

Teradata Analyst Pack

**More Power to Analyze and Tune
Your Data Warehouse for Optimal Performance**

By:

Rod Vandervort,
Jeff Shelton,
and Louis Burger

Teradata Analyst Pack

Table of Contents

| | |
|--|-------|
| <i>Executive Summary</i> | 2 |
| <i>Introduction</i> | 3 |
| <i>Getting the Most Out of Your Data Warehouse</i> | 3 |
| <i>Teradata Optimizer and the Physical Database Design</i> | 3-4 |
| <i>Picturing the Query Plan with Teradata Visual Explain</i> | 4-5 |
| <i>Generating Production Plans on Test Systems with Teradata System Emulation Tool</i> | 5-6 |
| <i>Collecting Useful Statistics with Teradata Statistics Wizard</i> | 6-7 |
| <i>Tuning Your Physical Design with Teradata Index Wizard</i> | 7-8 |
| <i>Analyzing and Improving Application Performance with Teradata Analyst Pack</i> | 8-10 |
| <i>Taking the Pain Out of Change with Teradata Analyst Pack</i> | 10-11 |
| <i>Summary</i> | 11 |
| <i>Endnotes</i> | 11 |

Executive Summary

This paper provides a product overview of the components of the Teradata® Analyst Pack in a Teradata solution. The Teradata Analyst Pack provides Teradata customers with a suite of tools for automating and easing the task of query performance analysis. Additionally, new tools introduced in Teradata solutions help automate the difficult task of tuning the various queries in an active data warehouse to achieve better workload performance.

Teradata Analyst Pack

Introduction

In a time when the need for productivity and return on information technology investment has reached a peak, more and more companies are replacing legacy departmental database systems with an enterprise-wide data warehouse. Many companies have discovered that tracking and mining detailed data with a data warehouse can provide tremendous return on investment. They are also discovering the benefits of migrating more traditional operational applications to the data warehouse to support tactical, as well as strategic business decisions. The need for improved system performance and performance manageability grows with the increased reliance and additional application load on the active data warehouse. Application environments are expanding to include mixed workloads¹ for both decision support and near-real-time analytic processing. Managing and maximizing the performance of the database then becomes a much more difficult task.

Getting the Most Out of Your Data Warehouse

To help overcome these challenges, Teradata Corporation provides the Teradata Analyst Pack, a suite of tools to

help IT professionals analyze and tune their data warehouse for better performance. Teradata Analyst Pack is targeted for query- or workload¹-based analysis that focuses on the execution performance at an individual query level².

The Teradata Analyst Pack comprises four components:

Teradata Visual Explain – enables capturing, graphically displaying, analyzing, and comparing query plans.

Teradata System Emulation Tool – automates the collection of information from a production environment and enables emulation of that environment on a smaller test system.

Teradata Statistics Wizard – identifies and helps automate the collection or re-collection of statistics to improve performance of a given workload.

Teradata Index Wizard – automates the selection of secondary indices for a given workload.

These tools can make the jobs of Teradata database administrators, application developers, and IT professionals easier by automating and simplifying some of the difficult tasks they perform.

Teradata Optimizer and the Physical Database Design

The Teradata Optimizer is the key component enabling scalability and performance of a Teradata solution. The parallel optimizer's function is to produce the most efficient access and execution plan to retrieve the data that satisfies the user's SQL query. This plan, henceforth referred to as query plan, determines the steps, the order, the parallelism, and data access method that will most efficiently deliver the result for the specified SQL query. Determining the best query plan depends upon a number of factors, including:

- > Database physical implementation
- > Current table statistics
- > System configuration and costing formulae (e.g., speed of CPU(s), disk I/O access times, etc.)
- > Algorithms to generate alternatives of interest
- > Heuristics to cope with combinatorial explosion of the search space

Some of these factors are predetermined algorithms and heuristics inside the parallel optimizer and beyond the user's control, but many of these factors can be controlled and manipulated by modifying the physical database design and the SQL queries.

Teradata Analyst Pack

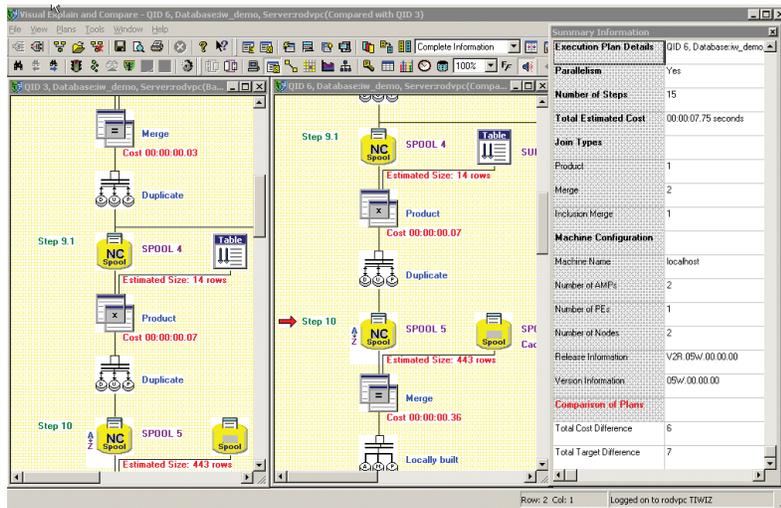


Figure 1: A visual comparison of two query plans using Teradata Visual Explain. In addition to side-by-side graphical comparison, a summary of the query plan differences displays in a pop-up window on the right.

The layout of tables, the choice of primary index, the selection and type of secondary indices, and the availability and accuracy of table statistics are all important aspects of the physical database design that must be managed and tuned by the database administrator. Because of the number and complexity of factors influencing the parallel optimizer's selection of a query plan, it can be a daunting task to analyze query performance.

The first steps in identifying opportunities to improve performance of the data warehouse are to analyze and understand the detailed steps involved in the query plan and the influences of the system configuration, data demographics³, and

index structure. Prior to the introduction of optimization tools in Teradata solutions, query performance analysis was a manual process. Typically, it required an experienced database administrator (DBA) or support analyst. You will see in the next section how Teradata Visual Explain is the first step in simplifying this task, and has been further enhanced in the latest Teradata solution release.

Picturing the Query Plan with Teradata Visual Explain

Teradata Visual Explain makes query plan analysis easier by enabling you to capture and graphically represent the steps of the query plan and to compare two or more

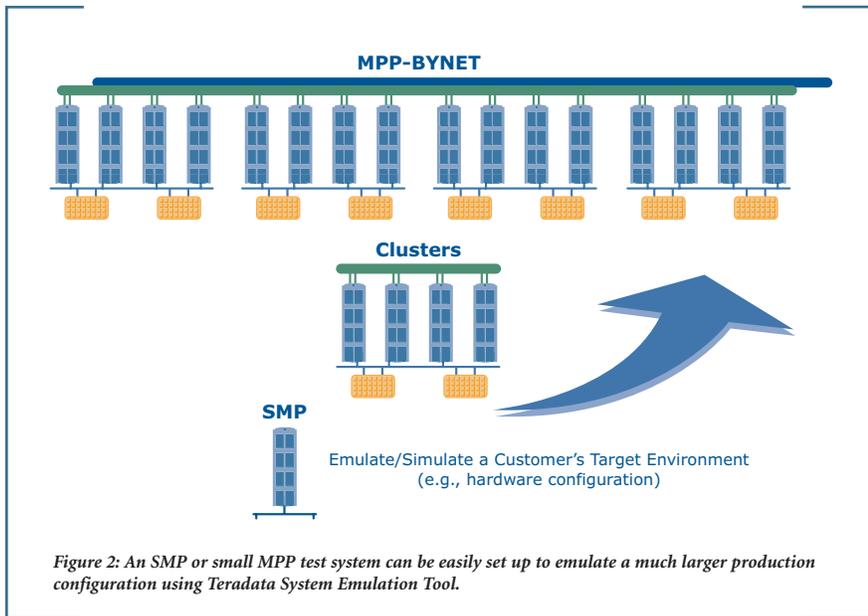
plans. Teradata Visual Explain consolidates all of the information required for query plan analysis:

- > Database object definitions (tables, views, macros, indexes)
- > Data demographics (statistics, data distribution)
- > Parallel optimizer cost
- > Cardinality estimates

Additionally, Teradata Visual Explain enables the user to submit an SQL query for execution or query plan generation and capture. It provides for graphically comparing two or more query plans side by side, highlighting any differences in the query plan, and summarizing the differences of the plans. (See Figure 1.) Detailed information about step level differences, data demographics, table definition, and index structure is easily accessed with a click of the mouse. Finally, for comparing a set of captured plans against another set of plans (e.g., comparing two workloads), there is a bulk comparison option. This option produces a summary report identifying which queries have differences. These differences can be examined for more detailed step level difference information.

Application developers, DBAs, and database support personnel can all use Teradata Visual Explain to gain a better understanding of how their SQL statements and physical database design impact the performance of a query.

Teradata Analyst Pack



The tool helps identify the performance implications of data skew, excessive data redistribution, and bad or missing statistics. Teradata Visual Explain can also capture query plans on a test system in an emulated database environment. Emulation on a test system enables the user to offload the analysis and SQL tuning tasks from the production environment. Emulation is also useful for comparing query plans for different configurations or row counts to proactively identify the impact of system expansion or table growth for a particular query. Emulation can be easily achieved using the Teradata System Emulation Tool, which is described in the next section.

Teradata Visual Explain adds the ability to compare the row count and time estimates for a query plan to the actual values captured in the Teradata Database Query Log⁴. The visualization of the query plan has been enhanced for usability, offering a new compressed view option, which depicts the plan as a set of inputs, intermediate results, and output(s), rather than displaying each procedural step of the query plan. This makes it easier to visually analyze complex queries consisting of many steps and intermediate results. Finally, support for additional data demographic information, including data skew⁵, has been added. Data skew can negatively impact parallel performance since it can cause one or a few parallel processing units to perform most of the work, creating a bottleneck for servicing the request.

Generating Production Plans on Test Systems with Teradata System Emulation Tool

The Teradata System Emulation Tool allows you to emulate a production system on a test system by enabling you to export all necessary information from the production to the test system. The parallel optimizer on the test system can then generate the same query plans as would be generated if the queries were executed on the production system. Teradata System Emulation Tool is a Microsoft® Windows®-based graphical tool that allows the user to capture system cost parameters, object definitions, random AMP samples, statistics, query execution plans, and demographics by database, by query, or by workload. This tool does not export user data, since user data is not necessary for query plan generation. This means that a small test system can be used to emulate a much larger production environment. (See Figure 2.)

Another feature of the Teradata System Emulation Tool allows the user to perform what-if scenarios relating to the data demography of the tables (statistics, random AMP samples) and system performance parameters. Upon import, the user can customize or edit object definitions, random AMP samples, statistics, and cost parameters.

Teradata Analyst Pack

A key feature of the Teradata System Emulation Tool is the ability to edit system cost parameters upon import. This is useful because it enables performing what-if scenarios for system expansions. By editing the cost parameter for number of nodes from 16 to 24 upon import to a test system, query plans can be generated as though they were generated on a 24-node system even though the production system is currently only 16 nodes. By using Teradata Visual Explain query plan capture in emulation mode and then performing bulk comparisons, you can easily identify any query plan changes that would occur if you expanded the production system by eight nodes.

Collecting Useful Statistics with Teradata Statistics Wizard

Statistics, or lack thereof, play an important role in influencing the optimizer's choice of a query plan. An important job of the DBA is to manage the collection of statistics. This ensures that the optimizer has good information about the data the queries are operating against so it can pick the best query plan. In the absence of statistics, the optimizer makes an educated guess about the data by randomly sampling a subset of the data from a single AMP, and then projecting the sampled demographics across all the AMPs for the table. Depending upon the data layout and skew, this may or may

not provide a good approximation of the full table statistics.

It shouldn't, however, be assumed that statistics should be collected on all tables and columns. Historically, collecting statistics is a manual operation that consumes human and system resources. There is a cost/benefit trade off for collecting statistics since the collection is not free, but the reward for good statistics on the right columns can be a great improvement in workload performance.

Teradata Statistics Wizard

New to the Teradata Analyst suite is the Teradata Statistics Wizard. This Microsoft® Windows®-based graphical tool helps automate the statistics management task by relieving the burden of determining when and where to collect

statistics. The result of this will be better query plans and improved DBA productivity. The Teradata Statistics Wizard can recommend a set of columns on which statistics should be collected to enable the optimizer to better estimate step costs, and thus, do a better job selecting the optimal query plan. Teradata Statistics Wizard works in two modes. The user can select a database or selection of tables for analysis using the tree view of the physical database design in the file explorer. A view of this interface is shown in Figure 3. Alternatively, the user may specify a workload to be analyzed for recommendations specific to improving the performance of the queries in the workload. This is done by analyzing the workload for indices and making statistics recommendations where appropriate.

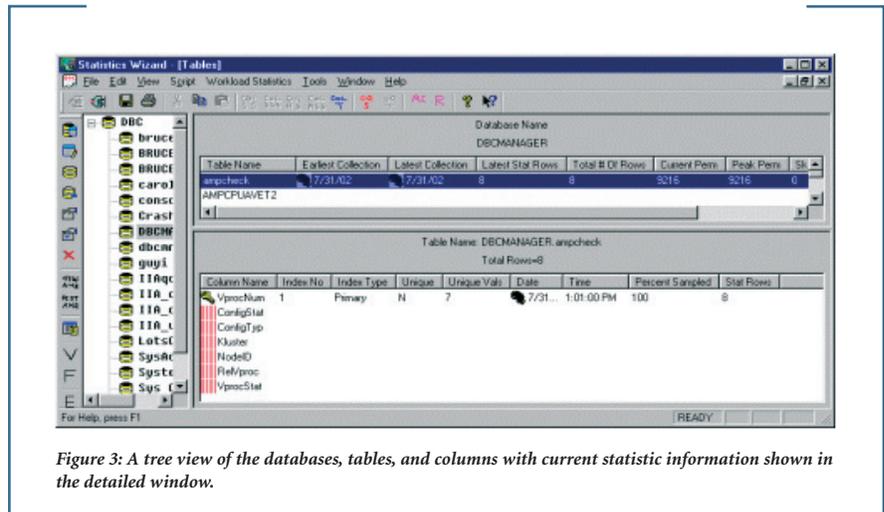


Figure 3: A tree view of the databases, tables, and columns with current statistic information shown in the detailed window.

Teradata Analyst Pack

The tool provides recommendations in the form of COLLECT STATISTICS statements. These can be executed directly from the tool or can be scheduled to execute at some point in the future.

In addition to recommending new opportunities for collecting statistics, the interface provides detailed information about where statistics have been previously collected, including the age and demographics information of the statistics collected. It is important to note that collecting statistics is a static operation; that is, the system takes a snapshot of the table and column data at the time the COLLECT STATISTICS statement is executed. It does not automatically update the statistics information when the table is updated via data manipulation (DML) statements or data load utilities. The Statistics Wizard can help DBAs keep statistics fresh by recommending re-collection when the statistics are identified as old either due to age or percentage of table growth since the last collection.

The Teradata Statistics Wizard automates statistics management, helping to ensure that appropriate columns have statistics collected and that they are kept up-to-date, ensuring the optimizer will identify the most efficient query plan.

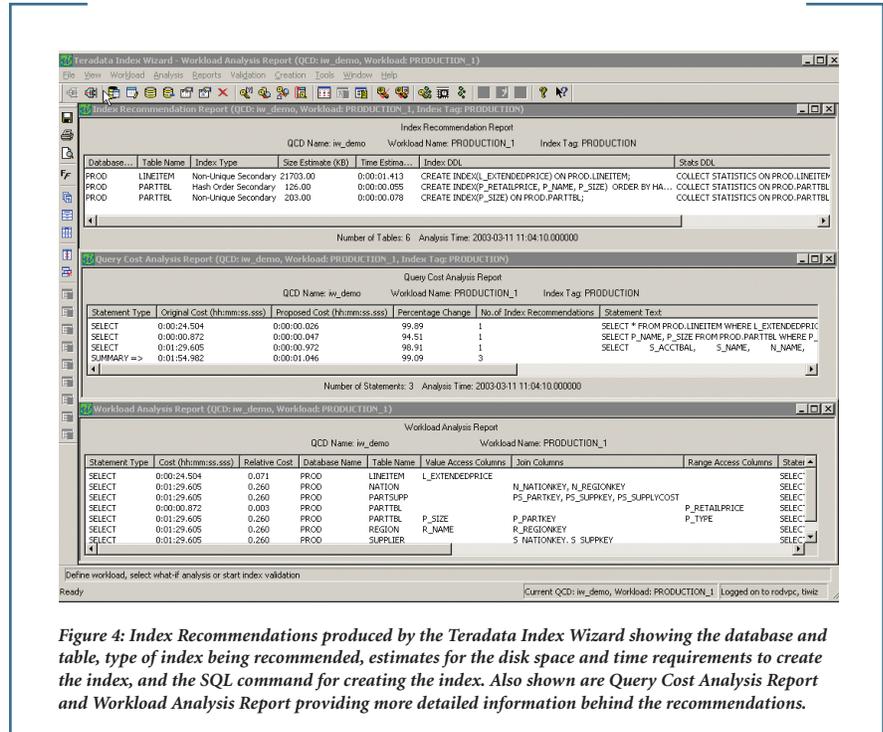


Figure 4: Index Recommendations produced by the Teradata Index Wizard showing the database and table, type of index being recommended, estimates for the disk space and time requirements to create the index, and the SQL command for creating the index. Also shown are Query Cost Analysis Report and Workload Analysis Report providing more detailed information behind the recommendations.

Tuning Your Physical Design with Teradata Index Wizard

Indexing is one of the most powerful tuning options available to a database designer or DBA. Traditionally, index selection has been a complex, manual process. It often requires the DBA to have detailed knowledge of the application workloads and data demography of the data warehouse. They must also understand parallel query plan optimization. As more application workloads are introduced to the data warehouse, analyzing the impact and derived benefit of an index becomes an increasingly difficult task. Indices come at a cost in terms of update

performance and resource utilization; thus there are cost/benefit trade-offs to consider when selecting a proper set of indices. To determine a good set of indices for the data warehouse, designers mostly rely on their application experience and intuition to make index design decisions.

Teradata Index Wizard is a new tool that automates the process of index design by recommending secondary indices for a particular workload. Teradata Index Wizard provides an easy-to-use graphical interface to assist the user in analyzing a workload, and provides recommendations for improving performance through the use of indices.

Teradata Analyst Pack

Teradata Index Wizard is available in Teradata Warehouse 7.0. Teradata Index Wizard consists of a database server component and a front-end client application. The index analysis engine is inside the Teradata parser, and works closely with the parallel optimizer to enumerate, simulate, and evaluate index selection candidates. The client front-end is a graphical Microsoft Windows Interface that provides step-by-step instructions for workload definition and index analysis. It also provides several reports for both workload and index analysis. Figure 4 shows the Index Recommendation Report, the Query Cost Report, and the Workload Analysis Report produced by Teradata Index Wizard.

Teradata Index Wizard works in two modes. In wizard mode, it contains menus that guide the user through the steps of identifying a workload, and then identifying and implementing a set of secondary indices that would help improve the performance of that workload. The typical wizard mode steps are:

1. **Workload Definition** – Specify the set of queries to consider for analysis
2. **Workload Analysis** – Identify the relevant tables, columns, and existing indices
3. **Index Analysis** – Enumerate candidate indices, simulate and measure impact of the candidates, and select the index recommendation set

4. **Analysis Reporting** – View summary and detailed information about the workload, the existing physical design, the proposed changes, and the cost/benefit of the recommendations relative to the workload
5. **Recommendation Validation** (optional) – An additional level of simulation and impact measurement with up-to-date statistics collected for all recommendations; particularly useful when Index Analysis is done on a test system using emulation
6. **Recommendation Implementation** (optional) – Tune the physical database design by creating the recommended secondary indices and collecting statistics on the newly created indices

The second mode supported by Teradata Index Wizard is the what-if analysis mode. This mode allows you to provide a set of recommended indices. The user-provided set of indices is then simulated, and a report is provided showing the overall estimated improvement to the workload if the proposed indexes were implemented. It also provides query-by-query details indicating the improvement for each query in the workload, and whether or not the simulated index was utilized in the new query plan.

Teradata Index Wizard can be used in conjunction with Teradata System Emulation Tool to enable index analysis on a test system. Both wizard and what-if

modes help automate the task of secondary index selection, increasing the DBA's productivity.

Analyzing and Improving Application Performance with Teradata Analyst Pack

Imagine your application is not performing up to expectation and users are complaining that response times are too long. Does it mean that you need to buy a bigger, more powerful database system? It might, but before you invest in more hardware, it is important to isolate and identify the issue. It could be that with a simple change to SQL, the addition of secondary indices, or just by collecting statistics on the right columns, you could dramatically improve performance without adding additional system resources.

The starting point for any analysis is the identification of which queries are contributing to the problem. The Teradata Database Query Log feature enables detailed query performance analysis by logging SQL, query start and stop times, and even step level information. This is helpful when trying to diagnose overall system performance issues in a complex active data warehouse with thousands of users and many applications running concurrently. Since query logging can log at either the user or account level, it is useful for tracking query performance for individual applications or sets of users.

Teradata Analyst Pack

This helps DBAs or IT managers identify the queries submitted by individual applications or groups of users.

Once the query or set of queries is identified, Teradata Visual Explain may be used to capture the optimizer's query plan for each query. This shows cost estimates for each step. Many times, the cause of the bad performance is due to excessive data movement to collocate data for table joins. This can happen when unnecessary or inefficient joins on large tables are required due to physical database design (choice of table layout, choice of primary and secondary indices), or due to poor estimation of row counts caused by stale or nonexistent statistics.

As a simple example of how index selection and statistics can affect performance consider the following query:

```
SELECT * FROM t1 WHERE c1 = 10;
```

Assume that table t1 has a secondary index c1. Even in this very simple query, the optimizer has a choice of query plans. Since there is an index defined on c1, one choice is whether to use index lookup or perform a full scan of all the rows in the table. With statistics available, the optimizer will determine approximately how many rows meet the condition of c1=10. If there are very few rows meeting this as compared to the total number of rows in the table, it is more efficient to use the

index lookup access. If there are many rows with a value of c1=10, it will likely be faster to just scan the whole table and filter out the rows which don't meet the condition. If we have no statistics, or worse yet, bad statistics, the parallel optimizer might make a bad estimate on the number of rows meeting the condition of c1=10 and pick the less than optimal query plan. So even in this simple example, you can see how the presence of a secondary index, and also the presence (or absence) of statistics, can influence the parallel optimizer's choice of a query plan.

Teradata Visual Explain is helpful in diagnosing individual query plans as in the example just described. An experienced DBA or support analyst can recognize opportunities for query plan improvement by collecting statistics on certain columns or introducing a secondary index to provide a possibly more efficient access method. In the cases where several queries must be analyzed as a workload, the one-query-at-a-time analysis can become a tedious and more complex job. In this case, Teradata Statistics Wizard and Teradata Index Wizard can be used to auto-matically perform the analysis for the workload and recommend the right set of physical database design tuning options that would improve the overall workload performance.

A summary of the analysis steps that may be performed is:

1. Identify the queries or workload to be analyzed for performance tuning. Logging to Teradata DBQL is one easy way to selectively capture SQL queries at query run time.
2. Capture the query plans in question into the Query Capture Database. For individual queries or a small set of queries, this can be done directly from the Teradata Visual Explain Launch QCF feature. For bigger workloads, Teradata Index Wizard Workload definition feature can be used to define the workload and capture query plans for all queries in the workload.
3. Optionally export the captured query plans, system cost parameters, statistics, and object definitions from the production or target system to a test system using Teradata System Emulation tool. This step is not necessary, but enables offloading the analysis activities from the production to the test environment.
4. Analyze the statistics and data demographics for the tables involved in the workload using Teradata Statistics Wizard. Analyze the workload for statistics recommendations. Execute the recommendations. At this point, you can optionally recapture the plans for the workload and compare them using Teradata Visual Explain's query

Teradata Analyst Pack

or bulk compare feature. Doing this allows you to visualize and measure the improvements achieved by implementing the Teradata Statistics Wizard recommendations.

5. Perform Index Analysis for the workload using Teradata Index Wizard. If the analysis is run on a test system, the index recommendations can be validated back on the production system using the Validate Recommendations feature. This will simulate the indices without actually creating them, and produce a query plan as though the indices were present. Validation also collects sampled statistics as part of the process; thus ensuring up-to-date statistics are available for the proposed indices.
6. Analyze the output of the Teradata Index Wizard reports to understand the recommendation's impact on performance (estimated performance improvement), and cost to implement (in terms of disk space and creation time). Utilize Teradata Visual Explain to compare pre- and post recommendation plans.

Post-recommendation plans are automatically generated and stored for each query in the workload during index validation. Implement the recommendations on the production system to achieve the performance improvement desired, or alternatively, use the what-if analysis

mode to identify potential performance improvements of your own proposed secondary indices.

We have seen how these tools can be utilized to help you isolate query performance problems, and provide the necessary tuning to the data warehouse to achieve performance enhancements for an application or workload. Next we'll discuss how you can use the tools to help ensure the data warehouse delivers consistent performance as changes occur over time.

Taking the Pain Out of Change with Teradata Analyst Pack

As previously mentioned, Teradata Visual Explain provides a graphical visualization of the query plan selected by the optimizer. It is also powerful in analyzing two or more query plans, with automated features for individual plan comparisons and also for bulk comparisons, where a set of query plans (i.e., a workload) can be compared with another set of plans for the same workload.

You can use Teradata Visual Explain, combined with Teradata System Emulation Tool, in change control planning to identify any impact a proposed change has on the optimizer's selection of the query plan for each query in the application workload. Using these tools together

with a test system environment, it is possible to emulate:

- > System expansion (editing the system cost parameters to include additional nodes / AMPS)
- > Database table growth (editing statistics for tables that are expected to grow significantly over time to emulate the fully populated tables)
- > Database software upgrade(s)

In each of these cases, the Teradata System Emulation Tool is used to capture the production environment as defined by its Database Object Definitions, Query Plans, Statistics, Cost Parameters, and Random AMP Samples. This information is imported to the test system. Because Teradata System Emulation Tool does not capture any user table data, the test system can be much smaller than the production system. Once the production system information is imported to the test system, the proposed change can be emulated on the test system. In the case of system expansion or table growth, these can be accomplished by editing the appropriate cost parameters or table statistics using the Teradata System Emulation Tool (e.g., change a 100,000 row table to contain 10,000,000 rows). In the case of a software upgrade, the new version of software can be installed on the test system (e.g., upgrade the test system from Teradata Database V2R4.1 to Teradata Database V2R5.0). After the appropriate

Teradata Analyst Pack

data migration (if any, as determined by the software upgrade process), and after editing the appropriate emulation parameters in the expansion or growth emulation cases, a new set of query plans can be captured using the Teradata Visual Explain tool. Once the new set of plans is captured, Teradata Visual Explain can perform a “bulk compare” to identify any differences between the query plans prior to the change. The plans with differences can be graphically viewed with Teradata Visual Explain to determine the nature of the difference and determine what impact, if any, the proposed change will produce. In this scenario, the tools are used to proactively identify and address any potential impacts to performance of the production system prior to implementing the proposed change.

Summary

The components of the Teradata Analyst Pack are:

Teradata Visual Explain – a Windows-based graphical tool for capturing, displaying, and comparing query plans.

Teradata System Emulation Tool – a Windows-based client tool for exporting system configuration, physical database design, and data demographic information, and, for enabling an emulation, on a much smaller test system, of the database environment.

Teradata Statistics Wizard – a Windows-based graphical tool for analyzing and managing statistics collections. This includes generating recommendations for new collections or re-collection of outdated statistics.

Teradata Index Wizard – a Windows-based graphical tool for automating the second-ary index selection process, or alternatively, enabling the DBA to simulate and test his own recommendations using the what-if analysis feature.

In addition to providing an overview of the Teradata Analyst Pack, examples of practical application of the tools have been offered for the cases of application performance tuning and change control performance management. Teradata Analyst Pack simplifies the complexities of performance management by providing easy access and consolidating the important data necessary for performance analysis. Also, the new Teradata Statistics and Teradata Index Wizards automate analysis and tuning of the physical design for improved application performance. In an ever growing, complex active data warehouse environment, the Teradata Analyst Pack enables you to successfully analyze, tune, and manage performance, helping you to get more out of your data warehouse.

Rod Vandervort is the Development Manager of Teradata Optimizer Tools group and oversees the development of the Teradata Analyst products.

Jeff Shelton and **Louis Burger** are software engineers in the Teradata Optimizer Tools group, focusing on developing database features to facilitate the process of query plan analysis and performance tuning.

¹ A workload is a set of SQL queries.

² This does not consider the impacts of resource contention due to concurrency. Concurrency, system resource bottlenecks, and overall system throughput analysis are best addressed by system resource level monitoring tools, such as Teradata Manager, and workload management tools, such as Dynamic Query Manager and Priority Scheduler.

³ Data demographics refers to information which describes the data such as number of rows in a table, average size of the row, row distribution across the Access Module Processor (AMP) units of parallelism

⁴ Teradata Database Query Log is a new feature in Teradata Database V2R5.0 that enables selective query logging including User / Acct. information, SQL text, start and completion time, and step time and row counts.

⁵ Data skew refers to an uneven or non-uniform distribution of data.