

DB2 Text Search

Overview



© 2013 IBM Corporation

DB2® Text Search version 10.5 provides a rich set of full-text search features. This presentation provides a short overview of some of the main features.

Text Search in DB2

- Search structured and unstructured data stored in a DB2 database
 - Integrated to DB2 search engine enabling text search within a single SQL statement
 - Significant performance benefits compared with SQL 'like' operators
 - XML is less like traditional data stored in a database and XML applications often rely on a full text index
 - XQuery fn:contains searches in "typed-value" of target element does not use an index
 - full text index searches in all text under path
 - DB2 provides two text search solutions that can co-exist in same database instance
 - DB2 Net Search Extender (deprecated)
 - DB2 Text Search (availability DB2 V9.5 FP1)
 - Text index management is (partially) integrated into DB2

DB2 Text Search provides extensive capabilities to search both structured and unstructured data stored in a DB2 database. The Text Search functions CONTAINS(), SCORE() and xmlcolumn-contains(), are tightly integrated with DB2 such that they can be used within SQL and XQuery statements.

Text indexes offer significant performance benefits compared to SQL 'LIKE' operator. DB2 Text Search supports pureXML[®] and searches the data using SQL, SQL/XML and XQuery statements. The XML search functionality provided by DB2 Text Search makes it highly attractive due to its performance benefits.

DB2 provides two text search solutions that can co-exist in the same database instance. One is the DB2 Net Search Extender, or DB2 NSE, and it is already deprecated from version 10.1. The second one is the DB2 Text Search which is available since DB2 V9.5 FP1 and is the new powerful Text Search engine.

Text index management is partially integrated with DB2, that is, certain table operations like drop table, are restricted when a text index is created and in an active state.

DB2 Text Search features

- Integration with DB2
 - Integrated install; no additional license required
 - Stored procedures for administration
- Document Indexing
 - Native XML support
 - Multiple document formats, including rich text
 - Incremental and asynchronous index updates
- Advanced search technology
 - SQL functions CONTAINS and SCORE
 - SQL, SQL/XML and XQuery support including XPath-like syntax to search XML docs
 - Built-in SQL functionality combined with DB2 optimizer
 - Linguistic processing for 20+ languages/locales

DB2 Text Search is an optionally installable component whose installation and configuration is fully integrated with DB2 server products. No additional license is required. It also provides several administrative stored procedures that can be called from any DB2 client machine.

The document formats supported are Text, HTML and XML. In addition, DB2 Text Search supports indexing and searching of rich text documents with the DB2 Accessories Suite.

The text index update is asynchronous and not part of the same transaction with which the data column in the base table has been updated.

DB2 Text Search provides advanced search technology by combining the SQL functionality with DB2 optimizer where least expensive plans are chosen.

SQL, SQL/XML and XQuery support is provided including XPath like syntax to search XML documents. It also supports advanced linguistic processing for 20 plus languages and locales.

Text Search installation and configuration

- Optionally installed during DB2 installation
- Configure per instance using the DB2 installer, db2icrt or the Configuration Tool
- Text Search server port per instance
- Authentication token for security
- Configuration directory
- Rich Text feature requires Accessories Suite for DB2

One DB2 instance is associated with only one DB2 Text Search server instance. The instance should be initially configured for Text Search using DB2 installer or using the `-j` option with `db2icrt` or `db2iupdt` or `db2iupgrade`.

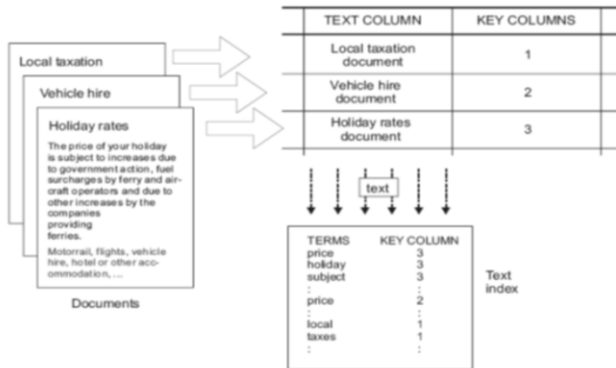
The configuration tool, referred to as `configTool`, can be used to do initial and subsequent configuration of DB2 Text Search. It authenticates with the DB2 Text Server for all administration commands and searches using authentication token, and the communication with the Text Server happens through HTTP Port. The DB2 Text Search instance should be configured by generating the authentication token and by setting a valid port number using the `configTool`.

All the configuration files are located under `$INSTHOME/sql/lib/db2tss/config` where, `INSTHOME` is the *instance home directory*.

The DB2 Accessories Suite should be downloaded separately to support rich text documents.

Text indexes

- Search data in plain text, HTML, XML and rich text documents
- Text index keeps track of significant terms from text documents



5

Overview

© 2013 IBM Corporation

Text index must be created in order to efficiently search the text documents. The parser extracts significant terms from the document, which are used to build a text index. Any changes to the table such as inserts, updates and deletes on the documents, are tracked and reflected in the text index with the next incremental update.

Full-text search features

- Combine full-text search on structured and unstructured data with relational predicates or XML expressions
 - SELECT author, story FROM books
WHERE CONTAINS (story, 'good fellow') = 1 AND YEAR >= 2000
 - SELECT chapter FROM books
WHERE CONTAINS(chapter,
'@xmlns:*/chapter/title
[. contains("good fellow")]') = 1
 - xquery for \$i in db2-fn:xmlcolumn-contains
('BOOKS.CHAPTER', '@xmlns:*/chapter
[./title contains("good fellow") and
./@number = 10] ')
return \$i
- Search in plain text, html, XML, and richtext/proprietary formats
 - Use wildcards, boolean expressions, phrase search, optional terms, synonyms, ...
 - Get scores: how well does the document match search criteria?
 - Set query language, result limits, ...

DB2 Text Search provides several features and functionality to search the structured and unstructured data. The examples displayed on this slide show how to perform searches using SQL/XML queries. Xmlxp is used to search specific portions of an XML document.

xmlcolumn-contains() is used in an XQuery expression to search for a specific term and attribute.

The search scope can either be extended or reduced using wild cards, Boolean expressions, phrase search, optional terms, synonyms, and so on. Other features include getting the score of a document, ability to set query language, result limit and so on.

Search examples

- Search with CONTAINS
 - SELECT author, story FROM books
WHERE CONTAINS (story, 'cat') = 1 AND YEAR >= 2000
- Optional QUERYLANGUAGE parameter
 - SELECT author, story FROM books
WHERE CONTAINS (story, 'cat', 'QUERYLANGUAGE=en_US') = 1
- Optional RESULTLIMIT parameter
 - SELECT author, story FROM books
WHERE CONTAINS (story, 'cat', 'RESULTLIMIT=10') = 1
- Using SCORE function
 - SELECT author, story FROM books
WHERE CONTAINS (story, 'cat') = 1
ORDER BY SCORE (story, 'cat') DESC
 - SCORE uses values between 0 and 1

A few more search examples are displayed on this slide. The CONTAINS() function searches the text index using the search criteria specified in the query and returns one when there is a match found. The QUERYLANGUAGE parameter is used to specify the locale of the search term. The RESULTLIMIT parameter is used to limit the number of results received. SCORE() returns a score of each document found indicating how well the document matches the given search criteria compared to other documents in the same column. It uses a value between zero and one.

Search XML documents

- SQL

- Find all books containing 'DB2' in title element of chapter

```
SELECT chapter FROM books
WHERE CONTAINS(chapter,
               '@xpath:''/chapter/title
               [. contains("DB2")]''') =1
```

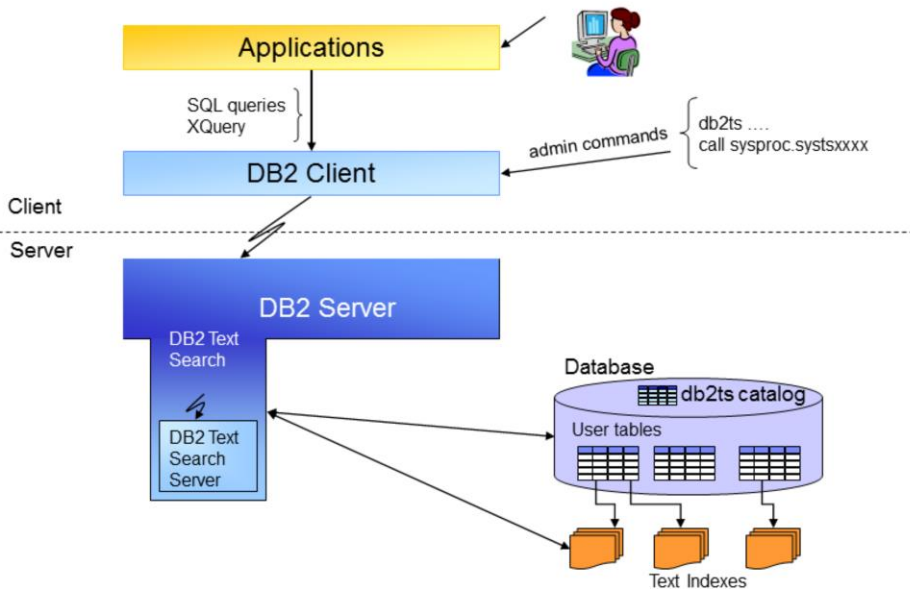
- XQuery:

- Find all books containing 'DB2' in chapter title and chapter number is 10

```
xquery for $i in db2-fn:xmlcolumn-contains
('BOOKS.CHAPTER', '@xpath:''/chapter
 [./title contains("DB2") and
 ./@number = 10] ''')
return $i
```

The first example displayed on this slide uses the XPath syntax within the SQL CONTAINS() function to search all books containing 'DB2' in the title element of a chapter. The second example uses XQuery expression with the xmlcolumn-contains() function to search all books containing 'DB2' in the chapter title and the chapter number is 10.

DB2 Text Search architecture



9

Overview

© 2013 IBM Corporation

The picture displayed on this slide represents the DB2 Text Search architecture. In this case, DB2 Server and DB2 Text Search Server are located on the same machine.

When a database is enabled for text search, catalog tables and views are created in 'SYSIBMTS' schema.

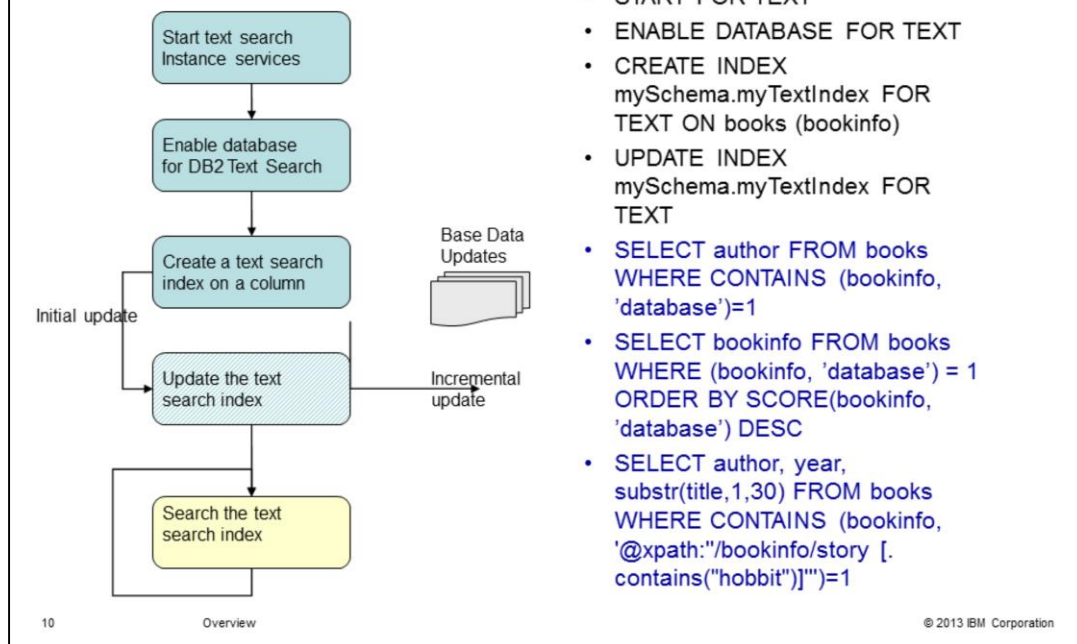
Text indexes are managed as file system objects and the low level representation of an index is known as a collection.

Once a text index is created on the user table, the staging table and log table are created.

The client applications can issue searches by way of SQL and XQuery statements.

The administration commands can be ran on the DB2 server by running the db2ts command or by calling stored procedure from the DB2 client.

Indexing and searching: A walkthrough



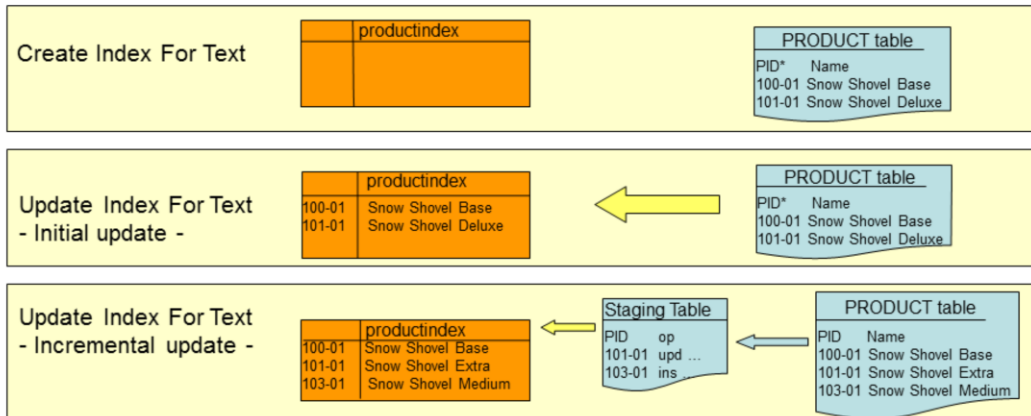
- START FOR TEXT
- ENABLE DATABASE FOR TEXT
- CREATE INDEX
mySchema.myTextIndex FOR TEXT ON books (bookinfo)
- UPDATE INDEX
mySchema.myTextIndex FOR TEXT
- `SELECT author FROM books WHERE CONTAINS (bookinfo, 'database')=1`
- `SELECT bookinfo FROM books WHERE (bookinfo, 'database') = 1 ORDER BY SCORE(bookinfo, 'database') DESC`
- `SELECT author, year, substr(title,1,30) FROM books WHERE CONTAINS (bookinfo, '@xpath:"/bookinfo/story [.contains("hobbit")]")=1`

This slide displays a view of a walkthrough example on indexing and searching. First, the DB2 Text Search instance service should be started and then, the database must be enabled so that the text search functionality is used. The text index should be created and updated so that the searches can be issued.

The first update, known as the initial index update, reads the data directly from the text column and adds it to the index.

The incremental updates run asynchronously and can be started manually or automatically through a schedule; that is, the text index update is not in the same transaction as the data update. Instead, with the default setup, for each data update an entry is inserted by way of a trigger into a staging table. This staging table is then used to control the incremental update and to synchronize the text index with the data column in the base table.

Creating and populating a text Index



11

Overview

© 2013 IBM Corporation

This slide clearly explains the concepts of initial and incremental update index. As displayed on this slide, once the index is created, it is empty, and no data is yet available. An empty collection is created by the Text Search server.

It is the initial index update which adds data about all the documents in the table column to the text index.

The subsequent changes happening on the table are tracked by triggers and captured in a staging table. Later, these updates are applied to an index by way of an incremental update.

Example: Create and update a Text Search index

- CALL SYSPROC.SYSTS_CREATE('MYSCHEMA', 'MYTEXTINDEX', 'MYSCHEMA.MYTABLE (MYCOLUMN)', '', 'en_US', ?)
 - String parameters are case-sensitive
 - message locale, for example, en_US
 - Default locale is US English

or

- db2ts "CREATE INDEX mySchema.myTextIndex FOR TEXT ON myTable(myColumn) CONNECT TO mydb"
 - Parameters are case-insensitive

- CALL SYSPROC.SYSTS_UPDATE('MYSCHEMA', 'MYTEXTINDEX', '', 'en_US', ?)

or

- db2ts "UPDATE INDEX mySchema.myTextIndex FOR TEXT CONNECT TO mydb"

DB2 Text Search provides both command-line interface and stored procedure interface to run administrative operations. This slide provides examples for creating and updating the text index. For a detailed description of options, see the DB2 Text Search User Guide.

Text Search catalog views

- Database-level views
 - SYSIBMTS.TSDEFAULTS: db-level default values for index, text and processing characteristics
 - SYSIBMTS.TSSERVERS: text server configuration data
 - SYSIBMTS.TSLOCKS: command lock info at db and index level
- Index-level views
 - SYSIBMTS.TSINDEXES: indexes and their settings
 - SYSIBMTS.TSCONFIGURATION: index configuration parameters
 - SYSIBMTS.TSCOLLECTIONNAMES: collection names
 - SYSIBMTS.TSEVENT_nnnnnn: events for each index
 - SYSIBMTS.TSSTAGING_nnnnnn: change log for each index

This slide displays the catalog views created in the database when it is enabled for Text Search. The Database-level views are SYSIBMTS.TSDEFAULTS, SYSIBMTS.TSSERVERS and SYSIBMTS.TSLOCKS.

SYSIBMTS.TSDEFAULTS shows all the default values for all text indexes in the database. For example, the default values of CODEPAGE, LANGUAGE, FORMAT and so on.

SYSIBMTS.TSSERVERS displays the Text Search server information like host name, port number, token and so on. SYSIBMTS.TSLOCKS shows the command lock information at the database and index level.

Search features

- CONTAINS, score and xmlcolumn-contains function
- SQL query, SQL xmlquery, XPath expressions, XQuery
- RESULTLIMIT
- Wildcard search
- Fuzzy Search
- Proximity Search
- Weight or boosting of score values for search terms
- Escape character
- Optional terms
- 'Fielded' search for XML documents
- N-gram and morphological indexing for CJK languages

This slide displays the search features available in DB2 Text Search. The DB2 Text Search functions available are CONTAINS, SCORE and xmlcolumn-contains. These are integrated with the DB2 optimizer to select the best possible plan.

Searches are performed by incorporating these functions in SQL query, SQL/XML query, XPath expressions and XQuery statements.

DB2 Text Search also supports wild card search, thus, increasing the size of the result set and decreasing the performance. Fuzzy search can be used where ever misspellings are possible. Proximity search is used to search for the terms located at a specified distance of 'x' words. Weight or boosting of a score value for search terms is possible. Special characters can be searched using the escape character.

Some of the query terms can be mentioned as 'optional' during the search. All XML fields and attributes can be searched. For CJK languages, it is possible to index them using the n-gram and the morphological option.

DB2 Text Search V10.1: New feature summary

- Search feature improvements
 - Fuzzy search
 - Proximity search
 - Search for special characters
 - Option for morphological indexing for CJK languages
- Support for text indexes in partitioned setups
 - Database partitioning
 - Table partitioning (Range partitions)
- Increased performance and scalability
 - Improved throughput for index update processing
 - More flexibility for workload distribution with stand-alone text server deployment option
 - More control for resource usage through additional configuration controls
- Security, Integrity and consistency
 - New system roles to manage text indexes
 - No requirement for the instance owner to have database privileges
 - Text indexes cataloged as system indexes

This slide lists a summary of new features provided with DB2 Text Search V10.1. There are new search features like Fuzzy search and Proximity search available and improvements are made to the special character search and CJK language indexing. This release also supports indexing on range partitioned tables and database partitioning environment. Due to the integration of the latest Text Server, it offers increased performance and scalability. A stand-alone server setup is supported now to have flexibility of work load distribution. New system roles have been introduced to manage the text indexes. Whenever a text index is created, an entry is created in DB2 system catalog tables as well, to prevent accidental deletion of a table on which an index is created.

Differentiators from DB2 Net Search Extender

- Integrated installation
- No additional license required
- Stored procedures for administration
- XPath-like syntax to search XML docs
- Text search from within XQuery
- Improved linguistic processing (for example, non-English)
- Text Indexes visible in system catalog
- ...

DB2 Text Search is the powerful new generation search engine offering various enhancements compared to its predecessor Net Search Extender.

Some of the key differentiators are integrated installation with DB2, no additional license required, stored procedures for administration, and XPath-like syntax to search XML documents. Also, text search can be used from within XQuery, improved linguistic processing for non-English languages, and text indexes are visible in the system catalog to prevent accidental deletion of the base table.

DB2 V10.5 committing batches for DB2 Text Search

- DB2 Text Search now provides more options for finer control of update processing
- Commit size is based on number of rows or time passed (in hours)
- Example
- `CREATE INDEX myidx1 FOR TEXT ON mytable(comment1) UPDATE FREQUENCY d(*) h(0) m(0) INDEX CONFIGURATION(UPDATEAUTOCOMMIT 3, COMMITTYPE hours, COMMITCYCLES 2)`
- `CREATE INDEX myidx1 FOR TEXT ON mytable(comment1) UPDATE FREQUENCY d(*) h(0) m(0) INDEX CONFIGURATION(UPDATEAUTOCOMMIT 300, COMMITTYPE row , COMMITCYCLES 2)`

The next few slides discuss new features available with DB2 Text Search V10.5. You can now specify the commit cycle in hours to simplify the specification of commit-batches. There are additional options available now to gracefully end the index update and continue the processing next time the update starts by enabling the use of a maximum time window for updates. Using these options, you can ensure that index processing does not run during peak workload times. The sample examples provided on this slide show how to create an index using the new options available with V10.5.

DB2 V10.5 set manual command locks for DB2 Text Search

- Use to prevent Text Index operation from interfering with non text search functions such as Backup operations
- **Example**
 - db2ts SET COMMAND LOCKS FOR INDEX i1 FOR TEXT CONNECT TO database_name user_id USING password

The command displayed on this slide is used to manually set a lock on an index or on all unlocked indexes in a database. It is used to prevent the text index operation once the lock is applied. Once the lock is applied, it is visible in the SYSIBMTS.TSLOCKS administrative view. An example is displayed on this slide for setting a manual lock on a text index.

DB2 Text Search V10.5 index configuration options

- Two new index configuration options, INITIALMODE and LOGTYPE, have been added to control update processing
 - INITIALMODE option
 - FIRST
 - Default value
 - Initial update happens on first update
 - SKIP
 - Incremental update is activated and provides first update
 - NOW
 - Initial updated is started at end of creating index operation
 - LOGTYPE option
 - BASIC
 - default
 - CUSTOM
 - No triggers are created
 - You should populate log table through any manual mechanism

A new option, 'INITIALMODE', is introduced in this release. With this option, you can run the initial update immediately when the index is created. Also, you can defer it to the first update operation or skip it altogether.

Another option, 'LOGTYPE', is introduced to set a custom log type so that the text index administrator can decide whether triggers should be added to populate the primary log table.

Additional resources

- IBM Education Assistant modules
 - DB2 Text Search - [Installation, configuration and upgrade](#)
 - DB2 Text Search - [Indexing](#)
 - DB2 Text Search - [Search features](#)
 - DB2 Text Search - [Administration commands](#)
 - DB2 Text Search - [Command line tools](#)
- IBM Information Center
 - <http://pic.dhe.ibm.com/infocenter/db2luw/v10r5/topic/com.ibm.db2.luw.admin.ts.doc/doc/c0051296.html>

This slide displays additional IBM Education Assistant modules related to DB2 Text Search and it provides the link to the IBM Information Center.

Trademarks, disclaimer, and copyright information

IBM, the IBM logo, ibm.com, DB2, and pureXML are trademarks or registered 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 other IBM trademarks is available on the web at "[Copyright and trademark information](http://www.ibm.com/legal/copytrade.shtml)" at <http://www.ibm.com/legal/copytrade.shtml>

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

THE INFORMATION CONTAINED IN THIS PRESENTATION IS PROVIDED FOR INFORMATIONAL PURPOSES ONLY. WHILE EFFORTS WERE MADE TO VERIFY THE COMPLETENESS AND ACCURACY OF THE INFORMATION CONTAINED IN THIS PRESENTATION, IT IS PROVIDED "AS IS" WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED. IN ADDITION, THIS INFORMATION IS BASED ON IBM'S CURRENT PRODUCT PLANS AND STRATEGY, WHICH ARE SUBJECT TO CHANGE BY IBM WITHOUT NOTICE. IBM SHALL NOT BE RESPONSIBLE FOR ANY DAMAGES ARISING OUT OF THE USE OF, OR OTHERWISE RELATED TO, THIS PRESENTATION OR ANY OTHER DOCUMENTATION. NOTHING CONTAINED IN THIS PRESENTATION IS INTENDED TO, NOR SHALL HAVE THE EFFECT OF, CREATING ANY WARRANTIES OR REPRESENTATIONS FROM IBM (OR ITS SUPPLIERS OR LICENSORS), OR ALTERING THE TERMS AND CONDITIONS OF ANY AGREEMENT OR LICENSE GOVERNING THE USE OF IBM PRODUCTS OR SOFTWARE.

© Copyright International Business Machines Corporation 2013. All rights reserved.

Trademarks, disclaimer, and copyright information

IBM, the IBM logo, ibm.com, DB2, and pureXML are trademarks or registered 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 other IBM trademarks is available on the web at "[Copyright and trademark information](http://www.ibm.com/legal/copytrade.shtml)" at <http://www.ibm.com/legal/copytrade.shtml>

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

THE INFORMATION CONTAINED IN THIS PRESENTATION IS PROVIDED FOR INFORMATIONAL PURPOSES ONLY. WHILE EFFORTS WERE MADE TO VERIFY THE COMPLETENESS AND ACCURACY OF THE INFORMATION CONTAINED IN THIS PRESENTATION, IT IS PROVIDED "AS IS" WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED. IN ADDITION, THIS INFORMATION IS BASED ON IBM'S CURRENT PRODUCT PLANS AND STRATEGY, WHICH ARE SUBJECT TO CHANGE BY IBM WITHOUT NOTICE. IBM SHALL NOT BE RESPONSIBLE FOR ANY DAMAGES ARISING OUT OF THE USE OF, OR OTHERWISE RELATED TO, THIS PRESENTATION OR ANY OTHER DOCUMENTATION. NOTHING CONTAINED IN THIS PRESENTATION IS INTENDED TO, NOR SHALL HAVE THE EFFECT OF, CREATING ANY WARRANTIES OR REPRESENTATIONS FROM IBM (OR ITS SUPPLIERS OR LICENSORS), OR ALTERING THE TERMS AND CONDITIONS OF ANY AGREEMENT OR LICENSE GOVERNING THE USE OF IBM PRODUCTS OR SOFTWARE.

© Copyright International Business Machines Corporation 2013. All rights reserved.