developers can

build on one template, then export it to different output formats," he says.

In addition, the IBM Smart Analytics System 1050 and 2050 (information management solutions for departments and midsize businesses) as well as the IBM Information Management in the Amazon Cloud (which enables DB2 and Informix database software to run on the Amazon Elastic Compute Cloud, or EC2) are available on Novell's Linux-based operating system, SUSE Linux Enterprise.

The IBM-Novell collaboration could extend the reach of DB2 applications on that platform. Bialek allows for the possibility that IBM will distribute other software for SUSE developers in the future. "I don't see why not," he says. "This has been a very productive project, and it brings new IBM information management capabilities to many more developers and customers."

### Availability

The IBM DB2 certified appliance is available within SUSE Gallery at <http://susegallery.com/a/CT44P8/db2-express-c-972-sles11-sp1-64-bit> or <http://susegallery.com/a/CT44P8/db2-express-c-972-sles11-sp1-32-bit>.

For more information about the SUSE Appliance Program, visit [www.novell.com/appliances](http://www.novell.com/appliances).

For a full range of SUSE Linux Enterprise appliances with IBM, visit [www.novell.com/ibm/appliance](http://www.novell.com/ibm/appliance).

## Data Quality Tools for IBM



IP Location



Dedupe



Property



International



Address Verification



Name Parse



Email Validation



Phone Verification



SmartMover



Free Form Parse

**Clean your database with tools that make it easy.**

Request a free trial at  
[MelissaData.com/myibm](http://MelissaData.com/myibm) or  
call 1-800-MELISSA (635-4772)

**MELISSA DATA**<sup>®</sup>  
Your Partner in Data Quality





## Get the Latest Updates

with Refreshed IBM Information Management Courses on Privacy, DB2 9, and Other Hot Topics

The IBM Information Management portfolio continues to be enhanced with courses to help data professionals update their skills. The following courses have recently been updated:

- ▶ Implementing Initiate HL7 Query Adapter
- ▶ Implementing Initiate Message Broker Suite
- ▶ Information Analysis
- ▶ Initiate Technical Boot Camp with Sample Implementation Project
- ▶ Initiate Fundamentals
- ▶ Initiate Technical Boot Camp
- ▶ DB2 9 for z/OS Utilities for Database Administrators
- ▶ DB2 9 for z/OS Database Administration Workshop Part 2
- ▶ z/OS and DB2 Basics for DB2 for z/OS DBA Beginners
- ▶ Introduction to DB2 for z/OS for Systems and Operations Personnel
- ▶ Using Optim Data Privacy Solution
- ▶ Informix 11.7 New Features

▶ **MORE INFORMATION**  
[bit.ly/IMCourses](http://bit.ly/IMCourses)



## Whaddaya Mean, It's a Software Problem?

ExtraHop application performance management solution now supports DB2

If you've ever been caught between two vendors while analyzing a network or application problem—or made the problem worse by using a performance monitor that consumes system resources or broadcasts even more traffic—ExtraHop Networks has good news.

The ExtraHop Application Delivery Assurance system—which now supports IBM DB2 as well as Informix—performs deep analysis of database protocols down to the statement level in real time. It provides statistics indicating how much network traffic results from each command or network process, how long transactions take to execute, and where failures occur.

This means that database and system administrators can often determine exactly where bottlenecks occur—and whether they're the result of database problems, other applications, or configuration problems in networks or servers. Such information can be invaluable when mobilizing resources to fix performance drags.

"When application performance begins to slow, the database administrator usually gets the blame," says Jesse Rothstein, CEO and co-founder of ExtraHop Networks.

ExtraHop is used to isolate problems to particular processes, which can then be debugged using profilers designed for even deeper analysis on IBM DB2 and Informix. (Other major database environments and applications are also supported.)

Unlike many probe- and agent-based systems, the ExtraHop solution doesn't

broadcast its own traffic or consume valuable system resources. Instead, the appliance makes a copy of network traffic and analyzes it on its own hardware.

"The problem with most of the database-specific troubleshooting tools on the market today is that they consume too many system resources, rendering them largely ineffective for running in a live production environment," says Rothstein. "Because the ExtraHop system takes a completely passive approach, database performance is not impacted, enabling the database administrator to pinpoint issues in a fraction of the time without having to worry about bogging down the entire system," he says.

Some of the typical problems that ExtraHop can troubleshoot include virtualization server misconfiguration, application problems, and some classes of hardware problems.

Several models of the ExtraHop Application Delivery Assurance system are available. The ExtraHop 2000—with a list price starting at \$59,000—processes up to 1 Gbps of traffic and provides real-time analysis for as many as 300 devices. The ExtraHop 5000 system is built for more demanding enterprise environments, and it processes up to 10 Gbps of traffic, providing real-time analysis for as many as 1,000 devices. The ExtraHop Central Manager supplies the ability to review, manage, and report on multiple devices.

▶ **MORE INFORMATION**  
[www.extrahop.com](http://www.extrahop.com)

## Prove Thy Informix Skills!

Certify your skills today by taking IBM's latest Informix exam. Passing this test of your Informix prowess shows employers that you have the important knowledge, skills, and abilities necessary to configure, install, monitor, and troubleshoot Informix Dynamic Server 11.7.

▶ **MORE INFORMATION**  
[ibm.com/certify/tests/ovr919.shtml](http://ibm.com/certify/tests/ovr919.shtml)

# [ Conferences ]



## Semantics in San Francisco

The 2011 Semantic Technology Conference (#SemTech) will be held June 5–9, 2011, at the Hilton Union Square in downtown San Francisco. Now in its sixth year, SemTech 2011 is the world's largest educational conference for executives, technologists, researchers, investors, and customers involved with semantic technologies, including Semantic Web, linked data, content management, text analytics, and search.

➤ **MORE INFORMATION**  
<http://semtech2011.semanticweb.com>

## Data Governance in San Diego

The Data Governance and Information Quality Conference (DGIQ 2011) addresses the needs of business and IT executives who are responsible for business performance, measurement, risk analysis, and accountability through effective data and information management. It will be held June 27–30, 2011, at the Catamaran Resort Hotel and Spa in San Diego.

DGIQ 2011 will feature six conference tracks and four full days of presentations and tutorials. Case studies will be presented by ConAgra Foods, Sallie Mae, Merck, Conoco Phillips, Cisco Systems, Alere Health, Swiss Re, Warner Brothers, Lexmark, MetLife, GM OnStar, and Prime Therapeutics.

➤ **MORE INFORMATION**  
<http://www.debtechint.com/dgiqconference2011>

## Want to Study at the Same School as Watson?

IBM opens analytics, big data “bootcamps” worldwide and online

IBM has announced free training resources on using IBM business analytics and information management software—in the form of more than 1,200 onsite skills “bootcamps” at client, partner, and university locations worldwide, as well as at 38 IBM Innovation Center facilities and online at DB2University.com.

The new resources are part of an initiative coming on the heels of The IBM *Jeopardy!* Challenge, where the IBM Watson system demonstrated the capability to understand natural language using a number of advanced technologies, many of which are commercially available today from IBM.

Much of the course material explains how to use IBM business analytics and information management software—and many of the underlying technologies of the Watson computing system—to capture information from new sources and use it to create business opportunities. The skills bootcamps will also cover topics such as big data—the huge data sets coming from sensors, mobile devices, social networks, cloud computing, and public sources of information—as well as closely related topics such as analytics, data management, and open source technologies including Hadoop and Eclipse Tools.

➤ **MORE INFORMATION**  
[ibm.com/developerworks/data/bootcamps](http://ibm.com/developerworks/data/bootcamps)

# The Power of IBM DB2® at Your Fingertips

## Simplify the Management of Your Complex DB2 Environment With Toad®



Toad® for IBM DB2 simplifies SQL development and administration, automates code optimization and delivers at-a-glance catalog browsing. And DB2 object management and reporting? It's a snap.

More than one million Toad users can't be wrong – after all, Toad also supports Oracle, SQL Server, Sybase, MySQL and other database platforms.

Learn about the best features of Toad for DB2 in “*The Top 10 Things a DBA Should Know About Toad for DB2*” at: [www.quest.com/ToadDB2forDBAs](http://www.quest.com/ToadDB2forDBAs)

 **QUEST SOFTWARE®**  
Simplicity At Work™



**Stuart Litel**  
is president of the  
International  
Informix Users

Group (IIUG; [www.iiug.org/president](http://www.iiug.org/president)),  
CTO of Kazer Technologies ([www.kazer.com](http://www.kazer.com)), an IBM Gold Consultant, member of  
the IBM Data Champion Inaugural 2008  
class, and recipient of the 2008 IBM Data  
Professional of the Year award.

# Big Data, Big Time

Series data, warehouse  
acceleration, and 4GLs

just got an e-mail from the editors of this magazine, saying “Stuart, you’re late again. Where’s your article?” I replied, “Dear Editors, sorry. I have been busy.” I needed a topic—fast—so when I found out that this issue of *IBM Data Management* focused on big data, I started to wonder what “big” really means. Last October at the IBM Information On Demand (IOD) conference, we heard from two speakers in the general session. The first was a Visa executive who shared that the Visa credit card system handles the transactions of 1.8 billion credit cards—and they were not even *all* his wife’s!

Even if each credit card is used only once a month—and I tried valiantly last month to up the average, using mine 27 times—that’s 1.8 billion transactions a month!

Another speaker discussed the challenges of managing time series data, meaning that you track incoming data according to the time interval when it was recorded. At IOD, an electrical company representative said his company plans to use Informix to collect electrical usage readings *every minute for each house*. Each meter will produce at least 1,440 readings per day, 30 days per month, 12 times per year. That’s more than half a million readings per customer, per year. And if you took readings every second, that’s about

31 million per customer each year! Is that big data? Sounds kinda big to me.

What do these companies have in common? They’re Informix customers. So if your big data infrastructure must cope with time series, Informix is more than up to the job.

And if you need to accelerate the flow of all that information, the recently announced Informix Ultimate Warehouse Edition (IUWE) can help in warehouse and mixed workload environments. Based on Informix 11.7, IUWE uses a powerful columns-based engine that accelerates warehouse queries a hundredfold.

The best part is that the accelerator is *completely* transparent to the application, providing fast-response, ad hoc query processing without relying on I/O

or partitioning. It compresses warehouse tables with an efficient algorithm and stores them in memory. The accelerator then accesses these tables without uncompressing data, so you don’t have index overhead. IUWE is also economical—for other people maxed out on their Visa cards—because it uses commodity hardware, requires no change to applications, and needs little database tuning.

And, here’s an announcement near to my heart because I started my career as an Informix 4GL programmer. IBM will soon resell Genero from Four Js Software ([www.4js.com](http://www.4js.com)). Think of Genero as Informix 4GL on steroids—the same language with a lot more graphical capabilities, including Windows or Web-based applications. Genero can even develop applications that run on an iPad or iPhone. Imagine running your old green-screen Informix 4GL in the palm of your hands! There’s now “an app for that.”

Finally, I recently asked my readers why they use Informix. Fewer than 20 people answered, and I have editors to feed. Please go to [www.iiug.org/president](http://www.iiug.org/president) and fill in the form. For just a few minutes of your day, you can help a columnist get on the straight-and-narrow with his editors, who are a mean and surly lot. Well, maybe not *really* mean and surly—just a little grumpy when I’m running late. \*

---

If your big data  
infrastructure  
must cope with  
time series,  
Informix is more  
than up to  
the job.

---



# GO CRUISING WITH INFORMIX IN 2011

## IIUG INFORMIX CONFERENCE

BE PART OF THE LARGEST INFORMIX GATHERING IN THE WORLD

- OVER 80 TECHNICAL SESSIONS
- FREE HANDS ON LABS
- HALF AND FULL DAY TUTORIALS
- FREE IBM CERTIFICATION TESTING
- PERFECT FOR DBAs AND DEVELOPERS
- LEARN ABOUT NEW INFORMIX 11.7 FEATURES



GET \$100 off REGISTRATION  
by using code AFJR3101

visit [iiug.org/conf](http://iiug.org/conf) FOR MORE INFORMATION

OVERLAND PARK MARRIOTT, KANSAS  
MAY 15-18<sup>TH</sup>, 2011



International  
Informix  
Users Group



**David Beulke**

(dave@davebeulke.com) is president of Pragmatic Solutions,

Inc. (PSI), a training and consulting company that specializes in designing and improving SQL, application, and system performance on DB2 for Linux, UNIX, and Windows, and z/OS. He has experience in the architecture, design, and performance tuning of large data warehouses and OLTP solutions. He is also a former president of the International DB2 Users Group (IDUG).

# Imagine What You Could Do

Free your mind, and your  
business will follow

What was once only imaginable is now possible. Big data helps a computer play and win at *Jeopardy!* against the best players in the world. Systems are smarter, bigger, and faster than ever. Our systems and data have gone from gigabytes and terabytes to petabytes and yottabytes. New data collection and analysis concepts applied to traffic patterns, DNA sequencing, cancer research, and even winning a game of *Jeopardy!* show how insights can be developed using big data.

What does this mean to you and your business? Everything. These large-scale examples aren't just fringe R&D experiments. Data retention and mining projects are now possible and available at scales that are more than just mind-boggling; they're game-changing.

Imagine identifying habits, buying trends, and government upheavals before they happen. In-depth analysis of petabytes of data is now easily achieved with all of the DB2 family platforms. The new mainframe z196 hardware delivers 50 BIPS of computing power. This new mainframe—with 96 of the fastest CPUs available on any platform in the world (5.2 GHz each) and its ability to virtualize almost any type of workload—provides one of the best private cloud environment alternatives available. Clusters of commodity boxes can scale out for extreme

parallelism to facilitate performance, data collection, and analysis. Any idea you or your executives imagine can be accomplished easily within the big data or DB2 environment because it can easily handle any amount of data.

The IDUG conference presentations show that the DB2 family can handle billions of rows on all of its platforms; DB2 pureScale and the recent Netezza acquisition and Apache Hadoop extensions extend the processing power available. Today, building a data warehouse with billions of rows is becoming commonplace for DB2 systems, whereas for other database systems, hundreds of millions of rows is an extremely big system. With other DBMS vendors playing catch-up on scalability, availability, and reliability, their processing is falling behind the times and can't

---

Imagining what is possible and then bringing the data together is what DB2 and big data are all about.

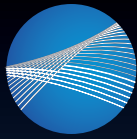
---

handle these big data systems. DB2—with its latest temporal tables, greater number of online analytical processing (OLAP) SQL functions, and Smart Analytics Optimizer integrated within a DB2 10 for z/OS solution—provides functionality several steps ahead of the DBMS competition.

One of my recent Web media projects processed more than a billion rows into a DB2 data warehouse and produced query insights in less than 12 hours—using only two commodity CPUs. Imagine if I had 96 of the new faster 5.2 GHz CPUs and DB2 for z/OS. The processing would have taken maybe minutes. Imagine the speed, the processing, and the fast answers.

Imagining what is possible and then bringing the data together is what DB2 and big data are all about. Forget about the petabytes and help your company imagine what information will make it better, more profitable, and a great company for your customers. Unleash your thinking and imagine no limits, because the DB2 family can handle it a variety of ways.

IDUG, IBM, and DB2 are at the forefront with these new ideas and many others. Come discuss your ideas at the IDUG North America conference this May in Anaheim ([www.idug.org/north-america-2011/index.html](http://www.idug.org/north-america-2011/index.html)), where you can let your imagination go wild. ✱



# IDUG

The Worldwide  
DB2 User Community

## Mentor Program Expanded

Now even more people qualify as IDUG Mentors!

Have you been to **three** IDUG DB2 technical conferences? That qualifies you to bring a first-time attendee with you at a huge 80% discount.

***A week of technical presentations and education for less than \$400!***

Get all the details at [www.IDUG.org/Mentor](http://www.IDUG.org/Mentor)

## Looking for the Tech Conference?

It's back – **at IDUG!** An IDUG DB2 Tech Conference is where DB2 is The MAIN EVENT, not a side show. Through these independent events, attendees get user-delivered, unbiased DB2 content featuring:

- User-driven technical sessions to sharpen skills
- Best practices to drive efficiencies and optimize investments
- Direction on new and existing database and related product issues
- Candid discussions with customers, IBM experts, developers and DB2 consultants

## IDUG DB2 Tech Conferences

Anaheim, California | 2 - 6 May 2011

Melbourne, Australia | 14-16 September 2011

Prague, Czech Rep. | 14-18 November 2011

MENTORING

CERTIFICATION

DB2  
EXPERTS

NETWORKING

EDUCATION

## Be a Part of the IDUG Community Strength in Numbers

more than **17,000** members representing **100** countries globally  
**thousands** of DB2 user experiences  
more than a **handful** of in-person events year-round

IDUG is an independent, user-driven community that helps DB2 users enhance their professional skills and network, drive efficiencies in their business, and optimize their DB2 technologies.

Experience IDUG year-round.

**Sign-up is FREE at [www.IDUG.org](http://www.IDUG.org)**



# TAMING BIG

By Lisa K. Stapleton



The realm of huge information flows is governed by new rules. What changes in the multi-petabyte world? And how will big data change what you do?

# DATA



**T**here's a rush of information terraforming the IT world. It flows from the data generated by 4.3 billion cell phones and 2 billion Internet users worldwide, and joins the roiling torrent of 30 billion RFID tags and hundreds of satellites incessantly sending more signals with each passing second. Now, nobody ever has to deal with all the world's data all at once. But when the whole pie grows, everyone's slices get larger. When you start measuring the pie in zettabytes, even a small piece starts to get pretty filling. Here's a sobering statistic: Twitter alone adds 12 terabytes of data every *day*—all text, and all added at a maximum of 140 characters at a time.



Dealing with data at this scale is a new frontier, and lots of different folks are approaching it in lots of different ways. But there's a growing sense that we're seeing the birth of a data challenge that's like nothing that's gone before. Some are calling it big data.

### Big data: The three V's

When they hear the term *big data*, most people immediately think of large data sets; when data volumes get into the multi-terabyte and multi-petabyte range, they require different treatment. Algorithms that work fine with smaller amounts of data are often not fast or efficient enough to process larger data sets, and there's no such thing as infinite capacity, even with storage media and management advances.

But volume is only the first dimension of the big data challenge; the other two are velocity and variety. Velocity refers to the speed requirement for collecting, processing, and using the data. Many analytical algorithms can process vast quantities of information—if you let the job run overnight. But if there's a real-time need (such as national security or the health of a child), overnight isn't good enough anymore.

Variety signifies the increasing array of data types—audio, video, and image data, as well as the mixing of information collected from sources as diverse as retail transactions, text messages, and genetic codes. Traditional

analytics and database methods are excellent at handling data that can easily be represented in rows and columns and manipulated by commands such as select and join. But many of the artifacts that describe our world can neither be shoehorned into rows and columns, nor easily analyzed by software that depends on performing a series of selects, joins, or other relational commands.

When you add volume, variety, and velocity together, you get data that doesn't play nice. And as a result, dealing with big data demands a level of database agility and changeability that is difficult or impossible to achieve using today's techniques alone. "In a traditional database, design is everything," says Tom Deutsch, IBM Information Management program director. "It's all about structure. If the data changes or if what you want to know changes—or if you want to combine the data with information from another stream or warehouse—you have to change the whole structure of the warehouse. With big data, you're often dealing with evolving needs—and lots of sources of data, only some of which you produce yourself—and you want to be able to change the job you're running, not the database design."

### Learning from extremes

Because traditional database managers and warehouses alone are often inadequate when dealing with big data, many organizations are adapting their systems to cope with a world of "badly behaved" data. These solutions vary according to the precise nature of the problems they attempt to solve—some are coping with high-velocity, high-volume information, while others must process enormous volumes of high-variability information.

But it's also possible to discern some common strategies and techniques that either reduce the magnitude of the information that needs to be stored or processed, or process it using newer, high-powered techniques that can handle the new, heavy-duty needs.

One company that's coping with all three V's is TerraEchos, a leading provider of covert intelligence and surveillance sensor systems that uses streaming data to monitor high-security facilities, national borders, and oil pipeline breaks. The TerraEchos Adelos S4 sensor knowledge system combines acoustical readings from miles of buried fiber-optic sensor arrays with data coming from diverse sensor sources such as security cameras and satellites. This enormous volume of high-variability, high-velocity data—sometimes terabytes in just a few



# SEARCH LESS. DEVELOP MORE. WITH THE INTELLIGENT LIBRARY IN THE CLOUD.



## REGISTER > YOUR TEAM FOR A FREE TRIAL NOW

Find all the latest and most relevant resources for IBM developers and IT professionals at Safari Books Online.

LEARN MORE AT: [safaribooksonline.com/ibmmagQ2](http://safaribooksonline.com/ibmmagQ2)  
...and get access to the world's most popular, fully searchable digital library.

**Safari** >  
Books Online

Move your library to the cloud and get instant, unlimited, fully searchable access to the tech books you need – including exclusive online access to books from O'Reilly Media, Addison-Wesley, Prentice Hall and more!

See for yourself why more than 15 million IT and business professionals, developers and web designers from corporations, government agencies and academic institutions use Safari Books Online.

hours—must be collected, combined with information coming from other streams, and analyzed at breakneck speeds to look for intruders, detect seismic events, or find equipment breaks.

“We’re faced with analyzing data as it passes by on a high-speed conveyor belt. We don’t have the luxury of structuring it and putting it into a database first, because we want to be able to classify it within 2 to 3 seconds,” says TerraEchos CEO Alex Philp. “With digital signal processors sampling at a rate of 12,000 readings per second—and potentially thousands of different data streams—we have to use a totally different approach so that we can respond quickly,” Philp says.

For TerraEchos, the first casualty of this nearly overwhelming data onslaught is the “extract-transform-load” paradigm that has dominated data processing for decades: extracting data from its source, performing numerous time-consuming operations to

transform it so that it fits neatly into a row-and-column format in a predetermined schema, and finally, loading it into a data warehouse. Increasingly, companies are transforming—and analyzing—incoming information as it arrives. If it meets certain conditions—for instance, if the audio stream shows a pattern that sounds like a vehicle approaching—it’s immediately flagged for more analysis and often triggers other data-collection or data-storage efforts.

“We are constantly analyzing just a few seconds of data at a time,” says Philp. “If we find something, we can trigger processes that look for the corresponding video stream or look for something interesting and, if necessary, quickly save just a few frames of the video surveillance camera data for that particular area. It’s still a massive amount of streaming data, but that really cuts down on what we have to process and store.”

#### **Filter first, ask questions immediately**

To process the incoming torrent, TerraEchos uses analytics that are designed specifically for the types of data streams



## 5 Five skill upgrades for big data opportunities

**The big picture: Companies will probably spend less time and money defining, scrubbing, and managing the structure of data and data warehouses. Conversely, they’ll spend more time figuring out how to capture, verify, and use data quickly, so these are the skills to master.**

“Today, DBAs and other IT people spend a lot of time creating cubes and stuffing data into them,” says IBM’s Roger Rea, product manager for IBM InfoSphere Streams. “That’s going to change. In the future, instead of reading data, transforming it, and then loading, you’ll just load it as fast as you can and transform it as you do your queries. This

new approach is more agile, but it means a shift in the way we think about data. It’s very different from managing according to the traditional relational model.”

What can you do to be ready to seize new opportunities? Consider the following skills upgrades.

### **1 Learn to use new big data analytics**

Some experts predict that data-mining software such as BigSheets—a spreadsheet-like interface used in IBM InfoSphere BigInsights—will make big data analytics more accessible to IT professionals and business analysts. Getting familiar with these tools and what they can do will probably benefit workers in a variety of IT disciplines.

that the company uses. The company has incorporated IBM InfoSphere Streams into its own Adelos S4 sensor knowledge system. IBM InfoSphere Streams parses incoming data and distributes the computational work involved to a myriad of processors, and its analytics packages are designed to deal with specific types of data, such as audio and video. For example, some of the analysis involves rigorous statistical analysis on the incoming waveforms to determine the probable nature of possible threats.

The trend toward specialized analytics that are tailor-made for special types of data is already accelerating. For example, analytics with algorithms for textual understanding are already being used to pore through the vast streams of tweets and e-mails produced each day to look for such things as terrorism threats and shifts in the way that a product is perceived.

The TerraEchos system combines tailored analytics—in this case, from IBM InfoSphere Streams—with advancements in parallel-processing hardware to perform millions of simultaneous, rapid calculations on the binary acoustic data coming from thousands of sensors.

“With big data, you’re often dealing with evolving needs—and lots of sources of data, only some of which you produce yourself—and you want to be able to change the job you’re running, not the database design.”

—Tom Deutsch

*IBM Information Management Program Director*

Many experts say that these techniques—filtering and analyzing data on the fly, using tailored analytics that understand how to process a variety of data in its “native” format, and bringing huge arrays of parallel processors to bear on incoming data—will soon dominate the data-processing landscape, as IT tries to cope with the special problems of high-volume, high-variety data moving at astounding speeds.

## 2 **Develop fluency in Java programming and related scripting tools**

Many of the programs used to handle big data—such as Hadoop and MapReduce—are Java-based, so learning how to program in Java is an important skill. If you already know Java, you can probably start working through online tutorials or books on Hadoop.

## 3 **Learn marketing and business fundamentals, with a focus on how to use new data sources**

Already, affinity programs are exploring the complex factors that influence customer loyalty by mining such diverse sources as customer call-center data and Twitter feeds. Understanding how to use different sources of data and to apply them to such business problems will become more important for a variety of positions, from marketing to IT.

## 4 **Develop a basic understanding of statistics**

At the core of analytical software are the fundamentals of statistics. Knowing the basics of populations, sampling, and statistical significance will help you to understand what’s possible, and to better understand and interpret what the results mean. Your best bet is a marketing or business operations statistics course, where the material is more likely to be immediately applicable.

## 5 **Learn how to combine data from different sources—especially public ones**

Much of the power of large data sets comes in combining proprietary information (such as sales data collected by companies) with publicly available data sources (such as map information or government data). Just knowing what data is available can often spark new ideas for profitable ways to combine that information.

**New technology for analyzing big data at rest**

Although better ways of handling streaming information “in motion” are a large part of solving many big data challenges, just processing extremely large amounts of data at rest can be tough, if there’s enough of it—especially if it’s high-variety data. One approach to handling this broad set of problems efficiently is through massively parallel computations on relatively inexpensive hardware. For example, IBM InfoSphere BigInsights analytics software starts with open-source project Apache Hadoop, but substitutes its own file system and adds other proprietary technology.

Hadoop is a Java-based framework that supports data-intensive distributed applications, enabling applications to work with thousands of processor nodes and petabytes of data. Optimized for the sequential reading of large files, it automatically manages data replication and recovery. Even if a failure occurs at a particular processor, data is replicated and processing continues without interruption or loss of the rest of a computation, making the system somewhat fault-tolerant and capable of sorting a terabyte of data very quickly.

“It’s hard to know what to look for in a stream of data if you haven’t already analyzed some historical data to look for patterns.”

—Tom Deutsch

*IBM Information Management Program Director*

To achieve speed and scalability, Hadoop relies on MapReduce, a simple but powerful framework for parallel computation. MapReduce breaks down a problem into millions of parallel computations in the Map phase, producing as its output a stream of key-value pairs. Then MapReduce shuffles the map output by key and does another parallel computation on the redistributed map output, writing the results to the file system in the Reduce phase of the computation. For example, when processing huge volumes of sales transaction data to determine how much of each product was sold, Hadoop would do a Map operation for each block of a file containing transactions, add up the count of each product sold in each transaction, and then “reduce” as it returned an answer.

Because it’s so simple to understand and use this technology—since it relies so heavily on just two steps, Map and Reduce—Hadoop-based systems have been used to handle a wide variety of problems, particularly in social media.

**Informing stream analysis with warehouse data**

Some observers predict that the data warehouse will go the way of the rotary phone dial, but rumors of the death of the data warehouse are greatly exaggerated. Data warehouses will continue to play a big role in many enterprises, says IBM’s Deutsch. But they’ll increasingly be used with other software to “tease out” relationships in data that can then be used to handle incoming stream data on the fly.

“It’s hard to know what to look for in a stream of data if you haven’t already analyzed some historical data to look for patterns,” says Deutsch. “But warehouse data can help you find those patterns.”

For example, Deutsch says that when University of Ontario Institute of Technology researchers first used stream-monitoring software on data captured from hospital neonatal wards, they were looking for patterns in unstructured data that might predict infant decline or recovery. They started by analyzing information from each infant, including audio recordings, heart rate, and other indicators, and eventually teased out a correlation between patterns in the audio recordings of the babies’ cries and the onset of newborn distress a few hours later.

These discoveries are being used to monitor new stream data to flag the change in cries and provide early warnings to doctors and nurses of impending problems. The ability to analyze huge amounts of high-variety warehouse data led to insights that have changed how new incoming streams are monitored.

**Bringing analytics to a wider class of users**

As data sets get bigger and the time allotted to their processing shrinks, look for ever more innovative technology to help organizations glean the insights they’ll need to face an increasingly data-driven future.

Just changing the way one views data can go a long way. “A lot of people don’t really think of unstructured data—such as video, audio, and images—as holding important information, but it does,” Deutsch says. “It’s really important to realize that this data can be just as valuable as the transactional data we’ve been collecting for years, and we have to look for new ways to put that information to work.”

One thing is clear—new ways of handling big data are accelerating almost as quickly as the flow of information that’s driving them. As TerraEchos’ Philp puts it, “I feel as if I’ve got a front row seat at the revolution.” \*

*Lisa K. Stapleton is a senior editor of IBM Data Management magazine.*



IBM invites the readers of *IBM Data Management* magazine to visit the Business Partner sites listed below.



DBAs know that managing a diverse IBM® DB2® environment can feel like doing a puzzle while blindfolded. See how Toad maximizes DB2 performance and simplifies object management, helps you perform maintenance tasks quickly and boosts database performance.

Learn more at [www.quest.com/ToadDB2forDBAs](http://www.quest.com/ToadDB2forDBAs)



### Avoiding the Spreadsheet Trap!

For finance departments, spreadsheets are like Swiss Army knives—they function as both primary analysis and reporting tools. Learn how IBM® Cognos® TM1® expands Excel's ability to support advanced analytics while removing sharing, security, and data storage obstacles.

Download our white paper at [www.applied-analytix.com/spreadsheet](http://www.applied-analytix.com/spreadsheet)



**DBI Software** is the leader in database performance tuning and optimization.

Our unique products enable DBAs to adopt a proactive methodology—escaping the trap of simple “reactive” tools and hardware addiction. Please also join us for our free educational webinar series *The DB2Night Show™*.

For more information visit [www.DBISoftware.com](http://www.DBISoftware.com)



### Patient Care and Meaningful Use Dashboards

**SmartPredict™**: an award-winning solution providing information integration, performance dashboards and scenario-based analytics for healthcare companies to:

- Provide relevant information to the care team
- Increase the efficiency of the care team
- Improve patient satisfaction and quality of care
- Demonstrate meaningful use

[www.niteo.com/SmartPredictHC](http://www.niteo.com/SmartPredictHC)

## Information On Demand 2011

### Save the Date!

Information On Demand 2011, October 23–27, 2011, Las Vegas, NV. This is the premier information and analytics event for business and IT professionals.

[ibm.com/events/InformationOnDemand](http://ibm.com/events/InformationOnDemand)

## Need to clean or enrich your database?

Get the IBM developer tools that make the job easy. Get your free trial at [www.MelissaData.com/IBM](http://www.MelissaData.com/IBM)



Your Partner in Data Quality



# Tuning SQL at the Senate:

## E PLURIBUS UNUM

When the SQL flow in the United States Senate went from static to dynamic, the database team had to see many queries, but tune them as one.

*By Ives Brant*

**Y**ou wouldn't think that U.S. senators and their staff would impose an overwhelming volume of queries on a financial system. Complex SQL and "senator" seem as oddly matched as, say, politics and rhetoric-free campaigns. Yet on an average weekday, the queries come by the thousands. DBAs everywhere know the score: the more responsive the database, the more ambitious user queries become. Budgets and projections for all the Senatorial offices, comparing them to actuals, plus other analysis, add up fast.

Within the U.S. Senate, the largest operations and support unit is the Sergeant at Arms, which provides a range of services to the Senate and includes the IT group. Lloyd Matthews, a principal DBA for the Senate Sergeant at Arms, has been a DBA for more than two decades. His group must meet performance goals and service level agreements set forth by his user community; the financial management system must respond to users' online queries in 3 seconds or less. "For us, any online dynamic query that consumes more than 4 seconds is a problem and we go after it," he says. The offenders that take more CPU and elapsed time get flagged for tuning.

In the past, Matthews' team could efficiently monitor and tune static SQL queries in both online and batch processes. The U.S. Senate operates its internal financial management system on IBM DB2 for z/OS. This 200 GB database has more than 700 tables and is queried constantly by Senate staff working for 100 senators, the Sergeant at Arms, and the Senate Disbursing Office.

As the financial management system transitioned from a batch environment to a dynamic Web-based platform, the existing tools no longer sufficed. Tuning dynamic SQL queries was taking more time and effort for the Senate DBA group. "The switch to dynamic queries was a scary transition," says Matthews.

What the Senate's IT group didn't realize is that while evolving their practices and tools to handle dynamic queries, they would adopt a new approach to addressing database performance.

“Tuning workloads for improvement turns out to be the best approach in our query-intensive environment.”

— **Lloyd Matthews**

*Principal DBA, U.S. Senate*

### **Incoming queries, take cover**

The Senate DBA team sees thousands of queries hit the financial management system in the peak period between 8 a.m. and 4 p.m. on any given weekday. Over recent years, the queries have shifted from static COBOL/CICS to 95 percent dynamic queries. Many are repeated with different literals and they can ask for "a lot of data, such as five years of history at a shot," Matthews says.

The DBAs and software developers searched for tools to help them handle query tuning, evaluating tools from several vendors including IBM, Computer Associates, and BMC Software. The financial management system runs on both mainframe and distributed platforms, so they looked at tools that support multiple platforms.

But when they looked at IBM InfoSphere Optim Query Workload Tuner software, they saw a tool that could change their entire process. "We did find other tools that capture and analyze both dynamic and static queries, but they were limited to single queries," Matthews says. "They don't handle multiple queries at once. The ability to optimize multiple queries (that is, a workload) emerged as the function that matters the most to us. It enables us to look at SQL statements by groups during an interval of time, then tune them as one. That's a powerful feature."

### **Analyze the workload, fix the workload**

Governments through history have learned they can line-item a massive budget, or consider the budget as a whole, which is much faster. Echoing these two options, in query tuning the traditional approach has been "line item." Catch an offending query, get a recommended solution, carry it out, and note the improvement. It's straightforward, but not very scalable in a dynamic query environment.



“What causes the most grief is really a lack of statistics. Data is skewed in most cases, not uniform. Once we apply the query tuner recommendations, in many cases we have immediately cut completion time from 30 seconds to 3 seconds.”

—**Lloyd Matthews**  
*Principal DBA, U.S. Senate*

The Senate’s DBAs now have a highly scalable approach: all queries from Monday to Friday are flagged for attention between an interval of time and that constitutes a “workload.” Then, on any given morning or afternoon at the U.S. Senate, a DBA specifies the workload to be analyzed and enters the criteria.

Typically, says Matthews, “I want to flag any SQL or query that executes in more than 5 seconds elapsed time or CPU time. The query tool gives me the ability to sort by different views. I normally look at views of elapsed times and CPU times for dynamic SQL. For batch SQL or large reports that usually run at night, I look at the view that shows any that read more than 100,000 pages of data within a transaction. This tells me that I may have an inefficient access path, where DB2 reads more data than it should. Normally I see this in batch or static packages.”

For dynamic workloads, Matthews has the query tuner “snap” SQL from the dynamic statement cache every 5 minutes and then, at the end of the workload, consolidate all the SQL into a final output of SQL statements. This is where reconciliation occurs. “I may have the same dynamic SQL statement in each interval, if it ran multiple times. In the final output version, however, or the reconciliation version, I see the statement only once, not duplicates. I then run my query tuning ‘advisors’ to get recommendations for that final version of SQL statements, which is more efficient since it includes no duplicates.”

### **The new work pattern**

The DBAs don’t want to lock up resources or implement a recommendation in the production environment without first testing in a production-like environment, because implementing recommendations directly into production could have a negative impact upon users. “After implementing the recommendations in our production-like environment, I run the workload again. If everything looks good—lower CPU costs, lower elapsed times, and less I/O—we implement those recommendations in our production environment over the weekend after a reorg of the data. At that point, I am confident that applying the recommendations to our production database will significantly improve performance for our dynamic queries that were flagged by the query tuner,” Matthews says.

The workload approach eliminates—actually, it reconciles—redundant statistics to produce efficient solutions. Typically, a problem query that takes 30 seconds of CPU time is immediately reduced to 2 or 3 seconds. “We’ve seen immediate improvement in CPU requirements, input/output and elapsed time, and better access paths,” Matthews says.

Matthews’ group can now proactively identify and resolve emerging issues before application performance is affected. “Throughout the day, we tune workloads that we filter out [from the entire mass of query traffic] for attention,” he says.

The DBAs take advantage of report features in Optim Query Workload Tuner to gain insight for improving database query performance. DBAs can visually compare access paths side by side, view relevant statistics, and conduct what-if analysis to determine whether a change in a new index will improve the performance of a SQL statement. A query might be too CPU-intensive, consume too much I/O, or ask for too much data to read. “If it runs infrequently or is not causing production issues and is part of the workload, we’ll tune it for next week during our weekly maintenance. However, if it’s running one hundred times per day and causing production issues, that’s a repeat offender and we’re going to take action immediately to tune it,” Matthews says.

Matthews recommends DBAs look for tools that include a function that converts long SQL statements into a format that is easy to parse. “If the statement is two pages long, a huge monstrosity, a feature of our query tool called Annotation puts all the elements into a format so you can read it. When we get a serious repeat offender and have to take it back to the developer who wrote it, the Annotation makes it easier to read and shows the DBAs where missing stats are, and skewed data. The Annotation also shows us cardinality of objects and predicates, which can be critical to the tuning process: “This is not a good SQL statement, and here is why.”

Access path graph (APG), commonly known as the visual explanation feature, is helpful as well, as it gives the DBA a visual view of the access path taken by the optimizer, along with cardinality of objects and costs of sorts.

### When data is all skewed up

“Data, in our experience, tends to be anything but uniform. Skewed data hurt us in the past; now we can attack it,” says Matthews. Optim Query Workload Tuner also analyzes and corrects for data skew. Without correct column distribution statistics, DB2 may not have enough information to make the best choice when choosing an access path to the data—especially if data is heavily skewed. That’s good.

Getting the correct statistics to DB2 is critical if the data is skewed, Matthews says, especially in an online environment where performance is essential. “Each day I run a workload against our dynamic cache. By the end of the week, by Friday, I have five days of statistics.” Applying just the Monday statistics recommendations—which might apply to 100 SQL statements captured—to the production-like environment often fixes most of the problems of all the workloads analyzed throughout the entire week and dramatically cuts execution time. With the new stats available to DB2’s optimizer, access to the data can be dramatically improved. “Within the workload, individual SQL statements that run for greater than 10 seconds or some queries that run more than 2 to 3 minutes can, with the appropriate statistics, drop to 2 or 3 seconds,” Matthews says. The gain in performance is achieved just by giving DB2 better statistics to compensate for data skew.

“Capture and analysis of both dynamic and static SQL queries was a key criterion, but not the most important one.”

—**Lloyd Matthews**  
*Principal DBA, U.S. Senate*



### Users appreciative

System users have noticed the team’s ability to identify and resolve poorly performing queries quickly. “Before, it would take us hours or days to identify problem queries. Now, in most cases, we can have a solution within minutes,” Matthews says.

The new approach has helped Matthews and his colleagues reduce CPU time usage and has let the DBAs prevent problem queries from becoming a drag on performance. Matthews’ experience has convinced him that generating new statistics on an entire workload produces better results because it reconciles all conflicting statistics.

### Best practices by the government, for the government

The Senate’s IT group moved from single-query fixes to ongoing “preventive” optimization and tuning. The new approach has cut down on the distress phone calls from users. For Matthews, it was a revelation to switch from diagnosing individual SQL statements to sift frequently through populations of queries and address them as a group.

For reasons outside the DBA’s view, the Senate tends to run different kinds of queries at different times of the day. The DBAs generally test the morning workload for two hours, then the afternoon workload for two hours, and finally the batched queries at night.

Matthews recommends that DBAs “not be afraid to use the query tuner’s variety of detailed reports. Reports generated by the workload tuner offer a wealth of detailed information about your objects.”

### When only a skilled DBA can fix it

When necessary, the Senate's DBAs still revert to the mode of catching a single major offending SQL statement. "We have all encountered that ad hoc user who uses the system once and somehow magically creates the worst possible query—so bad that it really affects performance—and we will apply a solution right there," Matthews says.

In these instances, Matthews may reorg that object, index, restat, and rebind (in the case of static queries), possibly creating a pseudo-outage. He is glad these scenarios are the exception now. "If we analyzed each query statement individually, it would require hundreds of analysis and tuning procedures, versus one," he says.

Matthews finds that query response times "are usually at their ugliest" after a new release or a migration. The DBAs anticipate that the database will react differently to the same workload. "We have to get ahead of the inevitable transition issues and wring them out. We capture stats, test them, and apply them to the production database over the weekend," he says.

### Does it always optimize?

Every method has its exceptions, and Matthews' team has learned to handle them. "Once we applied new stats, and to our surprise some non-problematic SQL statements got worse," says Matthews. "You have to be careful. These statistics will fix your offenders, but sometimes DBAs don't reorder data as frequently as they'd like to. If statistics haven't changed in a while and you have a desirable access path with current stats, DB2 will continue to use that path based on the old stats. However, when new statistics are applied, DB2 will now realize data may not be clustered as previously thought, therefore causing a change to an access path that is less desirable. Before you apply new stats to the database, I recommend you reorg the data first."

When the Senate DBAs apply the recommended statistics to the production database over a weekend, after testing them in the production-like environment, it's key to reorganize the database first, "to get back to baseline," Matthews says. He reminds his team to use the reorg utility, which

clusters the data for efficient access. "We all know the mantra: reorg, runstats, and rebind."

There are other situations where simulation in a production-like environment is not perfect. "We can simulate the performance of the queries with new stats, but we cannot truly simulate the number of concurrent users hitting at the same time," Matthews says. "Some queries may run faster due to more CPU availability when you don't have a lot of users."

### The new life of Senate DBAs: Less ad hoc

Matthews' team still gets user calls about the system response time, but much less often. "Usually they waste no time notifying us. We check our monitor, identify the resource-intensive query, and run it through the query tuner. We might find that we're missing an index or stats, and we make the fix. But our focus is on workload tuning, not ad hoc fixes."

Matthews points out that "we are not a huge shop." If they tuned one query at a time, "we'd be understaffed." Instead, they are able to spend more time doing other tasks rather than tuning queries all day. "Query workload tuning makes us proactive, not reactive," he says.

### Final vote

With the right preventive strategy and analytics tools in their repertoire, the U.S. Senate's DBA, developer, and quality assurance groups continually optimize database query design and seek to fix problems before users are affected.

"It's not just that now we can usually arrive at a solution for tuning queries within minutes versus hours or even days, which enables us to achieve our SLAs [service level agreements] in a timely fashion," Matthews says. "We have become skilled at tuning entire workloads rather than single queries. To the benefit of our users, we are now proactive versus reactive in tuning bad SQL." \*

---

*Ives Brant has experience in the database and analytics industries, and was formerly editor-in-chief at Tornado Insider magazine, a European competitor to Red Herring. These days, he writes technical and marketing content for a wide range of companies.*



# IBM, Intel Accelerate Terabyte-Class XML Database Processing

Enterprises with large XML processing needs don't have to sacrifice performance or scalability to take advantage of XML-based technology, according to recent Intel tests on servers based on the new Intel® Xeon® processor E7 family. Combined with IBM® DB2® technology, the new Intel Xeon processor E7-4870—which achieved more than 17,750 transactions per second on the Transaction Processing over XML (TPoX) 2.0 benchmark—creates a platform for fast, efficient processing of the type of XML database transactions that are increasingly common in enterprise applications.

## Strong Database Performance for New XML Processing Demands

These findings are good news for the continued success of web applications, service-oriented architectures (SOAs), and electronic data exchange between organizations. In particular, XML has emerged as a de facto standard format for electronic business records and messages, because this data format is self-describing and platform-independent, and therefore makes an easy platform for document exchange.

In addition, XML is an extensible and flexible data format, which makes it adaptable for evolving business needs and changing data processing requirements. The flexibility XML offers makes rapid, systematic content repurposing and reuse relatively easy.

## DB2 pureXML: New, Faster XML Database Storage

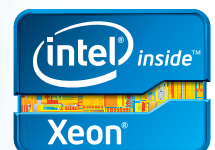
Moreover, the trend toward storing XML permanently in databases is rapidly accelerating, driven by new requirements for auditing and

compliance; the need for a more flexible and suitable data format than a rigid relational database schema can provide; and the desire to simplify applications and boost their efficiency with XML storage. Companies in virtually every industry—insurance, finance, publishing, engineering, and healthcare, to mention a few—are now using IBM DB2 pureXML® to align their database strategies with their web applications and SOAs, thereby achieving greater agility and interoperability.

DB2 pureXML also manages complex vehicle information in the automobile sector, as well as the integration of sales figures in retail companies and order management in the telecommunications industry. And government agencies employ DB2 pureXML to manage electronic forms where the variability is so great that only the XML data format can reasonably handle them.

## Speed and Agility with Intel Processors, IBM Technology

Even extremely performance-sensitive businesses such as those that use stock-trading programs—where a few microseconds could mean the loss or gain of millions in revenue—are implementing XML-based systems. The directors of these businesses know they can now have flexible systems that can be changed quickly and easily as business demands change. The speed and agility they need to meet shifting business demands is available thanks to continued breakthrough performance by Intel microprocessors and IBM DB2 pureXML technology.



|                                       | Intel Xeon X7460 | Intel Xeon X7560 | Intel Xeon X7560/SSDs | Intel Xeon E7-4870/SSDs |
|---------------------------------------|------------------|------------------|-----------------------|-------------------------|
| TPoX transactions per second (TTPS)   | 6,654            | 13,743           | 14,271                | 17,757                  |
| Users                                 | 220              | 420              | 420                   | 440                     |
| CPU %                                 | 94               | 96               | 98                    | 94                      |
| Average I/O latency <sup>1</sup> (ms) | 6.17             | 7.15             | 1.57                  | 3.72                    |
| Scalability                           | N/A              | 2.07             | 1.04                  | 1.24                    |
| Processors                            | 4                | 4                | 4                     | 4                       |
| L3 cache per processor (MB)           | 16               | 24               | 24                    | 30                      |
| Cores                                 | 24               | 32               | 32                    | 40                      |
| Threads                               | 24               | 64               | 64                    | 80                      |
| Frequency (GHz)                       | 2.67             | 2.27             | 2.27                  | 2.4                     |

**Table 1:** TPoX performance statistics for the Intel Xeon family of processors.

For example, some of the world's leading investment companies are storing and querying XML messages with financial trading information in DB2 pureXML—often using financial XML standards such as FpML or FIXML—which is exactly what the TPoX 2.0 benchmark simulates.

And when companies keep XML in persistent storage, they typically need to insert, index, query, and update XML, all with the same performance, scalability, and ACID properties (atomicity, consistency, isolation, and

durability) that relational databases have long offered for traditional relational data.

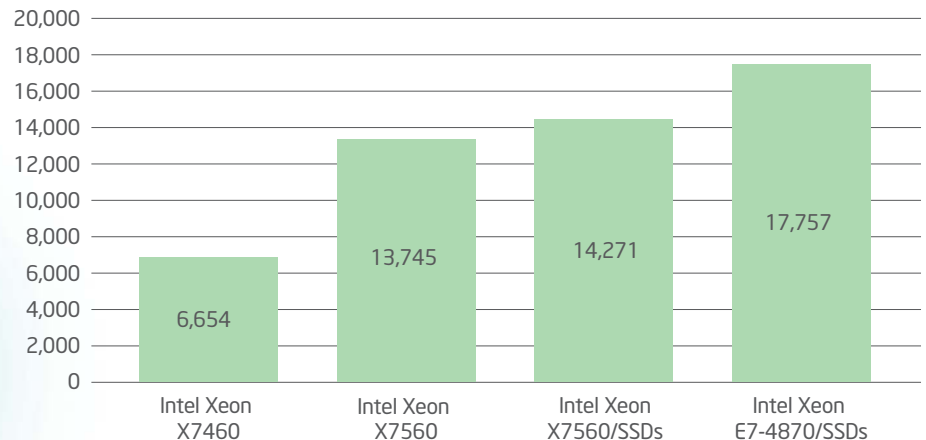
In response to these needs, DB2 9.1 and later versions offer sophisticated XML capabilities with support for SQL/XML, XQuery, XML Schemas, XSLT, and other XML-related standards. And as the latest tests show, DB2 pureXML—running on servers built on the latest Intel processors—can handle the extra demands of these new standards with speed and near-perfect scalability.

### The TPoX 2.0 Benchmark: Simulating a Real-World Trading Application

The TPoX 2.0 benchmark mimics a financial trading application, with traders placing buy and sell orders and checking the quickly changing prices and availability of securities. These TPoX tests operate on a terabyte of XML data—a midrange scale factor setting for XML transactional databases—and provide a formidable exercise of server and database technology. For more information on TPoX,

#### TPoX 2.0 Performance

TPoX transactions per second (TTPS)



**Figure 1:** Performance of different processors in the Intel Xeon processor family. The new Intel Xeon E7-4870, tested with solid-state drives, achieved more than 17,750 transactions per second on the Transaction Processing over XML (TPoX) 2.0 benchmark.

see the sidebar, “The TPoX 2.0 Benchmark: Simulating XML-Based Database Applications,” and check out the benchmark and associated documentation at <http://tpox.sourceforge.net>.

IBM and Intel conducted these tests for the Intel Xeon processor X7460, Intel Xeon processor X7560, and the latest addition to the Intel Xeon processor family, the Intel Xeon processor E7-4870. Because the throughput of the benchmark is pushed higher with the increasing performance of new generations of Intel processors, the performance demand on the storage system increases accordingly. To better manage this phenomenon, the Intel Xeon processor E7-4870 in this test uses Intel® Solid-State Drives (Intel® SSDs). To provide a fair comparison, the Intel Xeon processor X7560 was also tested with the Intel® X25-E SSD.

### Scalable, High Performance for XML-Based Transaction Processing

The new Intel Xeon processors faced stiff competition from their predecessors, which were already extremely fast. But the Intel Xeon processor family X7560—sporting a radically new architecture and simultaneous multi-threading for quicker processing—has only 25 percent more cores than its predecessor, the Intel Xeon processor family X7460, yet demonstrated twice the performance of the previous generation on the terabyte benchmark (see Table 1).

The Intel Xeon processor E7-4870 follows Intel’s tradition of continuous improvement by increasing performance an extra 24 percent, largely due to 25 percent more cores than the X7560, making it ideal for high-speed XML database access and database server consolidation. Figure 1 depicts these results.

## The TPoX 2.0 Benchmark: Simulating XML-Based Database Applications

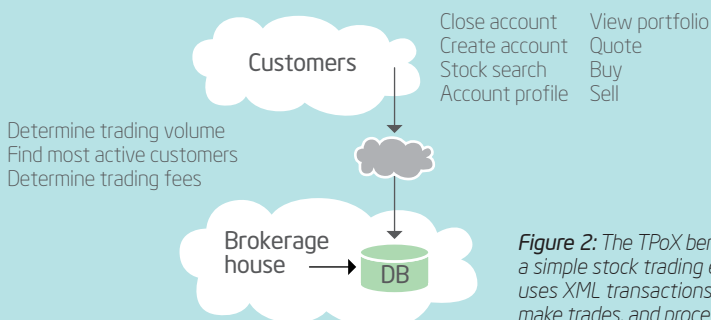
Like many enterprise applications, the TPoX 2.0 benchmark requires processing large amounts of complex information extremely quickly. Here are the benchmark’s most important characteristics:

- TPoX is an application-level XML database benchmark that executes a concurrent multiuser workload to simulate financial transaction processing. In TPoX, financial orders are represented as FIXML messages, a real-world XML format used in financial applications.
- TPoX is Java\*-based, database agnostic, and open source, with an extensible workload driver and an XML-based configuration to express the transaction mix.
- The database schema includes three tables, all with a single XML column: **security**, **custacc** (customer account), and **orders**.
- The workload mix comprises four types of transactions: queries, inserts, updates, and deletions. Seventy percent of the workload is queries; thirty percent is inserts, updates, and deletes.
- The listed score is taken in a steady state and reported as “TPoX transactions per second” (TTPS).
- TPoX includes a workload driver that spawns parallel threads that simulate concurrent database users. Each “user” connects to the database and submits a mix of transactions. The transactions are picked randomly from a set of transaction templates. At runtime, parameter markers in the templates are replaced by actual values drawn from configurable random value distributions. The workload driver collects and reports performance metrics, such as min/max/avg response time and overall throughput.

The following rules are used in defining the **update**, **delete**, and **insert** transactions:

- Customer accounts are updated to reflect trades (execution of orders).
- New orders arrive continuously, and old orders are pruned from the system eventually, at the same rate.
- Security prices are updated regularly during a business day.
- The turnover of customers is low, with few insertions and deletions of accounts.
- The number of securities remains fixed, with no deletions or insertions.

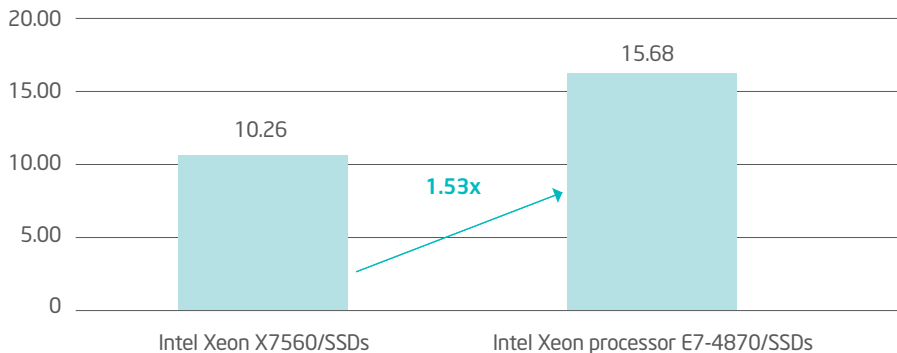
Visit <http://tpox.sourceforge.net> for more information.



*Figure 2: The TPoX benchmark simulates a simple stock trading example, which uses XML transactions to check prices, make trades, and process orders.*

## TPoX 2.0 Performance

Performance per watt (PPW)



**Figure 3:** Performance per watt for the Intel Xeon processor E7-4870 and its predecessor in the Intel Xeon processor family.

|  | Intel Xeon processor X7560/SSDs | Intel Xeon processor E7-4870/SSDs |
|--|---------------------------------|-----------------------------------|
| TPoX transactions per second (TTPS)            | 14,271                          | 17,757                            |
| Steady state power (watts)                     | 1,391                           | 1,133                             |
| Active idle power (watts)                      | 910                             | 716                               |
| Power per watt (PPW) (TTPS/steady state power) | 10.26                           | 15.68                             |
| PPW improvement                                | N/A                             | 1.53                              |

**Table 2:** Power statistics for the Intel Xeon processor family.

The advantages of the Intel Xeon processor E7 family are not limited to continuing the excellent performance and scalability demonstrated by its predecessors on XML workloads. They also include the power efficiency advantages of the latest Intel platforms.

As Figure 3 shows, on TPoX 2.0, the Intel Xeon processor E7-4870 demonstrates a 53 percent improvement on performance per watt (PPW)—also known as “performance per power”—over the previous generation of Intel Xeon processors. This is because the Intel Xeon processor E7-4870 maintains the same

thermal design point (TDP) as the Intel Xeon processor X7560, but also uses low-voltage memory due to improvements in the Intel® 7512 Scalable Memory Buffer, and so puts out less waste heat than previous generations of microprocessors. Other power statistics for the benchmark tests are also shown in Table 2.

## More Power for Industrial-Strength Database Workloads

The TPoX 2.0 benchmark results show that transactional XML workloads can benefit up to 24 percent in performance and up to 53 percent in performance per watt when moved from Intel Xeon processor X7560 to Intel Xeon processor E7-4870 platforms.

Even more important, these results show that IBM DB2 pureXML, running on servers based on the Intel Xeon processor family, readily exhibits the high performance levels needed to handle XML-heavy transaction processing on large databases. The results also confirm what many savvy enterprise IT managers already know from experience: each generation of Intel processors leaps to ever-higher processing speeds to handle the flexible, information-driven applications that run today's businesses, providing the power they have come to expect from IBM and Intel.

## Learn More

More about IBM DB2 pureXML:  
[www.ibm.com/software/data/db2/xml](http://www.ibm.com/software/data/db2/xml)

To learn more about the Intel Xeon processor E7 platform, visit [www.intel.com/itcenter/products/xeon](http://www.intel.com/itcenter/products/xeon)

<sup>1</sup> The I/O latency is the average time that a disk I/O request that DB2 makes to the operating system stays in the OS queue, plus the time for servicing by the disk device. It is obtained by using the await statistic from the iostat Linux\* tool, measured for all disks in the one-hour steady state interval. What is reported here for each run is the average await statistic across all logical drives in the system.

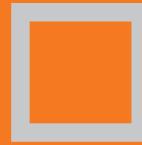
Any performance data contained herein was determined in a controlled environment. Therefore, the results obtained in other operating environments might vary significantly. Users of this document should verify the applicable data for their specific environment.

Intel, the Intel logo, and Xeon are trademarks or registered trademarks of Intel Corporation or its subsidiaries in the United States and other countries. IBM, the IBM logo, ibm.com, DB2, and pureXML are trademarks of International Business Machines Corp., registered in many jurisdictions worldwide. A current list of IBM trademarks is available on the web at “Copyright and trademark information” at [www.ibm.com/legal/copytrade.shtml](http://www.ibm.com/legal/copytrade.shtml).

\* Other product, company, or service names may be trademarks or service marks of others.



# The Man to See About Certification



By Howard Baldwin

The guru of DB2 certification tests talks about how he puts them together—and how they can help your career.

WHEN IBM EMPLOYEES HAVE A QUESTION about certification for DB2 for Linux, UNIX, or Windows, there's only one man they call. He's not in New York. He doesn't even work for IBM. He lives in the other hyphenated town in North Carolina (Fuquay-Varina, pop. 39,042), and his name is Roger E. Sanders.

In his day job, he's a consultant corporate systems engineer with EMC Corporation; he looks for ways to improve how DB2 works with his company's storage technology, including improvements in virtualization and disk-based

replication. His most recent book (his 21st on DB2 and his 9th on DB2 certification) is *DB2 9.7 for Linux, UNIX, and Windows Database Administration: Certification Study Notes* (MC Press Online, 2011).

And since 2002, he has helped IBM develop 17 DB2 certification exams, more than any other individual. Want to know why certification is important? Roger Sanders is the man to ask. We talked to him about how the tests are put together, how they can help a DBA's career, and—oh yes—about the certification test he failed.



**Q: How did you get started with DB2 certification?**

**Sanders:** My background is a little unusual. I started out as a chemist and kind of backed into computers. I was laid off in 1992 and ended up getting a job designing DB2 applications. That was when DB2 was the Database Manager part of OS/2 1.3 Extended Edition. I saw how lacking the documentation was for the product, so I decided to write a book that would help other developers build applications that interfaced with DB2 databases [*The Developer's Handbook to DB2 for Common Servers*, McGraw-Hill, 1999]. Eventually, I wrote four more books on database application development, including one on ODBC [Open Database Connectivity]. Then, my publisher [McGraw-Hill] asked me to write a certification study guide for DB2 7. So I contacted Susan Visser, who was the DB2 certification program manager at the time, and she got me involved in the exam development process.

**Q: Why is certification important?**

**Sanders:** There's a big debate about who's a better DBA: one who is certified or one who's been doing the job for 20 years, but has no certifications. Both can be equally good, but the individual who makes the effort to get certified will be forced to learn new things. And if you are constantly learning new things, you're improving your skill set.

Certification drags you out of your comfort zone. Let's say you're a DBA—your job is to keep the company database up and running, to do performance tuning, and to oversee back-up and recovery. You're comfortable performing these tasks because you do them on a day-to-day basis. But to get certi-

fied, you have to know a broad spectrum of things, and some of these things will involve activities you don't do on a daily basis. For instance, you may rarely, if ever, be required to analyze the access plans for SQL statements. But to pass a certification exam, you're going to have to know how to generate and analyze an access plan because you're going to be tested on it.

Plus, the technology changes every 18 to 24 months, and to take advantage of new features and functionality that gets introduced with each technology refresh, you have to keep up. Certification forces you to stay current.

**Q: What kinds of DBAs take certification tests?**

**Sanders:** All kinds. It's not just someone who's had a lot of hands-on experience with databases. It's someone who is intellectually curious. It's also someone who enjoys learning about the new features and functionality that are provided with each release of DB2, and who is looking at ways to incorporate those features in their own database environments. Deep compression is just one example—IBM introduced deep compression in DB2 9, but in order to use it you had to know how to enable it and how to build a compression dictionary. Plus, only tables could be compressed. In DB2 9.5, automatic dictionary creation was added so you no longer had to deal with building a dictionary yourself; in DB2 9.7, indexes and temporary tables can now be compressed. All of this results in a savings in storage costs and, in many cases, improved performance. If I were still a DBA, I would be investigating how I could put that technology to use in my company. And I probably would have become familiar with it by preparing for a certification exam.



The certified DBA  
who isn't  
a DBA

You can benefit from DBA certification even if you're not a DBA. Herb Vogel is a senior programming specialist for J.B. Hunt, the Lowell, Arkansas-based transportation and logistics company, and he's a certified DBA for both DB2 for z/OS and DB2 for Linux, UNIX, and Windows (LUW). "It helps me be a smarter programmer," he says. "But it's also like a driver's license. It's visible proof you know what you're doing."



“ [Certification] helps me be a smarter programmer. But it’s also like a driver’s license. It’s visible proof you know what you’re doing.”

—Herb Vogel

Senior Programming Specialist, J.B. Hunt

#### Q: How much do certifications help when someone is looking for a job?

**Sanders:** Some hiring managers do look for certifications. (I was such a manager; certification was something I always looked for on resumes and asked questions about in interviews.) It takes away the guesswork when they’re looking at someone who is otherwise an unknown entity. If applicants have certifications, it shows that they have some level of skill and knowledge needed to perform the job. When I got my first DB2 certification, I did so because I wanted to show my manager and any potential future employers that IBM thought I knew enough about DB2 application development to let me use their certification logo.

In terms of bringing value to a new job, having certifications in related areas can help a hiring manager get a feel for the depth and breadth of an individual’s knowledge. For example, if you are a certified DB2 DBA *and* you have been

certified on SAN [storage area network] and NAS [network attached storage] technology, a hiring manager can see that you have skills and knowledge that will allow you to help them deploy a large-scale data warehouse on an EMC storage array—and that may be just the set of skills they are looking for.

#### How certification tests get created

##### Q: How many people does it take to put together a certification test?

**Sanders:** Usually, about ten to twelve people, communicating via conference calls and live meeting sessions. Most are subject matter experts [SMEs] with IBM. Others are working DBAs, consultants, or members of the International DB2 Users Group [IDUG]. Each must specialize in some area related to DB2. Not everyone, including me, is an expert on everything.

Although he’s satisfied where he is, Vogel senses that the certification can open doors for him if he wants, if only because it shows his willingness to tackle new technology. “I’ve received phone calls from prospective employers that I wouldn’t have gotten if I didn’t have the certification,” he says, noting that he also has three other IBM certifications, all relating to DB2. He’s also been courted by the database team in his own company.

Interacting with that team is another advantage to being a programmer with a DBA certification, Vogel chuckles. “It also gives me credibility with the DBAs downstairs. If I disagree with their suggestions, I can counter with ideas of my own. They have to listen, because rather than just saying I’m good at SQL and I know a lot about databases, I can show them my certification.”



**Q: Walk us through the process of creating a certification test?**

**Sanders:** We start by building an exam blueprint, or outline. This is done by asking what the important topics are that a person at a particular certification level needs to know. For example, the DB2 Fundamentals exam tests a candidate's knowledge of security; creating and working with tables, indexes, and views; using SQL; and isolation levels and locking.

After we develop a list of topics, we decide how many questions are needed for each topic, what level of difficulty these questions should be, and who should write each question. We then separate for about three weeks to work on our writing assignments—it takes me anywhere from

“I learned something every time I prepared for an exam, and to me, that's one of the best things about certification.”

—Roger E. Sanders  
DB2 Certification Guru



one to three hours to write a single test question, and I usually end up writing between 14 to 24 questions for a new exam. Then, we get back together as a group and examine each question thoroughly to determine its validity, to make sure the correct answer indicated is indeed the correct answer, and to make sure none of the wrong answers are valid or are easy to identify as being wrong.

One thing that's very important to note: we need to provide a citation for every question we write. We must source exactly where in the technical documentation the answer appears, so that if anyone challenges the validity of the exam, its validity will stand up in court. A lot of times, I will write CLP [Command Line Processor] scripts to test my questions and answers, and these scripts are submitted along with my questions to serve as supporting material for the exam.

Finally, we perform what is known as an Angoff analysis to determine the passing score. Here, each SME goes through the entire test and identifies how many “minimally acceptable candidates” out of 100 candidates they believe will be able to answer each question correctly. The DB2 certification program manager combines the results of this analysis and determines the passing score, based on averages.

**Q: How do you determine the level of difficulty for each question?**

**Sanders:** The group establishes at the outset what skills we think someone at a particular level should possess. Then we adjust the difficulty by how we phrase the questions. You don't want a question whose answer is glaringly obvious, and you don't want a question on something that's so obscure that a test candidate would normally look it up in the documentation if they needed to use it.

A simple question might be a memory recall question, such as, “Which two communication protocols does DB2 support?” A more difficult question might have an exhibit to analyze or a scenario to work through. For instance, we might show output from a DB2 configuration file and ask, “What will happen when the current transaction log file becomes full?” To answer such a question correctly, you must know how to interpret the output shown and you must understand how DB2 logging works. The difficulty is really controlled in the wording of the question and the complexity of the exhibit or scenario.

**Q: Just out of curiosity, have you ever failed a certification test?**

**Sanders:** Yes. The very first time I went to take a DB2 certification exam, I was at an IDUG conference in Dallas where I also happened to be speaking for the very first time—and, at that time, I was terrified of public speaking. For some reason, I decided to take a DB2 6.1 Application Developer certification exam about an hour before I was scheduled to speak. Big mistake. I was too nervous to concentrate and I failed—by just one question. I gave my presentation, went by the bookstore and spent a few minutes reviewing a copy of one of my books on DB2 application development, and went back and took the exam again. This time I passed. To this day, I share that story with my students when I teach a DB2 DBA Certification Exam Crammer course and I tell them there's no shame in failing an exam. Just go back and study the areas you're weak in and try again. I've taken other certification exams over the years that I did not pass the first time. But I learned something every time I prepared for an exam, and to me, that's one of the best things about certification. It's a learning process. ✱

---

*Silicon Valley-based freelancer **Howard Baldwin** is old enough to remember when OS/2 Extended Edition was released.*

# Get Your Head in the Clouds

By Jin Zhang

Data pros are adopting cloud computing concepts to offer databases as a service—easing management burdens and sending users to cloud nine.

“

It takes weeks to set up a new database. I need it now!”

“Our dev/test databases are a mess. Why aren’t they ever cleaned up?”

Any of these complaints sound familiar? Chances are, if you’re a data professional in a large enterprise, they do. Today’s IT departments are plagued with a backlog of data administration demands. From requests for new application development and testing databases to the backup and restore of ever-growing data volumes, there’s never a shortage of busywork to keep DBAs on the run.

In an attempt to minimize the time that data professionals spend in reactive mode—responding to user requests with nonstop “database, clone, database, clone” tasks—some organizations are borrowing self-service concepts from the realm of cloud computing and moving toward a database-as-a-service or DBaaS model, where users can simply “reach into the cloud” and grab a database as needed.

It’s a tantalizing idea—especially for end users. System and software developers love the control they gain with the self-serve capabilities of DBaaS. When they are on a roll, rather than waiting for the IT department to come back a week later with a dev/test database, they can request and provision resources on the fly—keeping their momentum going and their ideas fresh.

To make this vision a reality, however, data pros behind the scenes must do a considerable amount of legwork on the back end. Building a private data cloud and successfully rolling out DBaaS to end users requires DBAs to consider a number of factors, among them the underlying hardware infrastructure of the cloud, the overarching data “best practices” to be implemented and replicated by the cloud, and finally, the services interface that will bring all of these items seamlessly to end users to complete the picture.<sup>1</sup>

<sup>1</sup>DBaaS is being deployed today both on public clouds as a type of public service and on private, on-premises enterprise clouds. For ease of explanation, this article will focus on private, on-premises implementation scenarios.

## Breaking through the clouds

Cloud computing refers to a category of technology solutions that lets users access computing resources (in this case, data resources) on demand, as needed, whether the resources are physical or virtual, dedicated or shared, and no matter how they are accessed (via a direct connection, local area network [LAN], wide area network [WAN], or the Internet).

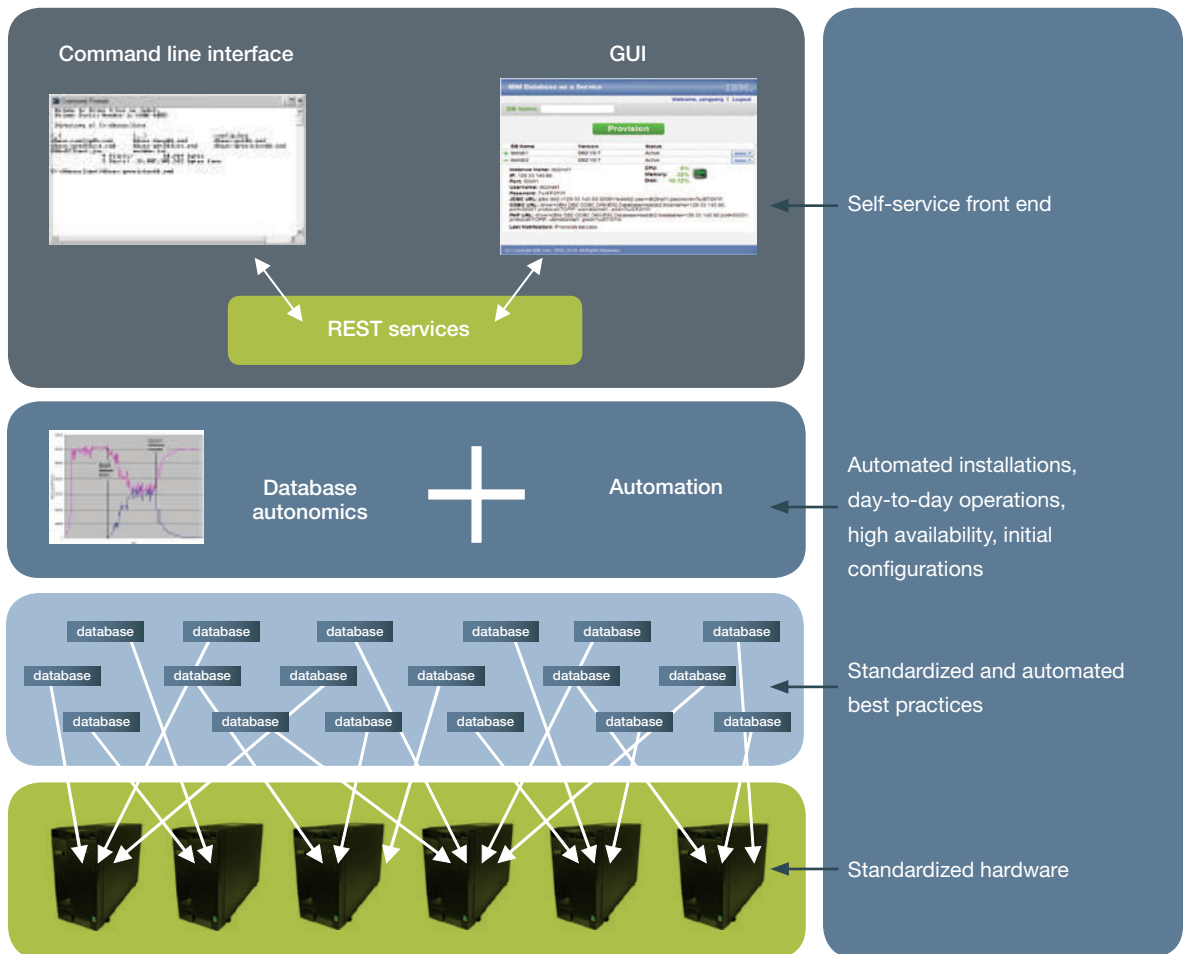
To offer DBaaS on the cloud, enterprise IT departments must construct and manage a private enterprise data cloud—a platform consisting of storage hardware, virtual images, database schemas, and more—and make that cloud available to users through a services interface.

Once this infrastructure is in place, as database needs arise, users can simply go to the cloud, request the resources they require, and gain instant access to their own personal database on demand. When they no longer need the data assets, the assets are recycled back into the cloud for reassignment, rather than being left wasted and idle.

## Step one: Build the cloud foundation

Your first stop on the way to constructing a cloud computing environment and delivering DBaaS will be to consider your underlying hardware infrastructure and ensure that it is aligned with DBaaS goals (see Figure 1). Because of the way most IT departments are structured, these hardware decisions aren't likely to take place in a vacuum. In reality, most DBAs will need to collaborate with system administrators and enterprise architecture counterparts to reach a consensus about what the hardware infrastructure will look like. This process may require compromises on all sides, so try to enter the conversation with a clear understanding of your top hardware priorities and your “nice to haves.” Not sure what those priorities should be? Read on.

As in any hardware purchase decision, many attributes will factor into the discussion—platform, storage size, speed, cost, and more. To support DBaaS on the cloud, above all you will want to ensure that your hardware is as



**Figure 1:** An infrastructure optimized for database cloud delivery emphasizes simplicity and efficiency through automation and hardware standardization.



## IBM database platforms and DBaaS on the cloud

The main article discusses how data pros can deliver DBaaS to their own end users by constructing private data clouds. IBM is committed to helping you build and deliver DBaaS via on-premises enterprise clouds, and is working to implement these capabilities on its cloud-enabled databases: IBM DB2 and IBM Informix.

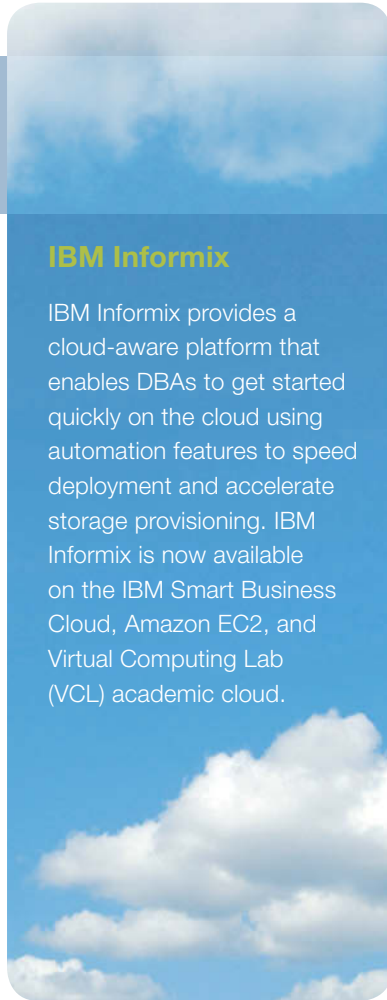
Today, IBM makes data management services on these platforms available on the cloud via an infrastructure-as-a-service (IaaS) model—or public cloud approach—with a choice of deployment options.

### IBM DB2 for z/OS and DB2 for Linux, UNIX, and Windows (LUW)

Because of the versatility of the DB2 interface and its mixed workload, multiplatform support, DB2 is well suited for cloud computing environments. Currently, DB2 LUW is available on the IBM Smart Business Cloud, IBM WebSphere Cloudburst Appliance, RightScale Cloud Management Platform, and Amazon Elastic Compute Cloud (EC2).

### IBM Informix

IBM Informix provides a cloud-aware platform that enables DBAs to get started quickly on the cloud using automation features to speed deployment and accelerate storage provisioning. IBM Informix is now available on the IBM Smart Business Cloud, Amazon EC2, and Virtual Computing Lab (VCL) academic cloud.



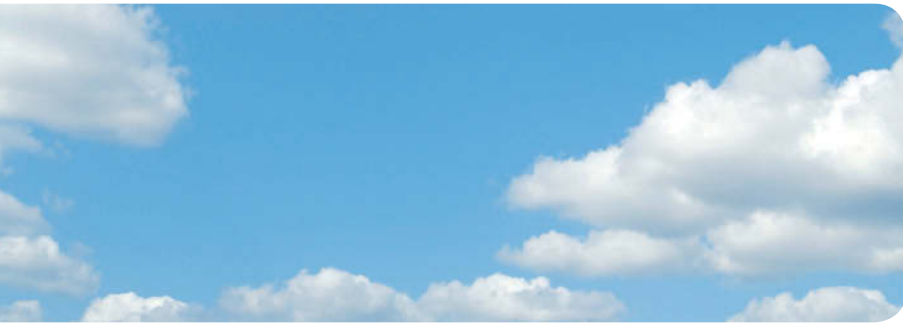
standardized as possible. Because it is far easier to automate one script running through an open, homogeneous system than many different scripts across a heterogeneous one, standardization is the key to automation. DBaaS at its heart is *nothing but* automation—the automation of the process of setting up and provisioning a database—so the more uniform your hardware platform is, the simpler it will be to set up DBaaS.

Next, take a look at the storage options available to support your database. Make sure you have a clear understanding of the types of features you will receive out of the box—including attributes such as high availability, disaster recovery, and autonomies—as well as the overall storage capacity and capabilities of your hardware infrastructure. Because this platform will ultimately form the foundation of your DBaaS offering, it is critical that you understand exactly what it is capable of—and what you can pass on to your end users. If you put in place a storage foundation that has exceptional

reliability, availability, and serviceability (RAS) capabilities, for example, you will be better equipped to provision databases on the cloud that are resilient and highly available as well.

### Step two: Identify common workloads and best practices

The next stage of DBaaS planning gives you the chance, as an experienced data pro with intimate knowledge of the inner workings of your organization and its data structures, to shine. The most critical step toward delivering DBaaS that truly brings value to your end users is to decide ahead of time what type of database templates and images should be made available on the cloud. To make such decisions, you must identify the common workloads and key processes that take place in your business environment, and collect best practices. These are the prime candidates for automation and delivery through DBaaS and the key to its successful rollout.



“Identify ‘must have’ data sets and use this information to create database templates.”

For example, DBAs can work hand-in-hand with line of business managers to identify “must have” data sets and use this information to create database templates that connect efficiently into front-end systems, work well with querying tools, and may be easily cloned for future provisioning via DBaaS. Then, personnel and systems can reach into the cloud and access entire templates that contain the latest data, up-to-the-minute information, and data structures—without creating the data administration hassles of schema changes, mapping, data migration, and more.

In other enterprise environments, DBAs may choose database images—often incorporating industry-specific metadata and reference data—as candidates for automation. A DBA familiar with business requirements can isolate an instance of a production database containing a critical set of tables, views, triggers, and stored procedures—as well as key reference data—to create a database image to be automated through DBaaS. When the business requests a database to support a new branch or test an application, there will be no need to wait for weeks while DBAs construct it. Rather, it will be instantly available via DBaaS on the cloud.

**Step three: Establish a delivery model**

Now that you’ve decided on your hardware infrastructure and identified the processes and procedures to be automated through DBaaS, your final step will be to work with end users to educate and help them select the interface through which these data services will be made available.

There are three main methods of accessing DBaaS: through a graphical user interface (GUI), command line interface (CLI), or directly via a standard representational state transfer (REST) interface. Which interface you ultimately employ will depend a great deal on end-user preference. For example, while GUI is the most user-friendly approach of the three, if end users already utilize applications that employ CLI, they may not wish to switch. Alternatively, users may wish to eliminate the need for human intervention entirely and promote tighter integration

with their environment by programming applications to communicate directly with DBaaS via REST. When you know the options, you can work with your users and help guide them to the DBaaS interface best suited to their particular wants and needs, and together select the wrapper that will pull the entire DBaaS package together.

**A cloud with a silver lining**

It’s no secret that managing the rapidly expanding data volumes and database administration needs of today’s large enterprises is no mean feat. DBAs have a tough job and there are no two ways about it. The good news is that with DBaaS, data pros are in a unique position not only to give end users new levels of freedom and service, but also to get off the hamster wheel of routine data tasks and on to the good stuff. And while it may take some groundwork to get there, as far as a cloud with a silver lining, that’s just about as good as it gets. ✱

*Jin Zhang is a program director at IBM’s Information Management CTO office.*

*To learn more about providing database services in the cloud and IBM-specific product offerings in this area, please contact the IBM DBaaS development team led by Mark Wilding (mwilding@ca.ibm.com) and Berthold Reinwald (reinwald@almaden.ibm.com).*

**RESOURCES**

**IBM products and technologies available to support cloud computing initiatives:**  
[ibm.com/developerworks/cloud](http://ibm.com/developerworks/cloud)

**DB2-specific:**  
[ibm.com/developerworks/downloads/im/udb/cloud.html](http://ibm.com/developerworks/downloads/im/udb/cloud.html)

**Informix-specific:**  
[ibm.com/developerworks/downloads/im/ids/cloud.html](http://ibm.com/developerworks/downloads/im/ids/cloud.html)

IBM Software

# Information On Demand 2011

October 23–27 | Mandalay Bay | Las Vegas, Nevada

Please join us at Information On Demand 2011, October 23 - 27 in Las Vegas to learn how IBM can help you optimize business performance and achieve breakaway results. Discover the latest IBM products and solutions and receive the very best technical and business education. You'll want to be among the more than 10,000 expected attendees learning new and innovative strategies to accelerate information initiatives and gain the practical know-how to maximize the value of your solutions.

## HIGHLIGHTS

- More than 700 Technical sessions
- Industry-focused Business & IT Leadership sessions
- Featuring Hardware, Software and Services Solutions
- 300 Customer Speakers
- Complimentary Certification Tests
- IBM's Largest EXPO
- IBM and Industry-renowned Speakers
- Networking Opportunities

## TOP 5 REASONS TO ATTEND!

- 1 Improve your skills -**  
Get deep technical education and the best strategic insight and analysis
- 2 Learn what's new -**  
Explore the latest advances in IBM Information Management, Business Analytics, and Enterprise Content Management software and solutions
- 3 Get best practices -**  
Hear from more than 300 industry leaders
- 4 Experience unrivaled networking -**  
Interact with peers, industry experts, Business Partners and IBM executives
- 5 Take action -**  
Make a direct impact on your organization with insights and actions you can quickly put to work for immediate return

**Register Today - Save \$300 with the early bird discount!**

Visit the web site today to register or for more information: [ibm.com/events/InformationOnDemand](http://ibm.com/events/InformationOnDemand)



# Securing DB2 Data

Grant privileges to a what, not a who



**Robert Catterall**  
([rfcatter@us.ibm.com](mailto:rfcatter@us.ibm.com))  
is an IBM DB2 specialist.

These days, executives are more concerned than ever about unauthorized access to data entrusted to their organizations. Their fears are justified: a recent survey showed that about a third of those polled would quit doing business with a company they perceived to be guilty of a data security breach. This is why there is such high demand for DB2 experts who can tighten up data access controls.

Role-based security is a great way to protect your organization's information assets, and it's probably easier to implement than you think. DB2 roles—and their close relatives, trusted contexts—have been available since the release of DB2 9.5 for Linux, UNIX, and Windows (LUW) and DB2 9 for z/OS. But although these DB2 releases became generally available more than three years ago, many users still seem confused regarding the purpose and advantages of roles and trusted contexts. I'll try to clear things up in this article.

## A better way to manage DB2 privileges

First, the introduction of roles and trusted contexts did not introduce any new DB2 privileges. Rather, this security capability provided a new way to *assign and manage* privileges. They can now be granted to roles instead of being assigned directly to users' authorization IDs. You can also limit the scope of granted privileges by restricting their use to trusted connections that conform to defined trusted contexts.

Managing DB2 security this way can be particularly useful when dealing with a common client/server computing scenario: An application running on a Java or a .NET (or some other) application server issues

SQL statements that are executed on a DB2 database server (in a DB2 for z/OS environment, the SQL statements would likely flow through the Distributed Data Facility, or DDF). Individual users authenticate themselves at the application server, but the application itself presents to DB2 a generic authorization ID and password that are hard-coded in a program.

If the SQL statements are dynamically prepared at the DB2 server—as is often the case for programs that use database interfaces such as JDBC or ODBC or ADO.NET—the application's generic authorization ID must be granted table privileges (**SELECT**, **INSERT**, **UPDATE**, **DELETE**) on target objects to enable successful statement execution.

But *lots* of programmers could know the application's DB2 authorization ID and password—because, as mentioned, these are embedded in program code. Someone could then use that ID and its privileges to access data in the database from outside the application, seriously weakening security.

## A useful analogy

To use an analogy from the real world, my eldest daughter is less than a year away from getting her driver's license. If only I could manage her driving with something like DB2 roles and trusted contexts, perhaps I could better control her access to our cars. I could set something up so that she could exercise the privilege of driving only between home and school, and only in the minivan (*not* the sport sedan). An impossible dream for me, of course, and probably a nightmare from my daughter's perspective.

## Using roles and trusted contexts

But you can do something very similar in DB2 if you use DB2 9 for z/OS in new function mode, or DB2 10 for z/OS. Administrators of DB2 9.5 for LUW, 9.7, and later versions have the same capability: instead of granting to a generic authorization ID a set of privileges required to execute an application's dynamic SQL statements, you could grant the privileges to a role.

Now, merely granting privileges to a role accomplishes next to nothing. Why? Because DB2 has no way of knowing either who can use the role's privileges, or the circumstances under which the role can be used at all.

That's where defining a trusted context comes in. The trusted context limits the exercise of a role's privileges to users connecting to DB2 from a particular application server—identified by an IP address—through an application that provides to DB2 a particular authorization ID, referred to as a “system” authorization ID.

Because the privileges needed to execute the dynamic SQL statements issued by the application are assigned to a role and not to an ID, the application's generic authorization ID is useless (in terms of providing someone with a means of accessing DB2 data) unless it has the privileges of the aforementioned role. And it can have those privileges only when it is used to connect to DB2 from the application server whose IP address is an attribute of the trusted context that specifies the conditions under which the role can be used. This way, security is much tighter than it would be if the application's generic ID had privileges that could be exercised regardless of the “come from” connection type.

### But wait, there's more!

That's pretty cool, but you can also set things up so that only certain individual user IDs can use the role in the defined connection context.

You have a few choices here:

- ▶ If you use the IBM WebSphere Application Server, you can propagate an end user's identity to DB2 by setting the database

property `propagateClientIdentityUsingTrustedContext` to 'true'.

- ▶ There are application programming interfaces (APIs) for JDBC (such as `getDB2Connection`), CLI (the `SQL_ATTR_TRUSTED_CONTEXT_USERID` attribute and the `SQLSetConnectAttr` functions), and .NET (where the connection string keyword `UserID` corresponds to the end user) that can be used by an application to establish a trusted connection to DB2, reuse a trusted connection with a different end-user ID, and propagate that end-user ID to the DB2 server.
- ▶ If the requester is a DB2 for z/OS system, you can provide the “system” authorization ID for a trusted connection in the requester's communications database (specifically, in the `SYSIBM.USERNAMES` table). End users' IDs will be propagated to the DB2 server as the trusted connection is reused.

Not only does this functionality let you restrict a role's privileges to designated users of a particular trusted connection, it also lets you get DB2 (and Resource Access Control Facility, or RACF) audit information that contains end users' individual IDs. This works even when those users are connecting to DB2 through an application that itself provides a single generic authorization ID when establishing connections to DB2.

If you do send end-user IDs to DB2 from an application server, you can get even more granular with respect to the roles associated with a given trusted context. For example, a trusted context could have a default role, `ROLE_A`. Assuming that the application for which the trusted context is defined propagates end-user identities to DB2, you could indicate that another role, `ROLE_B`, is usable by end user `SMITH` for a trusted connection by specifying `WITH USE FOR SMITH ROLE ROLE_B` on the `CREATE TRUSTED CONTEXT` statement. If you require authentication information—a password, for example—for `SMITH` to use `ROLE_B`, you'd add `WITH AUTHENTICATION` to the preceding `WITH USE FOR` clause.

Note that when you omit the `WITH USE FOR` clause of `CREATE TRUSTED CONTEXT`, it is as though you specified `WITH USE FOR PUBLIC WITHOUT AUTHENTICATION`. This means that the privileges of the default role associated with the trusted context are available to any individual who uses a trusted connection as defined by the trusted context.

You can even specify in a trusted context definition that a requester must communicate with DB2 using the Secure Sockets Layer (SSL) cryptographic protocol. Just make `ENCRYPTION 'HIGH'` one of the attributes of the trusted context. (`ENCRYPTION 'LOW'` corresponds to 64-bit DRDA encryption.)

Now, here are a couple of important things to remember about trusted contexts:

- ▶ For a mainframe DB2 server, a trusted context can also be defined for a local connection to DB2 through a batch job or a started task.
- ▶ A trusted context can be set up to make the context's default role the owner of any object created using the role's privileges.

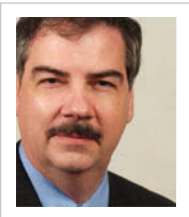
### The catch

When a user establishes a trusted connection with a DB2 subsystem—in accordance with a defined trusted context—he or she has the privileges of the associated role *plus* any privileges granted directly to his or her ID. The point here: roles and trusted contexts limit the exercise of DB2 privileges *only if* those privileges are not widely granted to users' DB2 authorization IDs. The assumption is that you'll begin to `REVOKE` privileges previously granted to individual user IDs (and/or to RACF—or equivalent—group IDs) as you phase in the use of roles and trusted contexts.

Setting up role-based security is easier than most people think. And as long as organizations seek better control over their information, there will be a demand for the increasingly fine control over data assets that DB2 provides. It's the wave of the future, folks. Catch it now, and you'll be ahead of the game. \*

# The DB2 Problem Determination Tool

How to use db2pd to find out what's really going on



## Roger E. Sanders

([roger\\_e\\_sanders@yahoo.com](mailto:roger_e_sanders@yahoo.com)), a consultant corporate systems engineer at EMC Corporation, is the author of 20 books on DB2 for Linux, UNIX, and Windows and a recipient of the 2010 IBM Information Champion award. His latest book is titled *From Idea to Print: How to Write a Technical Article or Book and Get It Published*.

Sooner or later, every DBA encounters problems. Consequently, a skill that every DBA must possess is the ability to perform a logical, systematic search of a database system for the source of any problems that might arise. Such a search often begins by answering basic questions like, “Where does the problem appear to be happening?” and “Under which conditions does the problem occur?” These types of questions help isolate the problem and provide a frame of reference in which you can limit your investigation. And in many cases, you must collect diagnostic data and then analyze that data to find a resolution.

One of the best ways to collect diagnostic data in DB2 for Linux, UNIX, and Windows (LUW) database environments is by using the DB2 Problem Determination tool (otherwise known as db2pd). In this column, I’ll describe how this tool works, and I’ll show you an example of how it can be used to troubleshoot a performance problem.

## The DB2 Problem Determination tool

The DB2 Problem Determination tool (db2pd) is designed to retrieve diagnostic information about a DB2 environment. Instead of utilizing snapshots, this tool attaches directly to DB2 shared memory sets to collect system and event monitor information. Because it doesn’t go through the DB2 engine, db2pd

doesn’t need to compete for resources, making it very lightweight and efficient. And since db2pd works directly with memory, it can retrieve data quickly and in a very non-intrusive manner. The only downside is that it may encounter data that is in the process of being changed at the same time that it’s being collected. Hence, data may not always be retrieved. (A signal handler is used to prevent db2pd from aborting abnormally when changing memory pointers are encountered; instead of aborting or reporting erroneous data, messages like “Changing data structure forced command termination” are returned in the output produced.)

Like the DB2 Command Line Processor, the DB2 Problem Determination tool can be run in interactive mode or directly from an operating system command prompt. To run db2pd in interactive mode, simply execute the command `db2pd` or `db2pd -interactive`. When either command is executed, you will be presented with a `db2pd>` command prompt, along with instructions on how to get information on db2pd’s use. From the `db2pd>` command prompt, enter one or more db2pd command options to collect the desired information; when you’re ready to exit the tool, enter `q`. On the other hand, to run the tool from an operating system command prompt, merely execute the command `db2pd`, followed by one or more of the options supported.

Alternatively, you can store the desired options in an ASCII-formatted file or the DB2PDOPT environment variable and have db2pd retrieve them during execution. For example, to retrieve a set of db2pd command options from a file named PD\_COMMANDS.TXT, you would enter a command (from the operating system command prompt) that looks like this:

```
db2pd -command pd_commands.txt
```

More than 50 commands and options are available and if you want to use them all, you can do so by executing the command `db2pd -everything`. This will cause db2pd to retrieve diagnostic information for all databases on all database partition servers. If you want to limit the type of information that gets collected, you must specify one of the following scope options:

- ▶ `-inst`  
Specifies that only instance-level information is to be collected and displayed
- ▶ `-database | -db | -d [DatabaseName]`  
Specifies that only database-level information for the database specified is to be collected and displayed
- ▶ `-alldatabases | -alldbs`  
Specifies that only database-level information for all available databases is to be collected and displayed
- ▶ `-dbpartitionnum [PartitionNumber]`  
Specifies that db2pd is to run on the database partition server specified
- ▶ `-alldbpartitionnums`  
Specifies that db2pd is to run on all active database partition servers in the instance (db2pd will report only information from database partition servers on the same physical machine that it is being run on)

Two other handy options are `-repeat` and `file=`. With the former, you can repeatedly collect information at regular intervals; the latter tells db2pd to send all output to a specific external file.

## Troubleshooting with db2pd

So just how can db2pd help you isolate the source of a problem? Suppose a database that resides on a storage area network (SAN) storage array exhibits poor performance every time a certain query is run against it. You suspect the poor performance is because a significant amount of disk I/O is being performed to satisfy the query. To test your hypothesis, you need to monitor buffer pool activity before, during, and after the query executes. One way to monitor this activity is by using the snapshot monitor; another is by using db2pd.

You can get a quick summary of all activity in all buffer pools for a database by executing a command that looks like this:

```
db2pd -db [DatabaseName] -bufferpools
```

The output produced by this command will contain, among other things, the following information:

- ▶ Buffer pool ID, name, and page size
- ▶ Number of table spaces using the buffer pool
- ▶ Current size of the buffer pool (in pages)
- ▶ Number of logical data page reads
- ▶ Number of physical data page reads
- ▶ Hit ratio for data pages
- ▶ Number of logical index page reads
- ▶ Number of physical index page reads
- ▶ Hit ratio for index pages
- ▶ Number of pages prefetched into the buffer pool, but never read

The hit ratios for data and index pages reflect the number of times a page request was handled by the buffer pool directly without requiring disk I/O. The more page requests that can be satisfied by a buffer pool, the better query performance will be. So, if hit ratios are low, a significant amount of disk I/O is taking place.

Some information db2pd can provide that you can't easily get from other monitoring tools is the actual contents of a database's buffer pools. Consequently, if a buffer

pool's hit ratio is lower than expected, you can examine the contents of that buffer pool by executing a command that looks like this:

```
db2pd -db [DatabaseName] -pages [BufferPoolID]
```

This command will tell you what objects are stored in the buffer pool specified, as well as how many data, index, long field, large objects (LOBs), and XML pages are currently in the buffer pool for a given object. To correlate table names to object IDs, execute the following command and note the IDs assigned to each table (only tables that have been accessed will be shown in the output produced):

```
db2pd -db [DatabaseName] -tcbstats
```

Along with table IDs, this command will show you the number of full table scans that have been executed, as well as the number of insert, update, and delete operations that have been performed on each table.

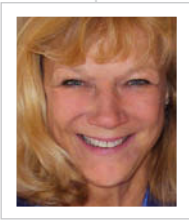
By comparing the objects being referenced by the query that's exhibiting poor performance with the information collected with db2pd, you can easily determine if a large amount of disk I/O is the source of the performance problem. And once you have identified the source, you can take the proper steps to remedy it.

## Conclusion

With more than 50 options to choose from, using db2pd can be a little intimidating at first. But with a little experimentation, and a close examination of the information provided for the `db2pd` command in the DB2 Information Center, you should be able to use this tool effectively in a relatively short amount of time. (The material in the Information Center will help you decipher the output produced when different options are used.) And hopefully, the next time a problem arises in your database environment, you'll be able to take advantage of this tool to resolve the problem quickly. \*

# New *ORDER BY* Information: Part 2

The impact of using the *RANDOM* index option on *ORDER BY* sort avoidance



**Bonnie Baker**

([bkbaker@bonniebaker.com](mailto:bkbaker@bonniebaker.com)) specializes in teaching onsite classes for corporations, agencies, and DB2 user groups. She is an IBM DB2 Gold Consultant, an IBM Information Champion, a five-time winner of the International DB2 Users Group (IDUG) Best Speaker award, and a member of the IDUG Speakers' Hall of Fame. She is best known for her ability to demystify complex concepts through analogies and war stories.

In the last issue, I began a series of columns concerning new aspects of *ORDER BY*. This column—Part 2—covers the impact of using the *CREATE/ALTER INDEX RANDOM* order option on sort avoidance. It also covers the issue of coding *only* *GROUP BY COL1, COL2* versus coding *both* *GROUP BY COL1, COL2* and a follow-up *ORDER BY COL1, COL2* in a SQL statement.

## A brief review of sort avoidance

Sort syntax such as *ORDER BY* and *GROUP BY* does not necessarily cause a data sort. With the appropriate access path, the *ORDER BY* or *GROUP BY* requirement can be met without sorting. (An appropriate access path means that an index is available and used so that the index drives the access to the data. DB2 can then return qualified rows to the program in the desired order, the order of the index.)

But what happens when the order of the index column is random? Before DB29 for z/OS, we could *CREATE* indexes with columns that had ascending key values and/or descending key values. As of DB2 9, we can also *CREATE* indexes with columns that we designate as *RANDOM*. When we use this option, key values are encoded so that the values are stored randomly. Note that identical key values will encode the same and will therefore be contiguous. Let's look at an example of index data where one of the columns is *RANDOM*. Assume an index on *DEPTNO ASCENDING, LASTNAME RANDOM, EMPNO ASCENDING* in which the decoded values are as follows:

| ASC<br>DEPTNO | RANDOM<br>LASTNAME | ASC<br>EMPNO |
|---------------|--------------------|--------------|
| A01           | Yothers            | 12345        |
| A01           | Josten             | 21234        |
| A01           | Josten             | 21235        |

|     |         |       |
|-----|---------|-------|
| A01 | Miller  | 14567 |
| A01 | Miller  | 14568 |
| A01 | Miller  | 14569 |
| B01 | Zagelow | 16657 |
| B01 | Bossman | 15678 |
| B01 | Bossman | 15679 |

Notice that *DEPTNO* is in ascending order. *LASTNAME* is in random order within each *DEPTNO*, but like values are contiguous. *EMPNO* is in ascending order within the like values of *LASTNAME*.

## The costs and benefits of *RANDOM*

Why would anyone want to keep index data in random order? Currently, online transaction processing (OLTP) workloads in a data sharing environment can experience contention on index pages, especially during *INSERTs*, when an application program *INSERTs* into indexes with columns that contain current timestamps or ever-increasing values. These column types often create insertion hot spots on index pages. Then applications must wait to acquire busy index pages.

We can use randomized index key columns to reduce contention, but at what cost? There may be more CPU usage and *getpage* operations, as well as more index page read and write I/Os. The *RANDOM* option is useful when ascending insertions or hot spots cause this contention but the resulting cost is not prohibitive.

Another possible issue is that while each distinct column value is stored randomly, like-kind values are contiguous. Therefore a randomized index column can relieve contention problems for sets of similar or sequential values, but it's no help with identical values. Identical values encode the same, and each is inserted at the same place on the index tree.

Here are some interesting facts (some beneficial and some not) about the `RANDOM` column option.

It allows equality predicate lookups, such as `LASTNAME = :LN`, but it does not support matching range lookups, such as `LASTNAME > :LN`. The option can be used in nonmatching index scans in a screening fashion and as part of index-only access. Even though values are stored as random encoded values, we can retrieve the original, decoded value of the column. The option causes `RUNSTATS` to populate `HIGH2KEY` and `LOW2KEY` with the original, decoded value of the column. Finally, the `RANDOM` column option cannot be specified in the following cases:

- ▶ For an index key column that is of variable length
- ▶ If the index is created with the `NOT PADDED` option
- ▶ As part of the `GENERATE KEY USING` clause

- ▶ If the `RANDOM` column is used to determine the partition location of a table row

Also, DB2 cannot use random order index columns as part of a sort merge join. If a join is between one table that has an ascending index on the join column and a second table that has a randomized index column, the indexes are not in the same order and cannot be merged.

### Sort avoidance with `RANDOM` columns

Now, what happens if the index that would have been ideal for DB2 to use to avoid a data sort is *not* an ordered list, but rather random?

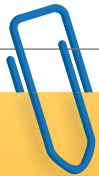
For the examples that follow, we will use our three-column index on `DEPTNO`, `LASTNAME RANDOM`, `EMPNO` to determine how DB2 might (or might not) use an index for sort avoidance.

```
1) SELECT DEPTNO, LASTNAME, ... FROM BIGTABLE
   WHERE DEPTNO in ('A01', 'B01')
   ORDER BY DEPTNO
```

For this `ORDER BY`, DB2 has no problem with the `RANDOM` column. The `RANDOM` column is not part of our `ORDER BY`. We can match on one column, and if we use the index the traditional way, the data will come back in `DEPTNO` order.

|     |         |     |
|-----|---------|-----|
| A01 | Yothers | ... |
| A01 | Josten  | ... |
| A01 | Josten  | ... |
| A01 | Miller  | ... |
| A01 | Miller  | ... |
| A01 | Miller  | ... |
| B01 | Zagelow | ... |
| B01 | Bossman | ... |
| B01 | Bossman | ... |

```
2) SELECT EMPNO, ... FROM BIGTABLE
   WHERE DEPTNO = 'A01'
   AND LASTNAME = 'MILLER'
   ORDER BY EMPNO
```



**IBM will tell your story.**

# GET NOTICED.

## Get started today.

Contact the IBM client references team at [cusref@us.ibm.com](mailto:cusref@us.ibm.com), and find out how we can help your company with:

- Case studies
- Speaking engagements
- Press releases
- Video testimonials
- Analyst interviews
- Advertising

Find more details at [www.ibm.com/ibm/clientreference/us/en/](http://www.ibm.com/ibm/clientreference/us/en/).

**We'll help tell your story. You'll reap the rewards.**

"Being a reference company for IBM gives Synopsis many opportunities to gain exposure with influential audiences, including reporters, IT analysts and potential customers."

—Ricardo Palma,  
General Manager,  
Synopsis



When Synopsis helped Peru's largest bank migrate from Oracle to DB2, the IBM client references team was there to tell the story to a global audience of decision makers.

It's a story about cutting data management costs in half—and it's a story about Synopsis. How well are you telling *your* story?

Again, there is no problem here with the `RANDOM` column. Because like-kind values of the `RANDOM` column are stored together, DB2 can do an equal lookup, matching on two columns of our index. And, if DB2 uses the index the traditional way, the data will come back in `EMPNO` order.

```
14567 ...
14568 ...
14569 ...
```

```
3) SELECT LASTNAME, EMPNO, ... FROM BIGTABLE
   WHERE DEPTNO = 'A01'
   ORDER BY LASTNAME, EMPNO
```

Here DB2 can match on only one column, and since `LASTNAME` is random, DB2 must do a data sort to bring our rows back in `LASTNAME, EMPNO` order.

```
Josten      21234 ...
Josten      21235 ...
Miller      14567 ...
Miller      14568 ...
Miller      14569 ...
Yothers     12345 ...
```

```
4) SELECT LASTNAME, EMPNO, ... FROM BIGTABLE
   WHERE DEPTNO = 'A01'
   AND LASTNAME BETWEEN 'MILLER' AND 'YOTHERS'
   ORDER BY LASTNAME, EMPNO
```

In this example, DB2 can match only on `DEPTNO` and must do an index scan, screening on `LASTNAME` within that one value of `DEPTNO` to find all of the index entries where `LASTNAME` is between our low value and our high value. The index rows will not be in `LASTNAME` order within `DEPTNO` A01. Therefore, DB2 must sort to satisfy the `ORDER BY`.

```
Miller      14567 ...
Miller      14568 ...
Miller      14569 ...
Yothers     12345 ...
```

```
5) SELECT LASTNAME
   FROM BIGTABLE
   WHERE DEPTNO = 'A01'
```

In this example, there is no `ORDER BY`. In what order will our rows be returned? Since both `LASTNAME` and `DEPTNO` are part of the index, we will get index-only access. We have one matching predicate on `DEPTNO`. Since DB2 will only use the index the “traditional” way for index-only access, the rows will be returned just like our index data, in groups of identical `LASTNAME`s in random order.

```
Yothers
Josten
Josten
Miller
Miller
Miller
```

```
6) SELECT DEPTNO, LASTNAME, COUNT(*)
   FROM BIGTABLE
   GROUP BY DEPTNO, LASTNAME
```

In this example, DB2 will *not* have to sort. Our index is in `DEPTNO` order, and within that order, in groups of identical (but unordered) values of `LASTNAME`. Therefore, using control break logic, DB2 can return distinct `DEPTNO, LASTNAME` combinations *without* sorting. The rows will be in `DEPTNO` order but will *not* be in `LASTNAME` order within each `DEPTNO`. The answer is accurate because we did *not* ask that the rows be returned in `ORDER`, only `GROUP`ed.

```
A01  Yothers  1
A01  Josten   2
A01  Miller   3
B01  Zagelow  1
B01  Bossman  2
```

Because `LASTNAME` is random in the index, if we want our rows to be in `LASTNAME` order within `DEPTNO`, we must add an `ORDER BY` clause. This requirement to add an `ORDER BY` clause in addition to the `GROUP BY` clause is new with the `RANDOM` option. Before DB2 9, we needed to add an `ORDER BY` only if we had more than one index that could be chosen to do the grouping. Remember, we do not have an `ASC` or `DESC` option on `GROUP BY`. Therefore, if we have an `INDEX` on `COL1`

`ASC, COL2 DESC`, as well as an index on `COL1 DESC, COL2 DESC`, either index can be used to avoid the `GROUP BY` sort, and the one chosen will determine the order of the output rows. To avoid unpredictable output in this situation, we would need to code the `ORDER BY` clause (even before DB2 9).

### A more realistic example

Consider the following realistic example of a situation where the `RANDOM` option might be preferable.

A user defines an index on the `INVOICE_NUMBER` column of the `INVOICE_MASTER` table in `ASCENDING` order. The user then inserts rows with an ascending sequence of values for the `INVOICE_NUMBER` column (0000100, 0000200, 0000300, ...). The index rows will be inserted at the end of the index, creating a hot spot. In this particular index, the user looks up specific invoices only by their `INVOICE_NUMBER`—in other words, only equality predicates are applied to the `INVOICE_NUMBER` column. To reduce contention, `DROP` and `re-CREATE` the index using the `RANDOM` option.

With the `RANDOM` specification, the unique `INVOICE_NUMBER` values will now be stored randomly throughout the index, preventing the contention that results from always inserting values at the end. Because the common use of this index is for looking up specific invoices by `INVOICE_NUMBER`, the new `RANDOM` ordering option provides a good solution.

### A caution regarding `RANDOM`

We must be careful with the `RANDOM` option and ensure that our solution is not creating *more* cost from increased CPU, I/O, and sort overhead. The option is best used when we are actually solving a contention problem, when our predicates on the random column are `EQUAL` predicates or, if not `EQUAL`, the random columns do not often appear in `ORDER BY` clauses. The benefit must outweigh the additional cost or the `RANDOM` option will cause more problems than it solves. And, once we begin using the `RANDOM` option, we must be very careful to code both a `GROUP BY` and an `ORDER BY` when order matters. ✱

# Fastest Informix DBA Contest III

## Performance tuning an OLTP system

Performance tuning is a continuous process for every DBA. My company has conducted three Fastest Informix DBA contests to highlight and learn what goes on in this process. Last year at the 2010 IIUG Informix Conference in Kansas City, we did something a little different that I have wanted to share with you for a while. Previous contests had focused on improving batch processing performance time; in 2010, we used transactions per minute as the benchmark. The results are very relevant to any DBA that must maintain and tune an online transaction processing (OLTP) system.

In the last contest, we used an open source OLTP benchmark, called BenchmarkSQL, that closely resembles the TPC-C standard for OLTP. It's a Java program that generates 100 sessions performing a mix of inserts, updates, and deletes against any database. The folks at AGS, creators of ServerStudio, converted the benchmark to run with Informix, and we converted the database to Informix. We challenged contestants to get the most transactions per minute during a 10-minute benchmark run.

### Concurrent user performance

DBAs are accustomed to seeing multiuser OLTP environments, but this is the first year our contest incorporated that challenge—and it was a big one! When the benchmark starts, it instantly creates 100 very active user sessions that select, add, update, and delete records

at the same time. With the default Informix setup, these sessions create lots of locking and concurrency errors. Tables created by Informix default to page-level locking, so when two or more users attempt to update a row on the same page, at least one will get a lock error.

The first thing every contestant needed to do was change every table to row-level locking. In any real OLTP database, implementing this change can be a problem. To determine which tables in a database have page-level locking and need to be modified to row-level locking, I use a helpful SQL script: `genlockmod.sql` (see Figure 1). This script will read the Informix system tables, find out which tables have page-level locking, and then generate another SQL script to change those tables to row-level locking.

Next, contestants needed to add appropriate indexes to each table, so that Informix would not scan the whole table to find the row it needed to update. This wasn't too difficult, as this benchmark database has only nine tables, and we provided the contestants a couple of hints: a list of the top SQL statements used by the benchmark, and a sample script with the primary keys for each table.

However, in the real world, adding the indexes can be difficult, and there is no magic SQL trick to do that—you must know the SQL statements being used and which fields are predicates in the `where` clauses for those SQL statements. In real life on large databases with thousands of tables, I use the `sysmaster`



**Lester Knutsen**

([lester@advancedatools.com](mailto:lester@advancedatools.com)) is president of Advanced DataTools Corporation, an IBM Informix consulting and training partner specializing in data warehouse development, database design, performance tuning, and Informix training and support. He is president of the Washington, D.C. Area Informix User Group, a founding member of IIUG, an IBM Gold Consultant, and an IBM Data Champion.



```

{
-- Author: Lester B. Knutsen
-- Email: lester@advancedatools.com
-- Advanced DataTools Corporation
-- Description: Generate SQL to set row level locking for all database
tables
}
output to lockmod.sql without headings -- Create SQL script and don't
include column headings
select "alter table " , -- Text to alter table
trim(tabname) , -- Table name
" lock mode (row);" -- Text to change the lock mode
from systables
where tabid > 99 -- Don't get the systables
and tabtype = "T" -- Get real tables not views
and locklevel = "P" -- Get tables with page level locking
order by tabname;

```

**Figure 1:** The SQL script `genlockmod.sql` generates another script to alter tables to row-level locking.

database table `systabprof` and the column `seqscans` to identify which tables may need an index. Figure 2 shows the code for `tabscans.sql`, which identifies the tables with the most sequential scans.

Sometimes a table is so small that it is faster to do a sequential scan on a table that fits in one or two pages, than to read index pages and data pages to find a row.

### Disk layout for performance

The contest machine ran Linux with four CPUs and 3 GB of memory, but only one disk drive. Databases are very disk I/O intensive, so the single drive was a significant limiting factor. It also meant that a lot of the sophisticated disk layout tuning that you can do with Informix did not help. In fact, the DBA with the fastest time did not make any changes to the baseline disk layout: four dbspaces, a rootdbs, a logdbs, a tempdbs, and a datadbs. The second fastest DBA added one dbspace to move the physical log from the rootdbs to a separate dbspace, but in my testing I have found that layout to be slightly slower since the dbspaces are all on the same physical disk. Configuring different page sizes did not help much either.

With only one disk, is there anything that will help I/O? Four of the top five contestants turned on `Direct I/O`, a new parameter in the `ONCONFIG` file that works on Linux and AIX systems to speed up disk I/O to file systems. Another strategy is to limit the number of processes writing to disk to avoid I/O contention between the processes: minimize the number of `CLEANERS`, `LRUs`, and `AIOVPs`.

```

{
-- Author: Lester B. Knutsen
-- Email: lester@advancedatools.com
-- Advanced DataTools Corporation
-- Description: Sysmaster query to identify tables with the most scans
}

database sysmaster;

select dbsname,
tabname,
sum(seqscans) total_scans
from sysptprof
where seqscans > 0
group by 1, 2
order by 3 desc;

```

**Figure 2:** The script `tabscans.sql` is a sysmaster query to identify tables with the most sequential scans.

### Memory tuning

In many cases, the biggest performance increases will come from adding as many buffers as you can. Informix uses buffers to store data on the first request, so that it can be shared and reused with other users without the cost of reading it again from disk. The benchmark machine was running a 32-bit version of Informix and was limited to 2 GB of memory (the 64-bit version of Informix does not have this limit). The top two entries used as much as possible for buffers, 1.6 GB and 1.5 GB respectively. They all also examined the amount of virtual memory and set `SHMVIRTSIZE` accordingly. The top two also set the `RESIDENT` parameter to -1, which tells Linux to keep all the Informix memory segments in memory and not swap them to disk.

Another key parameter used was `DS_NONPDQ_QUERY_MEM`. This parameter allows you to increase the default memory for user sort space. If the sort will fit in memory, it will be very fast; otherwise, it overflows to disk and will be much slower. On OLTP systems, where user sessions are performing quick queries, tuning this parameter can be very important.

### CPU tuning

Informix has an `ONCONFIG` parameter `VPCLASS` that controls how many CPUs will be used by the database server. This is also critical to get right. The default is to use only one CPU. The machine in the contest had four CPU cores, and the top configuration tuned Informix to use all four CPU cores. The only process running on this machine was the database server. Two of the contestants set the number of

## 2010 Fastest Informix DBA contest winners

The Fastest Informix DBA and Grand Prize Winner: **Tatiana Saltykova**

Fastest IBM Developer (IBM employee): **Spokey Wheeler**

Fastest DBA on Monday (first day of the contest): **Tom Girsch**

Fastest International DBA (non-U.S. resident): **Denis Zhuravlev**

Fastest Mid-aged DBA (30 to 50): **Andrew Ford**

Fastest Youngest DBA (under 30): **Pam Siebert**

Fastest Old-Timer (longest experience with Informix): **Wenching Chiang**

Fastest Domestic DBA (U.S. resident): **Eric Rowell**

The winners and the results of the contestants are available at [www.advanceddatatools.com/Informix/index.html](http://www.advanceddatatools.com/Informix/index.html).

We are planning an exciting new contest for the 2011 IIUG Informix Conference on May 15 to 18 of this year. For more information about the conference, visit [www.iiug.org/conf/2011/iiug](http://www.iiug.org/conf/2011/iiug).

CPUs for Informix to be greater than the number of physical CPUs on the machine. The CPUs on this machine were fast enough to handle this setting, which helped their performance.

### Conclusion

One important note: performance tuning is very system specific. I took the top five contestants' configurations and ran them on a different machine from the one we used for the contest, and I got very different results. Know the hardware configuration of your Informix server and tune accordingly.

Conducting this contest is one of the most exciting things I do every year, and the feedback I have received from the participants is very stimulating. Since the IIUG Informix Conference, 28 people have downloaded the contest and run it on their own. Congratulations to all the DBAs who worked hard on this and especially to the winners of the contest. The results are listed in the sidebar "2010 Fastest Informix DBA contest winners." \*

### RESOURCES

**BenchmarkSQL:** <http://sourceforge.net/projects/benchmarksql>



**Premier  
Business  
Partner**

# Fourth Millennium Technologies

*Your Trusted Partner for Database Products and Services*

### SERVICES:

DATABASE HEALTHCHECK  
DATABASE ADMINISTRATION  
DATABASE PERFORMANCE AND TUNING  
REMOTE (24x7) DBA SUPPORT  
LOGICAL AND PHYSICAL MODELING  
ETL PROCESSING  
BI/WAREHOUSE QUICKSTART  
BI/WAREHOUSE IMPLEMENTATION  
AND MORE...

### PRODUCTS:

IBM SMART ANALYTICS SYSTEMS  
IBM DATA MANAGEMENT PRODUCTS  
IBM PASSPORT ADVANTAGE SOFTWARE  
P SERIES  
X SERIES  
IBM STORAGE PRODUCTS  
NETEZZA™  
AND MORE...

[www.fmtusa.com](http://www.fmtusa.com)

Did you ever wish you could deploy your new DB2 database, datawarehouse/BI system just by pushing a button?

Would you like to improve the performance or viability of your existing data environment?

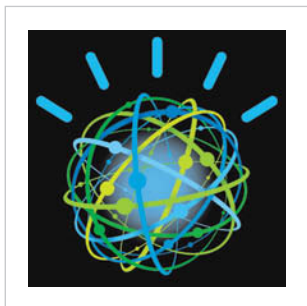
We can help...just push these ten buttons (281.481.6827) and ask for Bob.

Solutions



# Smarter is...

## Making Watson Smarter...Faster



*Howard Baldwin is a Silicon Valley-based freelancer who writes about business and technology issues.*

Last February saw a lot of excitement, as IBM's Watson supercomputer beat the humans Ken Jennings and Brad Rutter on *Jeopardy!* But if you look behind Watson's wonder, you'll find some database technicians whose jobs weren't that much different from the jobs of most DBAs. Without the long hours they worked tuning the DB2 database that stored metadata associated with Watson's thought process, the outcome might have been different—or at least taken a lot longer to deliver.

Although Watson made the result look effortless, a lot of work was going on behind the scenes. The question-and-answer data for Watson itself was stored in the UIMA format, which is highly suited for the analysis of unstructured data. However, only a subset of the analysis metadata was important to understand how Watson arrived at its answers; this data was pulled out of the UIMA analysis metadata and stored in a Derby open-source database.

A custom Web application, the Watson Error Analysis Tool (WEAT), was then used to visualize the data. For example, after a series of test matches, the team would use the WEAT tool to see how Watson was thinking. "We wanted to see why it chose the wrong answer from its top options," says Eddie Epstein, the IBM engineer responsible for Watson's scalability. "What caused a wrong answer to be ranked ahead of a better one?"

But the WEAT tool wasn't perfect. Running a single query could take minutes, and a group of developers all trying to use WEAT at the same time only made things worse. WEAT was working, and Watson was getting smarter, but the IBM team needed to make progress faster.

Enter Tong Fin, who implemented several key changes. First, he moved the data from the Derby database to DB2. This immediately improved performance, particularly for multiple WEAT instances accessing a common database. He also optimized the schema of the metadata within DB2 and achieved an order of magnitude speedup for the slowest queries.

Finally, when Fin looked at the WEAT metadata results, he realized that for a large number of queries, Watson was accessing a common set of data. Only for a smaller number of queries was it accessing a subset of that data. Fin separated the subset of less-common data into a second table, and then optimized the first table to run even faster.

The result of Fin's work? WEAT ran faster, the development process went faster, and well, you know the rest. Watson set a new mark for what computers can achieve—and a DBA helped it get there. \*

In the time it takes to  
 browse through this magazine, you could  
**download, install, and see remarkable breakthrough**  
 performance **results** from our software

We're the world's fastest route to DB2 LUW productivity and performance. In fact, you can download, install, and actually see results in minutes instead of days, weeks, or months.

Only DBI offers all these reasons to be your smart choice for performance monitoring and tuning tools:

| DBI                                 | IBM OPM*                 |   |
|-------------------------------------|--------------------------|---|
| <input checked="" type="checkbox"/> | <input type="checkbox"/> | Patented SQL Workload Cost Analysis             |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> | Tuning Results in 5 Mouse Clicks or Less        |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> | Automatically Tracks Database Changes           |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> | Trend Charts Plot Change Events on Graphs       |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> | Identifies SQL driving I/O to a Table           |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> | Summarizes Programs, Users, and their SQL Costs |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> | Shows All Users that ran a particular SQL       |

Contact DBI to experience our results-oriented Proof of Concept process and discover why **Smarter People choose Smarter Tools** for DB2 performance and cost optimization. Act now and receive a free white paper by Scott Hayes, IBM DB2 GOLD Consultant, which discusses critical performance measurements and provides complimentary SQL commands.

[www.DBISoftware.com/smarter](http://www.DBISoftware.com/smarter)



TDA GROUP  
800 West El Camino Real, Suite 101  
Mountain View, CA 94040  
U.S.A.

# Netezza. Up and running in 24 hours, not 24 days.



Get set up in hours instead of days, and start counting returns in minutes instead of hours. All with IBM's Netezza data warehouse appliance for high-performance analytics. It gives you analytics reports at supersonic speeds. At a fraction of the cost of Oracle Exadata. Get real, actionable business results fast.

[ibm.com/facts](http://ibm.com/facts)

COST comparison based on publicly available information as of 2/10/2011 for an Oracle Exadata X2-2 HP Full Rack and a full rack of Netezza TwinFin. The cost to acquire Netezza can be as low as 1/6 of Exadata if a client is acquiring new Oracle database licenses and as low as 1/2 if using existing Oracle database licenses. IBM, the IBM logo, ibm.com, Smarter Planet and the planet icon are trademarks of International Business Machines Corp, registered in many jurisdictions worldwide. Other product and service names might be trademarks of IBM or other companies. A current list of IBM trademarks is available on the Web at [www.ibm.com/legal/copytradeshtml](http://www.ibm.com/legal/copytradeshtml).  
© International Business Machines Corporation 2011.