

# The Oracle Warehouse

Gary Hallmark  
Oracle Corporation

Data warehouses are repositories that integrate and summarize historical and reference data from numerous sources. Warehoused data can be analyzed along several dimensions such as time, product, and geography to identify trends and gain competitive advantage.

Data warehousing is one of the fastest growing segments of the computer industry. Industry analysts estimate the current annual warehousing market at nearly \$2 billion and project growth to \$6.9 billion in 1999.

Up until now, businesses have used database technology and on-line transaction processing to track the flow of goods, people, and money. With data warehousing, companies can now analyze that same operational data to increase efficiency, achieve economy of scale, and actually make money.

## ***Any Source, any Data, any Tool***

The warehouse should accept data from any source, store it in any format, and present it to any tool.

The Oracle Warehouse supports "any source, any data, any tool" through a combination of products, services, and partnerships. This paper focuses on Oracle Warehouse products.

## ***Data Sources***

Common sources of record-oriented warehouse data are financial applications, order entry applications, credit reports, and mailing lists. Such data can be loaded into the warehouse in a number of ways:

- A number of Oracle's partner companies specialize in transforming external data and feeding it into Oracle's high speed parallel loader. The transformations can be very exotic. This is a good approach when the external data needs extensive cleaning, categorizing, or restructuring when moved into the warehouse.
- Oracle develops a family of SQL gateways to

relational and non-relational data stores. Because the warehouse integrates gateways and distributed SQL processing, the warehouse can load data from an operational system with an SQL *insert...select* statement.

- Oracle supports asynchronous replication between operational databases and the warehouse database. Replication can be periodic or event driven. Replication can be at the table level or at the procedure level. Procedure level replication can translate between possibly different schema structures.
- A common source of data is the warehouse itself. Many operations build specialized indexes or summaries to enable faster access to terabytes of data, or more efficient access for multiple users.

## ***Data Types***

The core of a data warehouse is an industrial strength relational database management system. Because only a fraction of company data is record-oriented, the warehouse must accommodate textual, spatial, and multimedia data as well.

## ***Record-oriented data***

Support for record-oriented data includes traditional query optimization and parallel execution techniques, and some non-traditional techniques such as bitmap indexes and Cartesian product joins.

## ***Optimization***

Oracle has two optimizer stages. The first stage flattens views and subqueries, and pushes predicates down. The second stage is cost based, trying all permutations of left-deep join orders and access paths. The optimizer considers Cartesian products because many warehouses use "star" shaped schemas.

Recent enhancements give the user the ability to trade resource for response time when optimizing for parallel execution, and provide the optimizer with detailed histograms on sets of columns.

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the VLDB copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Very Large Data Base Endowment. To copy otherwise, or to republish, requires a fee and/or special permission from the Endowment.  
Proceedings of the 21st VLDB Conference  
Zurich, Switzerland, 1995

## Parallel Execution

Several Oracle warehouses have over a terabyte of record-oriented data. Within two years, warehouses will be ten terabytes in size. This amount of data requires parallel processing. Oracle's Parallel Query Option is production on all major shared memory, shared disk, and shared nothing hardware platforms.

The parallel query option provides parallel loading, index creation, summary creation, and complex query execution. All operations may optionally disable logging. This is appropriate when the storage subsystem is itself fault-tolerant.

A data warehouse often contains summary tables such as "data cubes" for drill down analysis. A parallel version of *create table as <subquery>* creates summary tables quickly and efficiently.

Parallel execution of complex queries provides scalable implementations of table scans, index probes, aggregation and grouping, duplicate elimination, hash, sort-merge, and nested loop joins, unions, anti-joins, and user-defined function invocation.

Two "new" techniques for warehousing in a traditional relational database are bitmap indexes and Cartesian products.

## Bitmap Indexes

Bitmap indexes (introduced in Model 204) provide a segmented bitmap per distinct column value. Bit  $i$  is 1 when the value appears in row  $i$ 's column.

Bitmap indexes are most useful for columns with just a few distinct values (e.g. male/female). Warehouse schemas sometimes map continuous data to discrete ranges (e.g., salary ranges), making columns with few distinct values common. When the number of distinct values increases, the number of bitmaps increases and the number of ones in each bitmap decreases. Thus, bitmap compression is important. There is a tradeoff in amount of compression and efficiency of updating the index.

The efficiency of bitmap indexes is greatest when used to return only the count of matching rows. When more complex aggregation (e.g. moving averages) or actual detail data must be returned, the optimizer will often choose other access paths, such as full table scans.

All indexes, including bitmap and B-tree, have limits as access paths for general predicates. For example, a query involving a leading wildcard character will often use a full table scan.

In addition to record-oriented retrieval, bitmaps are useful for text retrieval. Oracle TextServer uses a bitmap per word. Bit  $i$  is on when the word is in document  $i$ .

## Cartesian Products

A common warehouse schema consists of a single large "fact" table, indexed by concatenated foreign keys that refer to dimension tables. As a simple example, consider a revenue table where each row contains the revenue for one item in one state on one day. To answer the query: *find the total revenue for umbrellas sold on sunny March weekends in the West*, it may be better to first form the Cartesian product of the item table, the state table, and the day table, and then access the fact table.

In some cases, bitmap indexes can be used instead of, or in addition to, Cartesian product joins. Bitmaps are most useful when the number of values in each dimension is small. Cartesian products are useful when the number of values in each dimension, after applying the query predicates, is small.

## Text Data

A large portion of corporate information assets is in the form of unstructured text. Examples of documents that can be integrated into the warehouse include data sheets, white papers, press releases, annual reports, legal contracts, and design specifications.

The Oracle Warehouse support for text extends traditional information retrieval capabilities. Word content searches use bitmap indexes. Documents are stored in compressed form, and proximity searches work on the compressed form.

ConText, a natural language parser and semantic network, performs thematic searches and summarization. One may search for articles about sports containing the word cricket and avoid paging through screens of articles about insects. Summary mode omits phrases in documents that do not contribute to the central theme.

## Spatial Data

Many companies have sales, support, or other operations spread over several regions. These companies can benefit from integrating geographical information with business data.

The Oracle Warehouse support for spatial data subsumes geographic information system capabilities. A new *HHcode* index clusters points in  $N$ -space. The *HHcode* recursively divides  $N$ -dimensional cubes into  $2^N$  subcubes.

The location of an object in  $N$ -space is given by the HHcode of the smallest containing subcube. The HHcode of a subcube is the concatenation of the containing subcube numbers.

## Multimedia Data

Finally, organizations have large volumes of data in the form of videotaped or audiotaped presentations, interviews, training classes, news clips, and advertisements that must be cross-referenced with other data in the warehouse. For example, what commercials aired 24 hours before the last three spikes in sneaker sales?

The Oracle media server stores and delivers broadcast quality video and CD quality audio.

## Tools

Tools are specialized for *ad hoc* browsing of multidimensional data or for building canned executive information system applications.

Tools for *ad hoc* access to the warehouse often have a familiar spreadsheet-like interface. The tools are more integrated with the database than a conventional spreadsheet, and have more sophisticated analysis functions for curve-fitting, trending, and extrapolation.

Oracle's Discoverer/2000 is an *ad hoc* warehouse browser with spreadsheet-like crosstabs and bar graphs. Drill down is accomplished by double-clicking on the cell or bar that represents an aggregate that is to be expanded. The tool generates SQL on-the-fly and retrieves data on demand, reducing client-side memory needs.

Oracle's IRI ExpressView is an *ad hoc* browser that analyzes existing warehouse data and extrapolates existing data into the future. The tool provides many different extrapolation models to enable "what-if" analysis. ExpressView interacts with warehouse data through an intermediate multidimensional cache, called ExpressDB. Many ExpressView clients share a common ExpressDB server.

Tools for development of EIS applications provide professional programmers with reusable components (often object-oriented) like forms, charts, and drop-down lists to isolate end users from the warehouse implementation details.

Oracle's Developer/2000 is a general purpose graphical database application development environment. It can be coupled with Designer/2000 to provide a formal CASE development methodology.

Oracle's IRI ExpressBuilder is an object-oriented development environment tailored to build trend analysis and forecasting applications. ExpressBuilder also works with ExpressDB as a shared cache for multidimensional warehouse data.

Many additional tools are available from partners. Most warehouse tools support summaries at various levels of detail. Some tools store summaries on the client machine and do not share them across users. Other tools use a specialized server for caching multidimensional summaries (e.g., IRI ExpressDB), or may store summaries as tables or snapshots in the RDBMS.

## Conclusion

Relational database technology, increased hardware power, and competition among hardware vendors have combined to drop the price/performance ratio of transaction processing to a point where even small companies can afford to automate their operations.

The same market forces that revolutionized transaction processing will make warehousing affordable to a wide range of organizations. The competitive advantage of consolidating and analyzing all corporate data will be so great that companies will not be able to afford to be without a data warehouse.