



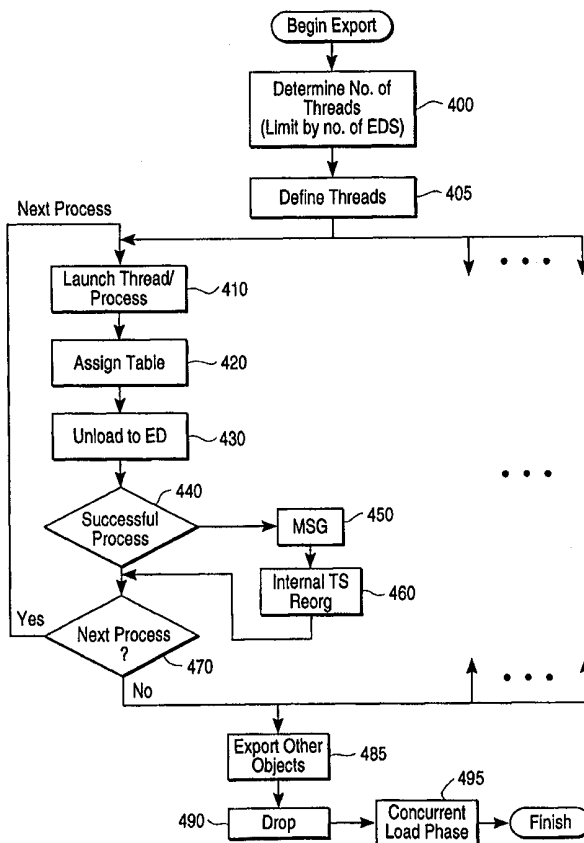
INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

<p>(51) International Patent Classification <sup>6</sup> : <b>G06F 9/46, 12/00, 17/30</b></p>	<p><b>A1</b></p>	<p>(11) International Publication Number: <b>WO 00/31635</b> (43) International Publication Date: 2 June 2000 (02.06.00)</p>
<p>(21) International Application Number: PCT/US99/27835 (22) International Filing Date: 24 November 1999 (24.11.99) (30) Priority Data: 09/200,254 25 November 1998 (25.11.98) US (71) Applicant: COMPUTER ASSOCIATES THINK, INC. [US/US]; One Computer Associates Plaza, Islandia, NY 11749 (US). (72) Inventor: MIRZADEH, Rosita; 18440 Hatteras Street #37, Tarzana, CA 91356 (US). (74) Agents: FLIESLER, Martin, C. et al.; Fliesler, Dubb, Meyer and Lovejoy LLP, Suite 400, Four Embarcadero Center, San Francisco, CA 94111-4156 (US).</p>		<p>(81) Designated States: AE, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, CA, CH, CN, CR, CU, CZ, DE, DK, DM, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TR, TT, TZ, UA, UG, UZ, VN, YU, ZA, ZW, ARIPO patent (GH, GM, KE, LS, MW, SD, SL, SZ, TZ, UG, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GW, ML, MR, NE, SN, TD, TG).</p> <p><b>Published</b> <i>With international search report.</i> <i>Before the expiration of the time limit for amending the claims and to be republished in the event of the receipt of amendments.</i></p>

(54) Title: METHOD AND APPARATUS FOR CONCURRENT DBMS TABLE OPERATIONS

(57) Abstract

Database tables are unloaded by launching a number of threads (400) corresponding to a number of export directories (320) located on separate storage devices (310) that maintain the database tables. Each thread is assigned a database table to unload (420). Data unloaded from each database table is stored in a corresponding export directory (320). The data is unloaded from each database table by reading blocks of data from each table and storing the data logs in the export directory (320). Each thread is handled by a separate process in a Symmetric Multi-Processing (SMP) environment. The process is repeated until each database table has been unloaded. The data is then loaded into database tables by first creating a number of temporary tables corresponding to the number of threads (600), reading a set of data stored in the export directory (320) and storing the data in the temporary tables.



**FOR THE PURPOSES OF INFORMATION ONLY**

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

AL	Albania	ES	Spain	LS	Lesotho	SI	Slovenia
AM	Armenia	FI	Finland	LT	Lithuania	SK	Slovakia
AT	Austria	FR	France	LU	Luxembourg	SN	Senegal
AU	Australia	GA	Gabon	LV	Latvia	SZ	Swaziland
AZ	Azerbaijan	GB	United Kingdom	MC	Monaco	TD	Chad
BA	Bosnia and Herzegovina	GE	Georgia	MD	Republic of Moldova	TG	Togo
BB	Barbados	GH	Ghana	MG	Madagascar	TJ	Tajikistan
BE	Belgium	GN	Guinea	MK	The former Yugoslav Republic of Macedonia	TM	Turkmenistan
BF	Burkina Faso	GR	Greece	ML	Mali	TR	Turkey
BG	Bulgaria	HU	Hungary	MN	Mongolia	TT	Trinidad and Tobago
BJ	Benin	IE	Ireland	MR	Mauritania	UA	Ukraine
BR	Brazil	IL	Israel	MW	Malawi	UG	Uganda
BY	Belarus	IS	Iceland	MX	Mexico	US	United States of America
CA	Canada	IT	Italy	NE	Niger	UZ	Uzbekistan
CF	Central African Republic	JP	Japan	NL	Netherlands	VN	Viet Nam
CG	Congo	KE	Kenya	NO	Norway	YU	Yugoslavia
CH	Switzerland	KG	Kyrgyzstan	NZ	New Zealand	ZW	Zimbabwe
CI	Côte d'Ivoire	KP	Democratic People's Republic of Korea	PL	Poland		
CM	Cameroon	KR	Republic of Korea	PT	Portugal		
CN	China	KZ	Kazakstan	RO	Romania		
CU	Cuba	LC	Saint Lucia	RU	Russian Federation		
CZ	Czech Republic	LI	Liechtenstein	SD	Sudan		
DE	Germany	LK	Sri Lanka	SE	Sweden		
DK	Denmark	LR	Liberia	SG	Singapore		
EE	Estonia						

- 1 -

**Method and Apparatus for Concurrent  
DBMS Table Operations**

5                   Background of the Invention

Field of the Invention

This invention relates to a process of unloading and loading a database table. The invention is more particularly related to the application of parallel processing techniques to increase the speed at which database tables are unloaded and loaded. The invention is still further related to parallel processing techniques applied to unloading and loading database tables in a Symmetrical Multi-Processing (SMP) environment.

Discussion of the Background

Modern database management systems are increasingly called upon to maintain larger stores of data. In addition to the increased size of databases, the structure within modern databases is increasingly complex.

Typically, a database maintains data in the form of tables, each table maintaining one or more rows of related data. As an example, a basic database table may maintain plural rows having, for example, name, social security number, address and telephone number of individuals belonging to an organization.

The database would be increased in size as the organization adds new members, and would increase in both size and complexity as additional information about members is included. For example, a larger and more complex database could maintain, in addition to the above information, a map, perhaps in graphical format, showing the club members' residence. The

- 2 -

database could again be increased in size and complexity by including a work address and an additional graphical map showing the location of the work place.

5           The database table may again be increased in complexity by maintaining pointers to other tables or rows of the database. For example, a pointer to a set of coworkers maintained in another table, pointers to nearby organization members, or a pointer(s) to any  
10       number of items to where additional member data may be maintained.

          Conventional Database Management Systems (DBMS) provide space for building database tables by allocating blocks. Once a table is defined, the DBMS  
15       will allocate blocks necessary for storing rows of the related data. For example, if a table is to be built to contain 100,000 rows, and the rows of the table were defined such that 100 rows could fit per block, the DBMS would allocate 1,000 blocks to build the table.

20           Generally, DBMS systems allocate blocks in sets of contiguous blocks. A contiguous set of allocated blocks is commonly referred to as an extent. As a general rule, extents vary in size. Using the above example, the DBMS may utilize a single extent of 1,000  
25       blocks, 2 extents of 500 blocks, or any other combination of extent sizes to allocate the required 1,000 blocks to build the table. Once the required blocks are allocated, the data is then stored in rows in the table utilizing the blocks allocated.

30           Over the course of time, additional data may be added to the table and the DBMS will allocate additional blocks as required. For example, if a user adds 250 rows to the table, using the above parameters, an additional 3 blocks would need to be allocated.

- 3 -

Also over the course of time, information in the database may be deleted. In this case, deletion of rows occurs through the use of SQL to delete rows from the table. For example, a user may delete fifty rows  
5 from block 1, 40 rows from block 20, and 30 rows from block 60. When this occurs, it causes there to be more blocks in the table than required to maintain the data present in the table.

In addition, data within the database will be  
10 updated. For example, using the above-described database tables, a organization member who has not yet entered the workforce would have a row in the table allocated that maintained his/her name, address, social security and telephone number. Upon entering the  
15 workforce, the row would be updated to include the work address and related information. However, if a substantial amount of information is added as a result of the update, the original row may not have enough blocks of data allocated to maintain the updated  
20 information.

Row migration occurs when a row gets updated and the original row does not have enough space to hold all of the updated information. When this occurs, the row is moved to a location with more space, and a pointer  
25 is placed in the block where the original row resided, the pointer being set to point to the location of the moved (migrated) row. A large amount of row migration is caused if there are substantial updates to a table and an inadequate amount of space was allocated for the  
30 original rows.

More often than not, there is insufficient space within a block to hold an updated row. In this case, the row is migrated to an entirely different block than where the original row resided, and the pointer is

- 4 -

placed in the original row position points to the updated row in the different block.

Anytime a row is migrated it causes what is known as fragmentation. Fragmentation causes greatly  
5 increased retrieval time of database information because in addition to reading a block, a pointer must be read and interpreted. When a row is migrated to another block, at least two blocks (the block containing the pointer, and the block containing the  
10 migrated/fragmented row) must be read to retrieve row specific information. Other structural changes within the database tables also cause fragmentation and related efficiency problems (row chaining, for example).

15 From time to time, the Database Administrator (DBA) will perform an analysis on the DBMS tables that provides information regarding the condition of the tables. For example, the database administrator may look at information regarding the number of deleted  
20 rows to ascertain efficiency information with respect to how many blocks are holding deleted rows. As another example, the Database Administrator may look to see how many rows of a table have been migrated or fragmented by other processes.

25 If a lot of fragmentation has occurred, it indicates that block size and row space could be more efficiently allocated and that table data is not being retrieved efficiently. When this occurs, the database administrator will likely decide to rebuild the tables.

30 When creating a table, the DBA makes a decision regarding the structure of a database table by setting a percentage of blocks free (PCTFREE) or percentage of blocks used (PCTUSED). As the DBMS fills up each block with row or table information, it will keep a

- 5 -

percentage of a block free at least equal to the percentage PCTFREE.

The DBA sets the PCTFREE variable depending on how the database table is to be used. For example, if a table is to have frequent updates, additional PCTFREE would be established so that enough space is available to allow any necessary row migration to occur within the same block. As discussed above, row migration within the same block does not cause a table to become fragmented. Migrated, but non-fragmented rows are retrieved with a single block read rather than the cumbersome process of reading a block, interpreting a pointer, and reading a second block (or more) as occurs when migrated rows are fragmented. Therefore, appropriate PCTFREE settings allow DBMS performance to be maintained although the database tables may be modified.

PCTUSED is another parameter that allows the DBA to control the structure of a DBMS table. The DBMS prevents additional rows to be placed in a block unless the percentage of that block has fallen below PCTUSED. PCTUSED is different from PCTFREE in the sense that although a block may be used if there is an update, it will not be used to insert a new row unless the percentage used in the block is below PCTUSED.

A DBMS table involved in heavy OLTP activity (inserts, updates and deletes) over time will likely experience row migration, fragmentation, row chaining, etc. Furthermore, various database tables may not necessarily have appropriate settings (PCTFREE, PCTUSED, for example) when first built, or the needs of the database table may have changed, resulting in additional migration, deletion or fragmentation of

- 6 -

tables. This in turn results in a degradation of data retrieval performance and space usage.

The DBA will perform an analysis to determine whether the tables are storing data efficiently. As a  
5 result, one or more of the DBMS tables may be determined to being inefficient in storing and retrieving data. Reorganization (rebuilding) of the table is a solution to this problem. In order to achieve maximum performance the table needs to be  
10 rebuilt (i.e., the data unloaded into secondary space and a fresh instance of the table rebuilt). This process gets rid of many of the unwanted effects mentioned above because the fragmented rows are unloaded and stored without fragmentation in the  
15 rebuilt table.

Structurally sound databases make efficient use of disk space. They require less time to access data, reduce the time required for normal transactions, and provide better response time to the user. Even though  
20 Oracle and other modern database systems use efficient logic for data placement, normal activity over time causes the physical layout of the data on disk to degrade and space requirements to grow. This results in excessive space usage and extra time needed to  
25 perform table scans, database backups, and other functions. Partial empty pages and unused extent space contribute to the additional space usage. Also, nonsequential rows and extent interleaving seriously degrade performance if they are not resolved  
30 periodically.

One method to ensure that databases stay efficient (increasing productivity) is to regularly perform reorganizations on the databases' data. Currently, products are available to reorganize DBMS tables.



- 7 -

However, even with automated reorganization tools, reorganization of database tables can require substantial amounts of time. The time required to perform a reorganization can have substantial effects  
5 on revenue or productivity of a shop that requires database access. For example, some shops can incur costs of approximately \$100,000 an hour when a database is off-line (See Dec. 1997 issue of Oracle magazine, "Reorgs in a Non-Stop Shop," for example). Therefore,  
10 any improvement in efficiency or speed at which a reorganization is performed would increase competitiveness and profitability.

#### Summary of the Invention

15 The present inventor has realized the need to provide fully parallel operations supporting database table unloading and loading which will increase the speed of any database operations performing either one or both of table unloading and table loading.

20 Accordingly, it is an object of the present invention to provide a parallel processing technique that allows for parallel database table unloads.

It is another object of this invention to provide a method for unloading and loading database tables  
25 utilizing parallel processing techniques in a Symmetric Multi-Processing (SMP) environment.

It is yet another object of this invention to provide a method for preventing bottlenecks in parallel unloading and loading of database tables.

30 It is still yet another object of the present invention to increase the speed at which a reorganization of database tables is performed by utilizing the parallel processing techniques of the present invention.

- 8 -

These and other objects are accomplished by a method for performing parallel unloading of database tables, including the steps of launching a number of threads to process the database tables, assigning a  
5 respective one of the database tables to a corresponding of said threads and unloading each respective database table by a process of the corresponding thread. The method allows the table  
10 unload process to take advantage of a Symmetric Multi-Processing (SMP) environment to significantly improve the speed of database table unloads.

The invention includes a method of parallel loading of table data including the steps of loading data into database tables, including the steps of  
15 determining X threads for loading data into database tables, creating X temporary tables, each temporary table corresponding to a set of data stored in an export directory, launching an SQL\*LOADER™ process for loading each temporary table, and loading each  
20 temporary table with the data stored in the corresponding export directory via the corresponding SQL\*LOADER™ process.

#### Brief Description of the Drawings

25 A more complete appreciation of the invention and many of the attendant advantages thereof will be readily obtained as the same becomes better understood by reference to the following detailed description when considered in connection with the accompanying  
30 drawings, wherein:

Figure 1 is a flow chart illustrating a database fast unload/load(PDL)procedure;

Figure 2 is a flowchart illustrating the iterative nature of non-concurrent database unload procedure;

- 9 -

Figure 3 is a flowchart illustrating high level processes for concurrent load and unload of database tables according to the present invention;

5 Figure 4 is an illustration of parallel processes performed on a single export disk;

Figure 5 is an illustration of parallel processes performed using multiple export disks according to the present invention;

10 Figure 6 is an illustration of plural storage devices maintaining database tables and export directories according to the present invention;

Figure 7 is a flowchart illustrating export (unload) process according to the present invention;

15 Figure 8 is an illustration of a Destination Tab Property Page according to the present invention;

Figure 9 illustrates a Parallel Query Tab page according to the present invention;

Figure 10 is a flowchart illustrating a parallel load process according to the present invention;

20 Figure 11 illustrates a Concurrent Table Reorganization with a parallel export phase and a parallel direct load according to the present invention;

25 Figure 12 is a flowchart illustrating a fail safe/fault recovery system for the unload/load processes according to the present invention; and

Figure 13 is a flowchart illustrating high level table name and loading processes according to the present invention.

30

- 10 -

Detailed Description of the Invention

## CONCURRENT UNLOAD/LOAD OPTION IN TS REORGANIZATION

5           The present inventor has developed a new Concurrent Unload/Load methodology. This methodology will be utilized in an Unload/Load option to be introduced in Platinum TS reorg V2.1.0 to enhance the speed of tablespace reorganization. Previously, the  
10           fastest way to reorganize a tablespace in TS reorg was to use Fast Unload in conjunction with the Parallel Direct Load option.

          Referring now to the drawings, wherein like reference numerals designate identical or corresponding  
15           parts throughout the several views, and more particularly to Figure 1 thereof, is a flowchart illustrating a Fast Unload/Load process. At step 100, a single Fast Unload process unloads each table, and steps 100a . . . 110n illustrate parallel direct loads  
20           utilized to load tables.

          Figure 2 is flow diagram detailing the Fast Unload process. At step 200, a Fast Unload is performed on a table. The Fast Unload reads and stores data from the table. At step 210, it is determined whether a next  
25           table in the database(s) is to be unloaded and the process is repeated until each table is unloaded.

          As illustrated by the process in Figures 1 and 2, in Parallel Direct Load, the tables are unloaded sequentially using one single process on the unload and  
30           multiple SQL\*LOADER™ (an ORACLE utility) processes on the load, while working on only one table at a time. For example, when unloading a tablespace containing a large number of tables, TS reorg must wait for one

- 11 -

table to be completely unloaded before it can process the next table.

Figure 3 illustrates the Concurrent Unload/Load option of the present invention that enables a user to take advantage of multi-CPU machines to unload (step 300) and load (step 305) multiple tables concurrently. This method is particularly efficient when the tablespace contains a combination of both large and small tables. Then, while TS reorg unloads and loads one large table, it can also unload and load several small tables simultaneously. The principles of the present invention are also extended to the concurrent building of multiple indexes and/or constraints.

The Concurrent Unload/Load process of the present invention are best practiced in a computing and database environment having the following characteristics:

(1) Symmetric Multi Processing (SMP) computing environment - Appropriate computing hardware contains multiple CPUs and must accommodate a symmetrical multiprocessing environment. The higher the symmetrical multiprocessing capacity of a host machine, the higher performance potential.

One type of multiprocessor computer is a *symmetric* multiprocessor (SMP) computer. An SMP computer usually has between 2 and 16 processors, all of which share the computer's single memory source and shared storage devices. The SMP capacity depends on the number of processors available. With more processors, the present invention allows more queries and loads to be performed in parallel.

(2) Datafiles partitioned across multiple disks - The datafiles of the tablespace to unload reorganization must be stored on different physical

- 12 -

disks. This requirement is needed to obtain full benefit from the Concurrent Unload/Load processes of the present invention.

5 (3) Defined export directories on disks - A number of threads used for the Concurrent Unload/Load option processes will be less than or equal to the number of export directories. These directories are best utilized when residing on separate physical disks.

10 Even with a high SMP capacity and efficient parallel operations, data movement can experience a *bottleneck*, or a halting reduction in speed, because of the limited bandwidth of physical storage disks. One of the few limitations of SMP occurs when physical disks cannot quickly accommodate the simultaneous read  
15 and write requests made by the numerous CPUs.

Another source of interrupted data, especially in parallel processing, is *disk contention*. When multiple processes attempt to access and change a database, the disk head can serve the request of only one process at  
20 a time, causing the other operations to wait. This situation of two processes simultaneously trying to access the same disk is called disk contention. The result of disk contention is clogged data, or a bottleneck.

25 Since parallel processes use multiple CPUs to move data between memory and disks, it is important to have multiple disks defined, so that the numerous CPUs can quickly move the data without interruptions or waiting.

30 One way to eliminate bottlenecks is by *partitioning* the data to multiple disks. Partitioning data is the process of physically spreading the data across multiple disk drives to reduce the limiting effects of disk I/O bandwidths and disk contention. The more disks you define (partition) for data export,

- 13 -

the more you reduce I/O bottlenecks, which results in faster parallel operations.

Figure 4 illustrates parallel processes using only one export disk. Even though the server's query coordinator breaks the query into two separate operations, the disk head can serve the request of only one scan at a time, causing the other operations to wait, or bottleneck. In this example, the single disk's I/O limitations defeat the purpose of the server's parallel process features. Resolving this problem depends on proper disk allocation.

Figure 5 illustrates one embodiment of a proper disk allocation, and parallel processes using numerous export disks (partitioned data). The parallel processes run simultaneously, rather than one process waiting for the other.

#### Using Concurrent Threads in TS reorg

The Concurrent Load/Unload processes of the present invention utilize a number of defined threads that determines a number of tables to be unloaded/loaded at the same time. In other words, the number of concurrent threads is equal to the number of processes that TS reorg launches in the unload and load phase. Each of these processes works on one table.

The number of defined threads is limited by a number of export directories. Figure 6 illustrates an SMP computing environment having three separate disks, 310a, 310b, and 310c, each disk storing datafiles and each having an export directory 320a, 320b, and 320c. In the illustrated environment the Concurrent Load/Unload would define three threads (one per export directory) for unloading tables.

- 14 -

#### Concurrent Unload Phase

Figure 7 is a flowchart illustrating the Concurrent Unload phase of the present invention. At step 400, a number of threads to be utilized in a current computing environment is determined. The number of threads to utilize is provided by a user via a GUI or other input device. However, the number of threads is also limited to a maximum equivalent to a number of export directories residing on separate disk drives. For example, TS reorg now utilizes a destination Tab Property Page to enter, change or delete the export directory paths designating the location for the unloaded data and other export files during a tablespace reorg (see Fig. 8).

During the export phase of a tablespace reorg, the reorg utility unloads the Data Description Language (DDL) and table data from the tablespace and exports it to a specified directory. The Host Directory Path(s) for Export is the directory path designating the location for the unloaded DDL and data during a reorg.

The Destination tab property page is utilized to do the following:

- Add a host export directory.
- Delete a host export directory.
- Modify a host export directory.

If insufficient space available in the export directory or the export file size reaches the *ulimit* while the reorg utility is unloading data, the reorg utility exports the remaining data to the next specified export directory. If another export directory is not specified, the TS reorg will not perform the reorganization.

This removes the 2-gigabyte limit common on most UNIX platforms. Since a compressed file cannot be



- 15 -

split, this does not apply to a compressed export method (described hereinafter).

At step 410, TS reorg launches a same number of Fast Unload processes as the number of defined threads.  
5 This is referred to as the first set of Fast Unload processes.

At step 420, a table is assigned to each process. Each Fast Unload process is executed by the corresponding thread and unloads the assigned table  
10 into one of the separate export directories (step 430).

To search for a directory on which to unload data, each thread checks all of the export directories and chooses the smallest available directory that can accommodate the created export file. Such a method is  
15 known as finding the *best fit* for the export directories. Once a thread finds a *best fit* directory, it claims that directory, so the next thread must choose another.

In the event that other applications are using space on the same file system, a thread might not load  
20 all its data into its chosen *best fit* directory. If this occurs, the reorg utility splits the export file into multiple export files and attempts to write it on another export directory. If there is no alternate  
25 export directory, or if the disk is full, the thread terminates and one of the remaining threads unloads the remaining data.

A user may estimate the size of the user's export files in order to assign enough space in an export  
30 directory for those files. To estimate the size of an export file, divide the current allocation of an object by the number of threads to be utilized.

The reorg utility utilizes a Parallel Query tab property page, as shown in Figure 9, to fill in fields

- 16 -

along with any existing parallel parameters for the object. For example, a table's parallel parameters might exist if that table was initially created or revised using parallel parameters. Those parameters  
5 are retained within the table's Data Description Language (DDL) and the Data Dictionary. If parallel parameters already exist for an object, the reorg utility splits the data query automatically, and a user need not specify values in the Parallel Query tab  
10 property page. Otherwise, the reorg utility splits the query according to the parameters set in the Parallel Query tab property page.

For full advantage of the present invention, the number of export directories is equal to or greater  
15 than the number of threads. Also, each export directory should reside on a separate physical disk to avoid disk I/O bottlenecks. As discussed above, I/O bottlenecks will likely occur if more than one export directory is located on a same disk because of the  
20 possibility of more than one process writing (unloading) to the same disk at the same time.

When the first Fast Unload process of the first set of threads finishes, TS reorg checks if the process was successfully terminated (step 440), if that was the  
25 case, the next Fast Unload process to unload a next table is launched (step 470).

In this manner, each table is unloaded by a process executed in it's respective thread. Each thread is independent, therefore if one thread has a  
30 process assigned to a large table, the other processes need not wait for the large process to finish before launching the process for the next table.

Otherwise, if a problem occurred during the unload of the table, TS reorg displays a reason why the

- 17 -

process failed (step 450), along with eventually the Fast Unload logfile in the job logfile. The table that a problem occurred is unloaded using an internal unload program of TS reorg. This unload program is not  
5 launched as a separate process and therefore it has to be finished before TS reorg can continue with the concurrent unload of the rest of the tables.

Once all the tables have been successfully unloaded, the unload process is completed. A  
10 performance meter (estimated, or interactive) may also be implemented to display the speed of unload/load operations and an amount of disk space being consumed.

Once the unload process is completed, TSreorg continues with the next step of the tablespace reorganization which is the export of other objects of  
15 the tablespace (step 485).

Export of other objects includes the creation of a files for storing the DDL of the tablespace including all indexes, triggers, constraints, and primary or  
20 unique keys (everything except the table and clusters).

Finally, TS reorganization performs a drop on the unloaded tables (step 490), and then performs a concurrent load (step 495) loading the unloaded data  
25 into fresh tablespace.

#### Concurrent Load Phase

After dropping all of the tablespace objects, TS reorg begins the load (import) phase of the reorganization. The concurrent load phase is  
30 illustrated in the flowchart of Figure 10. At step 600, TS reorg begins the load phase by creating tables to store the data. TS reorg sequentially creates as many tables as the number of threads specified. The

- 18 -

specified number of threads is equivalent to the number of threads in the unload phase.

After the first set of tables is created, TS reorg simultaneously launches the SQL\*LOADER™ processes used to load the data of each corresponding table. The number of processes launched is equal to the number of threads. Each SQL\*LOADER™ process loads data for one corresponding table, reading table data from an export directory and loading that data into the corresponding table (step 610).

The reorg utility recreates the objects, eliminating fragmentation and optimizing storage parameters, using the DDL in the export files. The reorg utility coordinates multiple CPUs in parallel threads to load the data back into the objects, resulting in a reorganized tablespace.

After one of these processes is completed, TS reorg checks if the data was successfully loaded back into the table and if the number of rows inserted was correct (Step 620). If an error occurred during the load or the number of rows inserted by SQL\*LOADER™ was not correct, TS reorg switches automatically to the internal TS reorg load function to load sequentially the data before continuing with the next table (step 630). If there was no error with the SQL\*LOADER™ process, or after the internal load terminated without error, TS reorg creates the next table (step 650) and launches concurrently still another SQL\*LOADER™ process to load this table's data (repeating step 610).

As illustrated in the concurrent Table Reorganization diagram shown in Figure 11, a concurrent direct load invokes multiple CPUs, each of which executes a separate SQL\*LOADER™ session to synchronously load data from the export files back into

- 19 -

the object's datafiles. When SQL\*LOADER™ writes to these datafiles, the reorg utility uses the best fit method, which selects those partitioning directories that have enough space to accommodate the incoming data. The reorg utility then continues using the best fit method within that selected list of directories and chooses the smallest datafile that can accommodate all of the incoming data for the thread. Each thread loads its data into as much free space in a datafile as it can.

If a fatal error occurs during the load phase, or if the reorganization job is canceled or killed for any reason, the failed job becomes a job that needs recovery. As shown in Figure 12, after a job has failed (step 1105), TS reorg automatically skips all the tables that have been already successfully created and loaded before the failure (step 1110) and only loads concurrently the non-existing tables or the tables that were not completely loaded (step 1115).

As shown in Figure 13, during the load phase, TS reorg first creates the tables under a temporary name (step 1300). After the data has been reloaded the temp table is renamed to the original name of the table (step 1320). Finally, the primary constraints and indexes on the table are created (step 1330). This allows TS reorg to recognize the tables that were not completely imported before the failure and to drop all temporary tables and restart the load for those tables (see process 1115A, Fig. 11).

When using the Concurrent Unload/Load option, TS reorg automatically selects the option: Create all Indexes/Constraints after all of the Tables have been Created. The indexes and constraints of this tablespace are created after all of the tables are

- 20 -

successfully created and their data reloaded into the tables.

5       CONCURRENT INDEX CREATION IN TS REORG

          Concurrent index creation is an option can be used  
in a table or a tablespace reorganization to increase  
the speed of index creation during the reorganization.  
10       In a regular table or tablespace reorganization, the  
indexes are created one after another after the table  
is created and the data loaded back into the table.

          In a tablespace reorg, the user may select a  
concurrent index creation option. In this case, TS  
15       reorg will create the indexes concurrently (sequential  
index creation is the default behavior). This will  
allow TS reorg to gather all the indexes and  
constraints in one export file and when all the tables  
of the tablespace (in case of a tablespace  
20       reorganization) or the table to reorganize (in case of  
the table reorganization) has been created, it starts  
to create these indexes and constraints concurrently.

          When selecting this option, the user has to  
specify a number of threads which would be used as the  
25       number of process to launch concurrently during the  
index/constraint creation phase.

          When the import phase for indexes begins, TS reorg  
launches the same number of Index creation processes as  
the number of specified threads. When the first index  
30       process of the first set of threads finishes, a next  
index creation process is launched and subsequent index  
creation processes are similarly launched upon  
completion of other threads until all the indexes and  
constraints are created.

- 21 -

If an error occurs on the creation of one index or constraint, TS reorg logs the error and writes the DDL of the failed index into a file. The user can then manually edit and fix the problem.

5           The present invention has been described with reference and in terms consistent with an implementation in conjunction with a Oracle database. However, the processes described are equally applicable to other known database products and custom database  
10 installations. For example, instead of utilizing the Oracle utility, SQL\*LOADER™, another program capable of reading table data stored in an export directory and loading the data into fresh table space may be utilized.

15           The present invention may be conveniently implemented using a conventional general purpose or a specialized digital computer or microprocessor programmed according to the teachings of the present disclosure, as will be apparent to those skilled in the  
20 computer art.

          Appropriate software coding can readily be prepared by skilled programmers based on the teachings of the present disclosure, as will be apparent to those skilled in the software art. The invention may also be  
25 implemented by the preparation of application specific integrated circuits or by interconnecting an appropriate network of conventional component circuits, as will be readily apparent to those skilled in the art.

30           The present invention includes a computer program product which is a storage medium (media) having instructions stored thereon/in which can be used to program a computer to perform any of the processes of the present invention. The storage medium can include,

- 22 -

but is not limited to, any type of disk including floppy disks, optical discs, DVD, CD-ROMs, microdrive, and magneto-optical disks, ROMs, RAMs, EPROMs, EEPROMs, DRAMs, VRAMs, flash memory devices, magnetic or optical  
5 cards, nanosystems (including molecular memory ICs), or any type of media or device suitable for storing instructions and/or data.

Stored on any one of the computer readable medium (media), the present invention includes software for  
10 controlling both the hardware of the general purpose/specialized computer or microprocessor, and for enabling the computer or microprocessor to interact with a human user or other mechanism utilizing the results of the present invention. Such software may  
15 include, but is not limited to, device drivers, operating systems, database engines and products, and user applications. Ultimately, such computer readable media further includes software for performing the present invention, as described above.

20 Included in the programming (software) of the general/specialized computer or microprocessor are software modules for implementing the teachings of the present invention, including, but not limited to, retrieval of user inputs and the determination of a  
25 number of threads for parallel processing, launching threads, unloading database tables, determining success, initiating internal reorganization processes, exporting database objects, compressing unloaded data, monitoring the processes of the present invention, and  
30 setting up a concurrent load environment utilizing a table load utility, and the display, storage, or communication of results according to the processes of the present invention.



- 23 -

Obviously, numerous modifications and variations of the present invention are possible in light of the above teachings. It is therefore to be understood that within the scope of the appended claims, the invention  
5 may be practiced otherwise than as specifically described herein.

- 24 -

CLAIMS

What is claimed is:

1. A method for unloading database tables,  
5 comprising the steps of:  
    launching a number of threads to process the database  
    tables;  
    assigning a respective one of said database tables to  
    a corresponding of said threads; and  
10      unloading each respective database table by a process  
    of the corresponding thread.
  
2. The method according to Claim 1, wherein said  
step of launching includes the steps of:  
15      retrieving a number X of threads input by a user;  
    identifying a number of export directories located on  
    separate storage devices that maintain said database  
    tables;  
    limiting X to the number of export directories  
20      identified; and  
    utilizing X as said number of threads.
  
3. The method according to Claim 2, wherein said  
storage devices are disk drives.  
25
  
4. The method according to Claim 1, wherein said  
step of unloading comprises the steps of:  
    reading blocks of data from a respective database  
    table, and storing the data blocks read in an export  
30      directory.
  
5. The method according to Claim 4, wherein said  
step of unloading further comprises the steps of:

- 25 -

determining success of completion said steps of reading and storing; and

(1) stopping each of said threads, if said determining success step indicates non-completion,

5 (2) performing an internal TS Reorganization on the table.

6. The method according to Claim 1, further comprising the step of:

10 repeating said steps of assigning and unloading until each of the database tables is unloaded.

7. The method according to Claim 1, further comprising the steps of:

15 exporting other objects related to each of said tables; and

dropping said tables.

8. A method of loading data into database tables, comprising the steps of:

20 determining X threads for loading data into database tables;

25 creating X temporary tables, each temporary table corresponding to a set of data stored in an export directory;

launching an SQL\* Loader process in conjunction with each thread for loading each temporary table;

30 loading each temporary table with the data stored in the corresponding export directory via the corresponding SQL\* Loader process.

9. The method according to Claim 8, further comprising the steps of:

- 26 -

determining success of the loading step for a respective database table; and

if said step of loading was unsuccessful, performing the steps of:

- 5           (1) stopping each of said threads, if said determining success step indicates non-completion,
- (2) performing an internal TS Reorganization on the table.

10           10. The method according to Claim 8, further comprising the steps of creating, launching and loading until each table is loaded.

15           11. The method according to Claim 8, further comprising the steps of:

            recognizing a fault in said step of loading, and performing the steps of:

- (1) recognizing unsuccessfully loaded tables; and
- 20           (2) performing an internal TS Reorganization on the unsuccessfully loaded tables.

25           12. The method according to Claim 8, wherein said step of determining includes the steps of:

- retrieving a number of threads input by a user;
- determining a number of export directories; and
- establishing a number of threads equal to the lesser of the number of threads retrieved and the number of export directories.

30

13. A computer readable medium having computer instructions stored thereon that, when loaded into a computer, cause the computer to perform the steps of:

- 27 -

launching a number of threads to process the database tables;

assigning a respective one of said database tables to a corresponding of said threads;

5 unloading each respective database table by a process of the corresponding thread.

10 14. The computer readable medium according to Claim 13, wherein said step of launching comprises the steps of:

identifying a number X of export directories located on separate storage devices that maintain said database tables; and

15 utilizing X as said number of threads.

15 15. The computer readable medium according to Claim 14, wherein said storage devices are disk drives.

20 16. The computer readable medium according to Claim 13, wherein said step of unloading comprises the steps of:

reading blocks of data from a respective database table, and storing the data blocks read in an export directory.

25 17. The computer readable medium according to Claim 16, wherein said step of unloading further comprises the steps of:

30 determining success of completion said steps of reading and storing; and

(1) stopping each of said threads, if said determining success step indicates non-completion,

(2) performing an internal TS Reorganization on the table.

- 28 -

18. The method according to Claim 13, further comprising the step of:

repeating said steps of assigning and unloading until each of the database tables is unloaded.

5

19. The method according to Claim 13, further comprising the steps of:

exporting other objects related to each of said tables; and

10 dropping said tables.

20. A computer readable medium having computer instructions stored thereon that, when loaded into a computer, cause the computer to perform the steps of:

15 determining X threads for loading data into database tables;

creating X temporary tables, each temporary table corresponding to a set of data stored in an export directory;

20 launching an SQL\* Loader process for loading each temporary table;

loading each temporary table with the data stored in the corresponding export directory via the corresponding SQL\* Loader process.

25

21. The computer readable medium according to Claim 20, wherein said instructions stored thereon, when loaded into a computer, further cause the computer to perform the steps of:

30 determining success of the loading step for a respective database table; and

if said step of loading was unsuccessful, performing the steps of:

- 29 -

(1) stopping each of said threads, if said determining success step indicates non-completion,  
(2) performing an internal TS Reorganization on the table.

5

22. The computer readable medium according to Claim 20, wherein said instructions stored thereon, when loaded into a computer, further cause the computer to perform the step of repeating said steps of creating,  
10 launching and loading until each table is loaded.

23. The computer readable medium according to Claim 20, wherein said instructions stored thereon, when loaded into a computer, further cause the computer to perform the step of:  
15

recognizing a fault in said step of loading, and performing the steps of:

(1) recognizing unsuccessfully loaded tables;  
and  
20 (2) performing an internal TS Reorganization on the unsuccessfully loaded tables.

24. The computer readable medium according to Claim 20, wherein said step of determining includes the steps of:  
25

retrieving a number of threads input by a user;  
determining a number of export directories; and  
establishing a number of threads equal to the lesser  
of the number of threads retrieved and the number of  
30 export directories.

25. An apparatus for unloading database tables maintained in a system, comprising:

- 30 -

means for launching a number of threads to process the database tables;

means for assigning a respective one of said database tables to a corresponding of said threads; and

5 means for unloading each respective database table by a process of the corresponding thread.

26. The apparatus according to Claim 25, wherein said means for launching further includes means for  
10 retrieving a number of threads to launch from a user.

27. The apparatus according to Claim 25, wherein said means for launching includes means for limiting the number of threads launched to a number of export  
15 directories located on separate storage devices of said system.

28. The apparatus according to Claim 25, wherein said means for unloading comprises:  
20 means for reading blocks of data from a respective database table, and storing the data blocks read in an export directory associated with the corresponding thread.

29. The apparatus according to Claim 25, further  
25 comprising means for repeatedly assigning respective database tables to corresponding threads and unloading each respective table until all of said database tables have been unloaded.

30. An apparatus for unloading database tables,  
30 comprising:  
a launching device configured to initiate a number of threads to process said database tables;



- 31 -

an assignor configured to assign each of said database tables to corresponding of said threads launched; and

5 an unloader installed on each corresponding thread, each respective unloader configured to unload database tables assigned to the corresponding thread of the unloader.

10 31. The apparatus according to Claim 30, wherein said launching device includes:

a user interface configured to retrieve a number of threads to launch from a user; and

15 a limit device configured to limit the number of threads to launch to a number of export directories located on separate storage devices of said system.

32. The apparatus according to Claim 30, wherein each respective unloader comprises:

20 a read device configured to read blocks of data maintained within tables assigned to a same thread as the respective unloader is installed; and

25 a write device configured to store the data blocks read into an export directory corresponding to the same thread as the respective unloader is installed.

30 33. The apparatus according to Claim 30, wherein said assignor assigns a database table to each of said threads, and, upon completion of one of said threads, assigns another of said database tables to the completed thread until each database table has been assigned.

34. The apparatus according to Claim 32, further comprising:

- 32 -

a loader, loaded and executed on each of said threads after each database table has been unloaded, each loader configured to,

5 read data blocks stored in an export directory corresponding to the thread executing the loader, and save the datablocks in fresh tablespace.

FAST UNLOAD/PDL

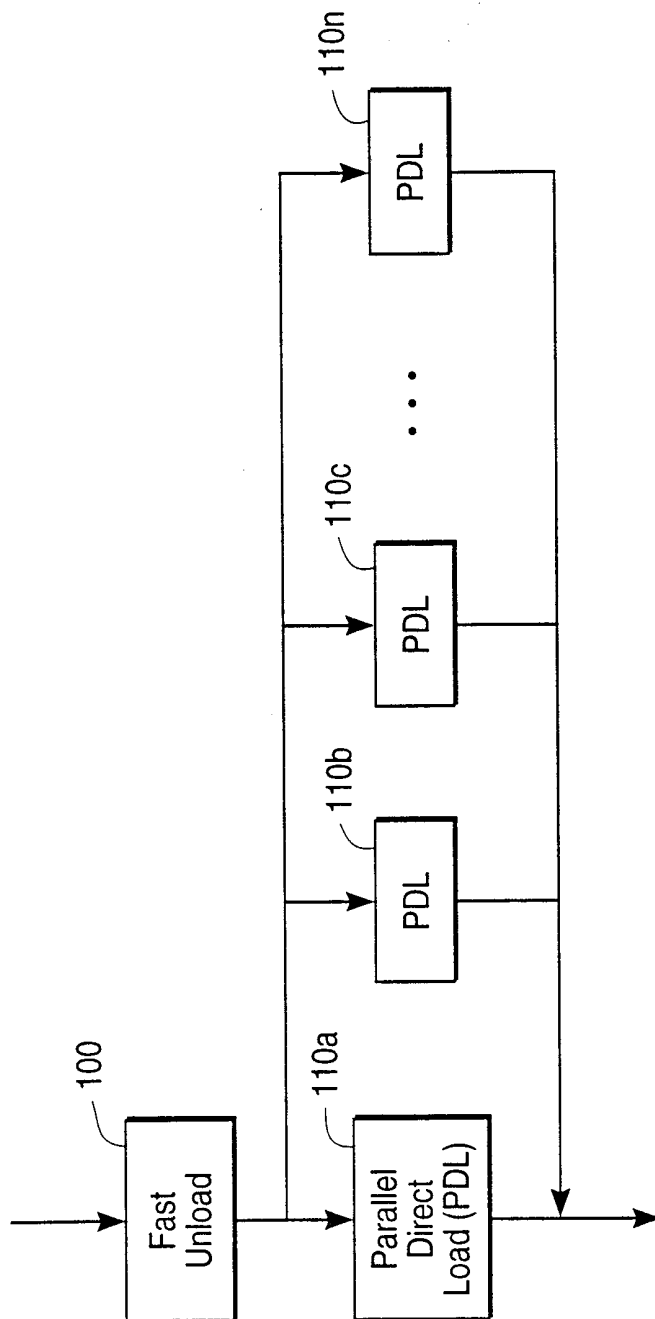


FIG. 1

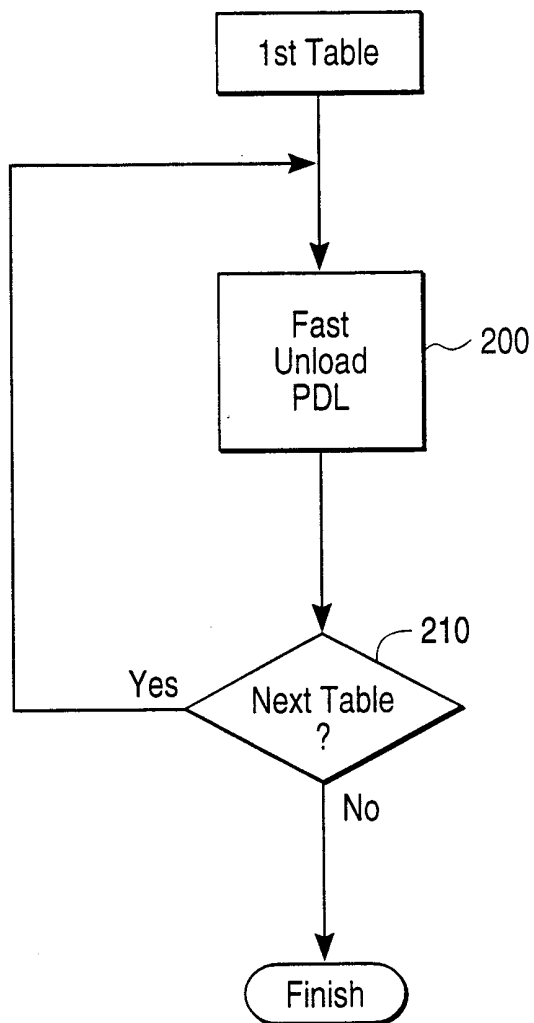


FIG. 2

3/13

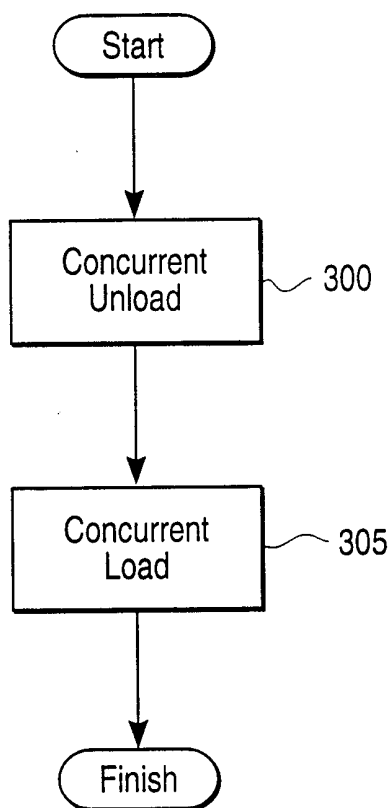


FIG. 3

4/13

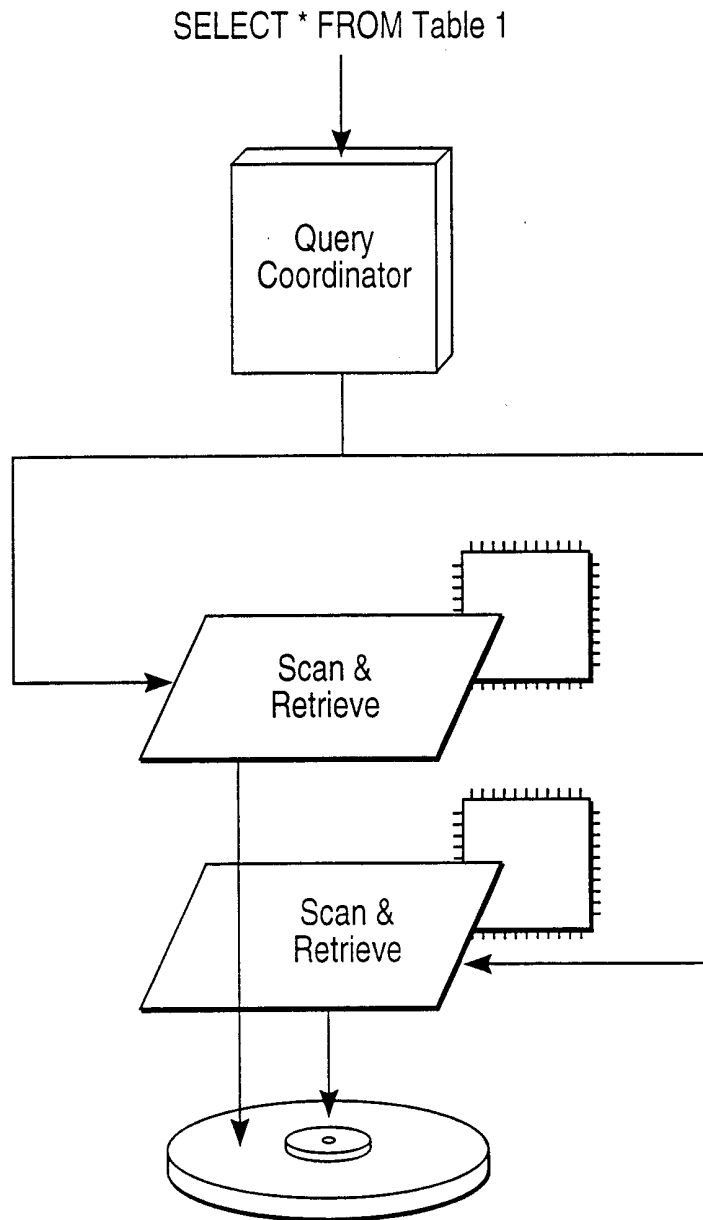


FIG. 4

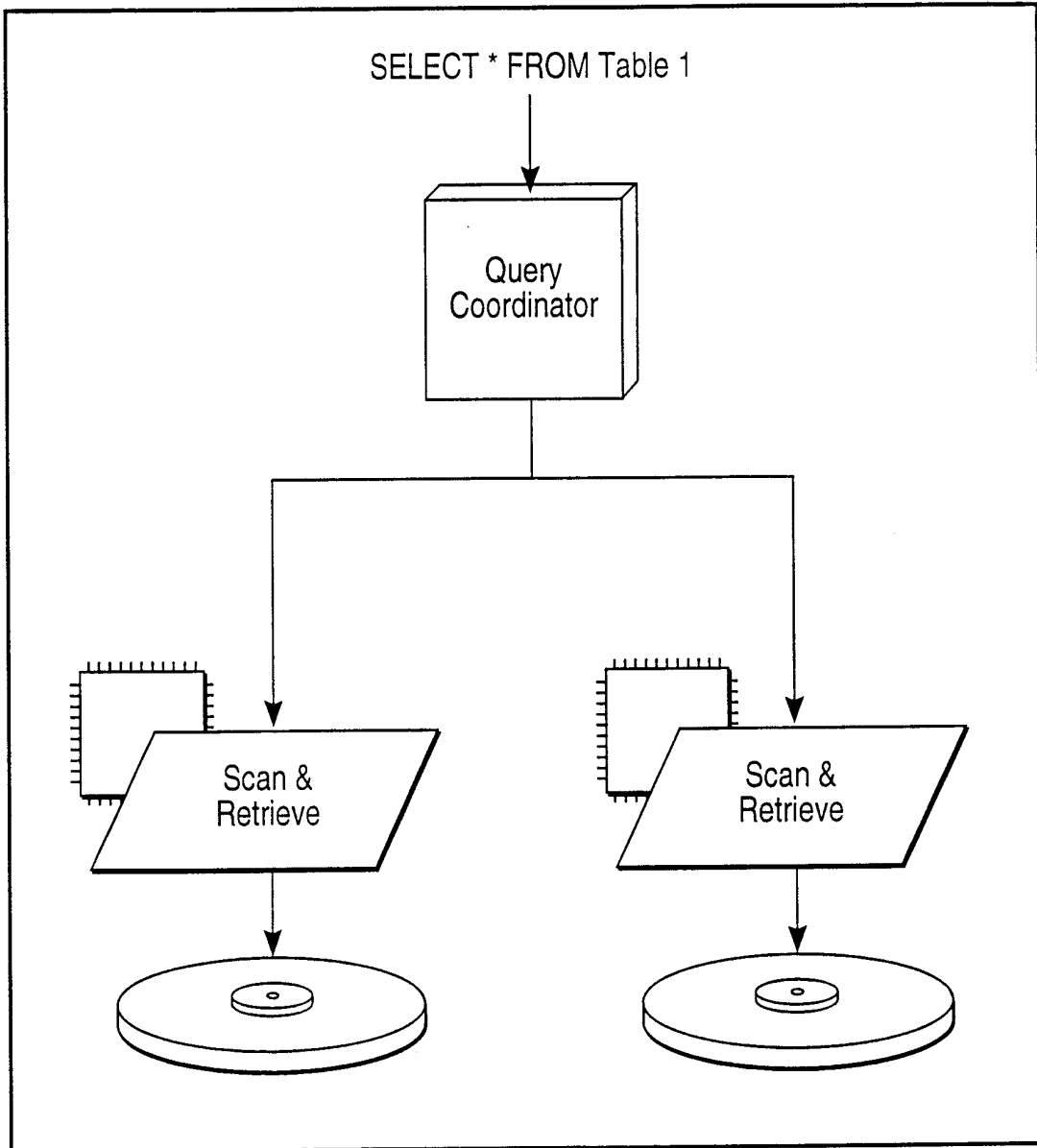


FIG. 5

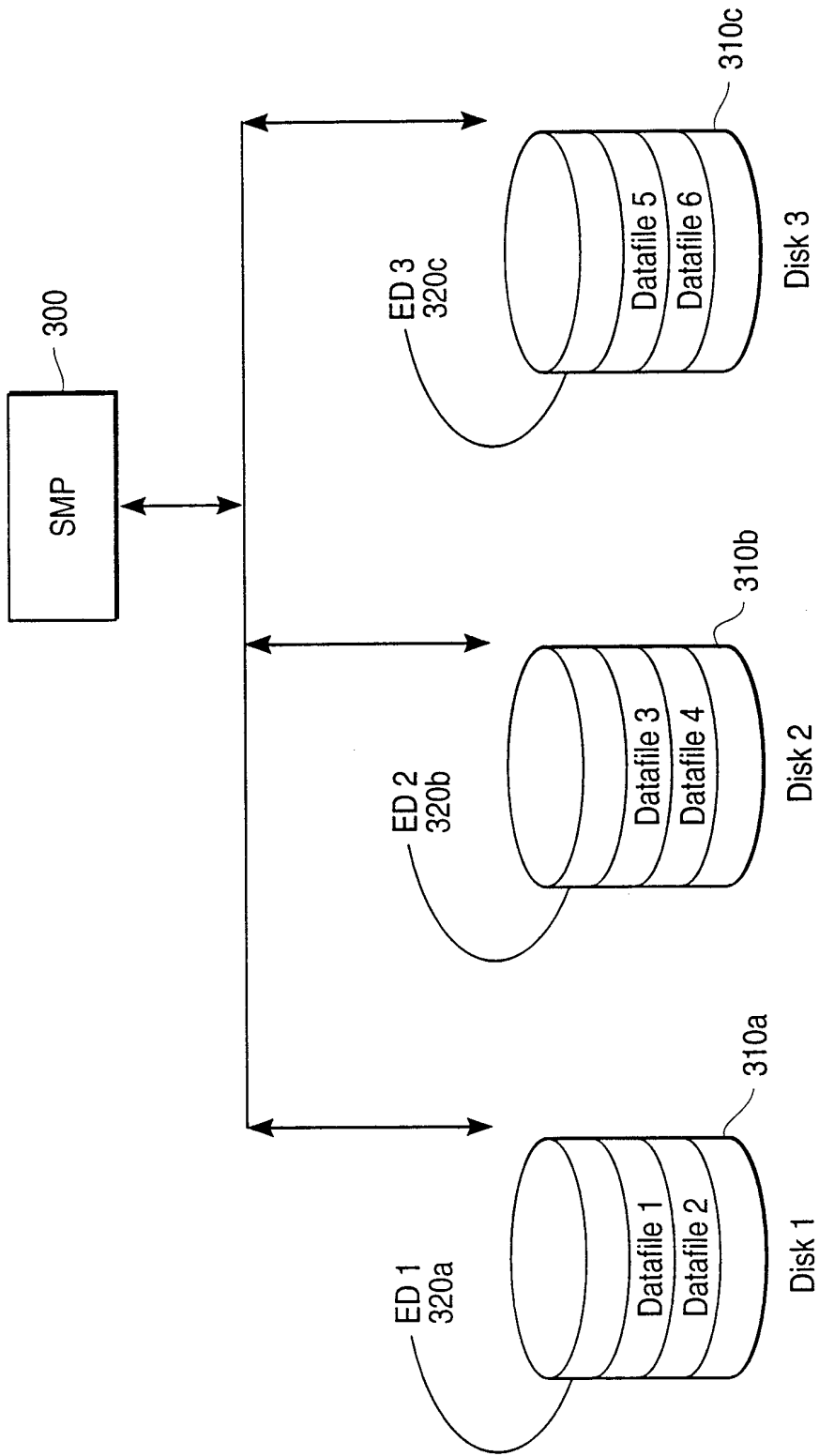


FIG. 6



7/13

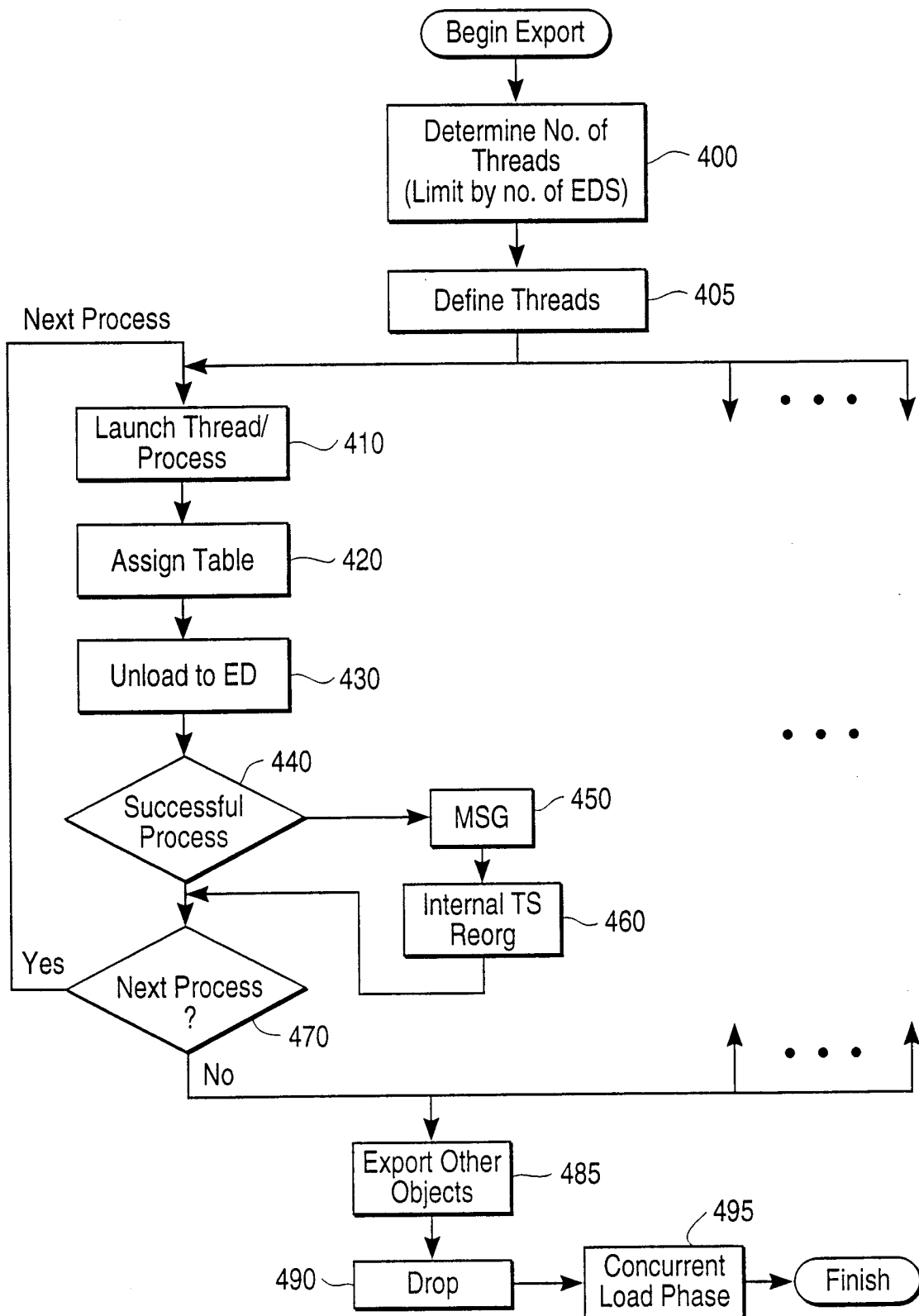


FIG. 7

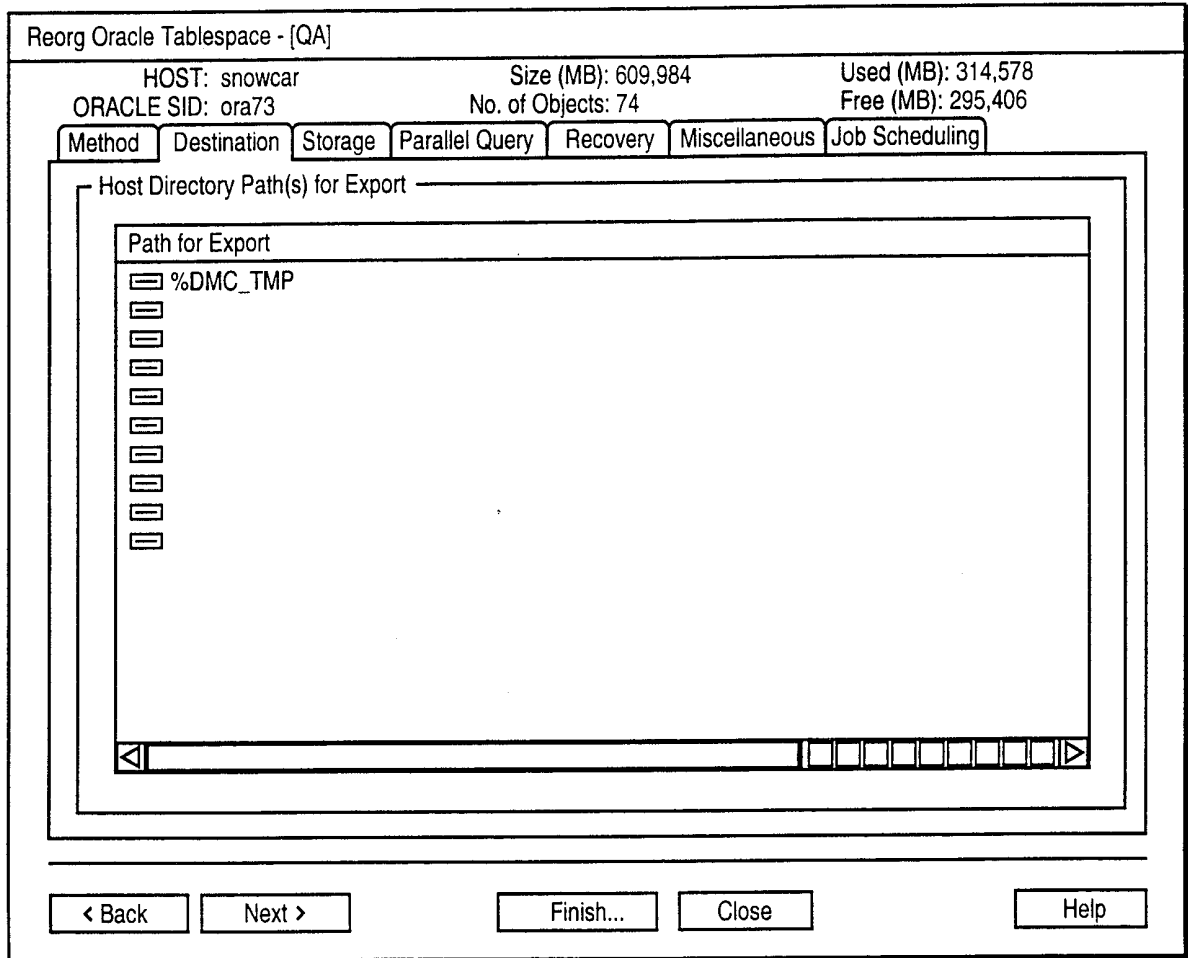


FIG. 8

Reorg Oracle Tablespace - [QA]

HOST: snowcar                      Size (MB): 609,984                      Used (MB): 314,578  
ORACLE SID: ora73                      No. of Objects: 74                      Free (MB): 295,406

Method   Destination   Storage   Parallel Query   Recovery   Miscellaneous   Job Scheduling

Parameter Options

Use Existing Parallel Parameters (if available)

Degree of Parallelism

Default

Other    1

Number of Instances

Default

Other    1

< Back    Next >                      Finish...    Close                      Help

FIG. 9

10/13

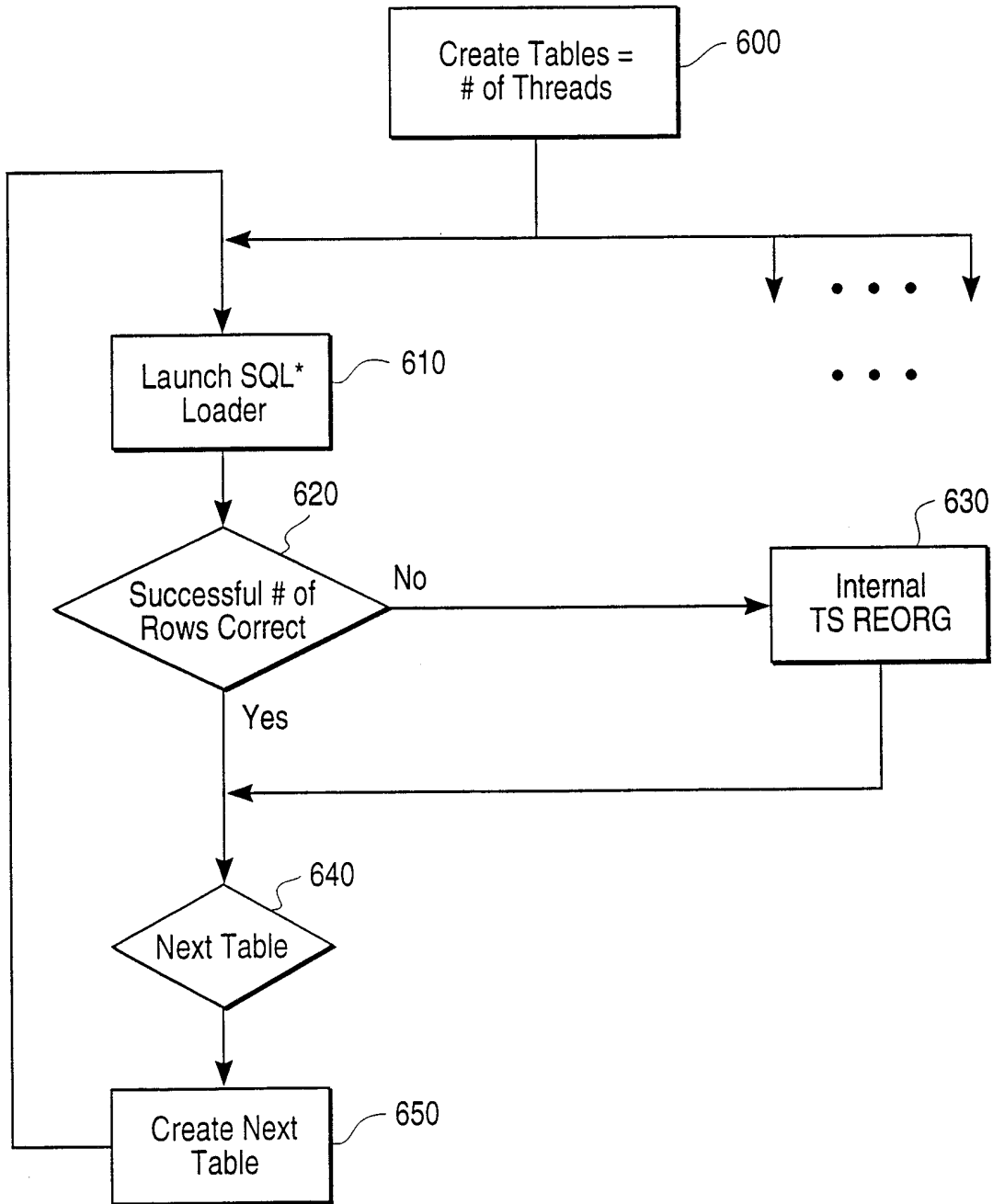


FIG. 10

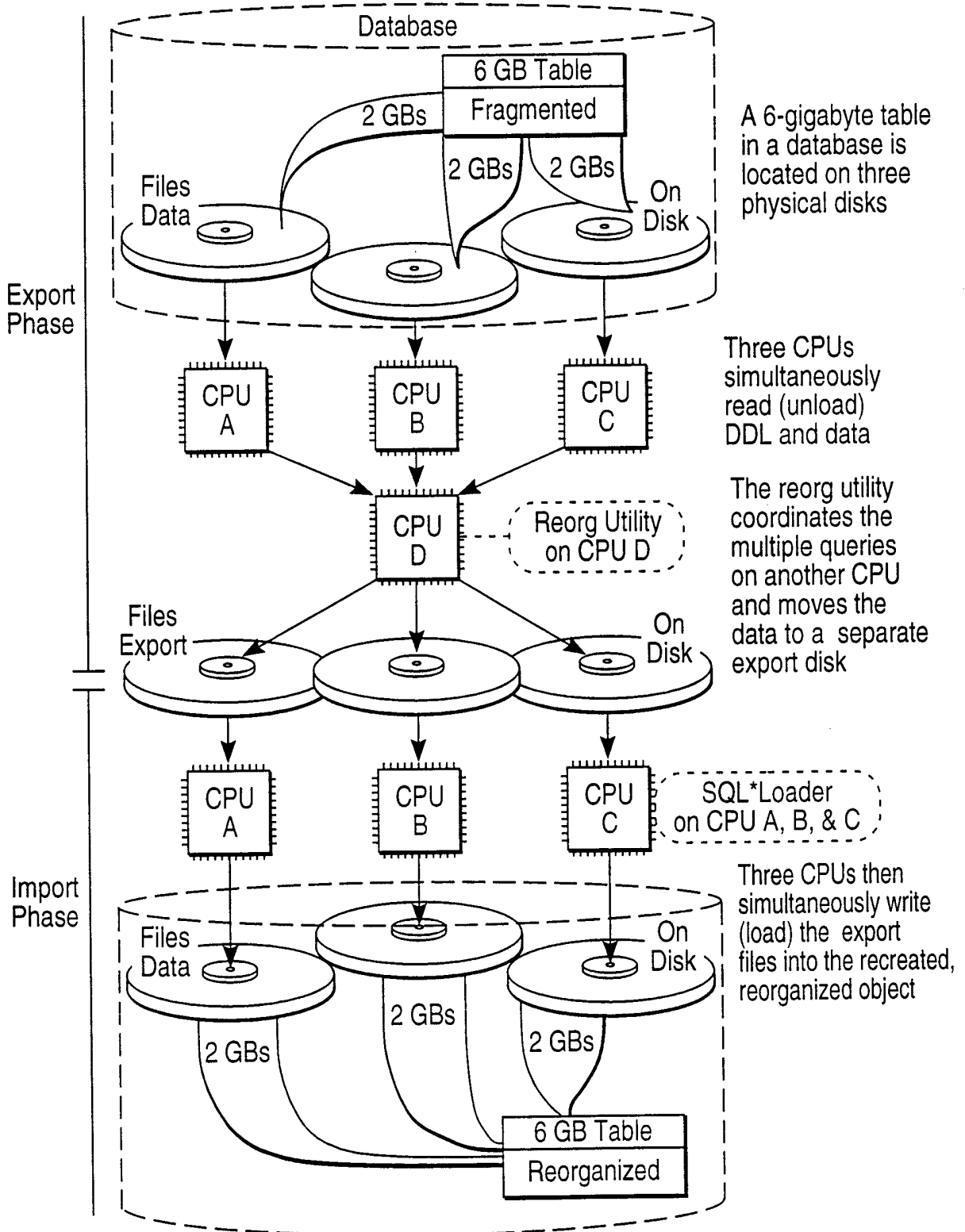


FIG. 11

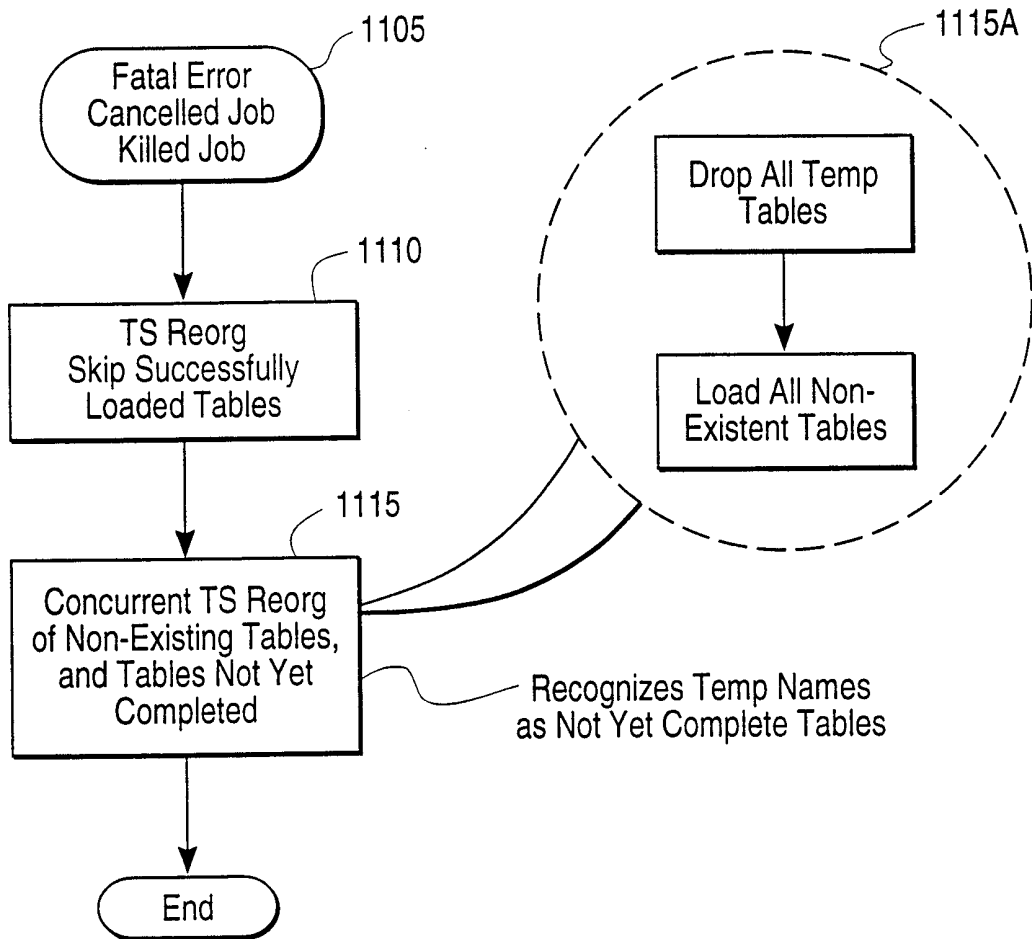


FIG. 12

TS REORG LOAD PROCESS

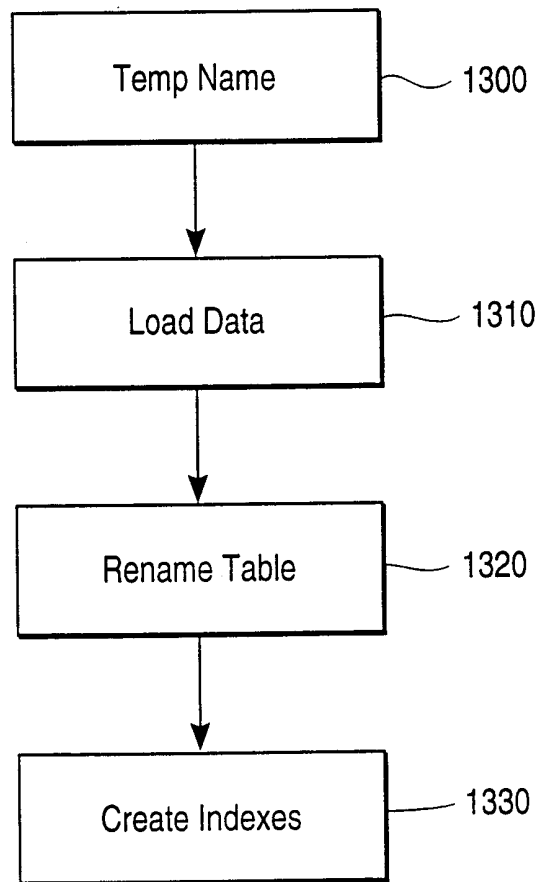


FIG. 13

# INTERNATIONAL SEARCH REPORT

International application No. PCT/US99/27835
---

**A. CLASSIFICATION OF SUBJECT MATTER**  
 IPC(6) : G06F 9/46, 12/00, 17/30  
 US CL : 707/200, 202, 205, 206  
 According to International Patent Classification (IPC) or to both national classification and IPC

**B. FIELDS SEARCHED**  
 Minimum documentation searched (classification system followed by classification symbols)  
 U.S. : 707/200, 202, 205, 206

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)  
 EAST, WEST

**C. DOCUMENTS CONSIDERED TO BE RELEVANT**

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	US 5,754,771 A (EPPERSON ET AL) 19 May 1998, Figure 2.	1
Y	US 5,437,032 A (WOLF ET AL) 25 July 1995, column 1, lines 24-33	1-34
A,P	US 5,884,310 A (BRICHTA ET AL) 16 March 1999, Figure 1.	1-34
A,P	US 5,860,070 A (TOW ET AL) 12 January 1999.	1-34

Further documents are listed in the continuation of Box C.       See patent family annex.

- |   |  |
|---|--|
| <ul style="list-style-type: none"> <li>* Special categories of cited documents:</li> <li>*A* document defining the general state of the art which is not considered to be of particular relevance</li> <li>*B* earlier document published on or after the international filing date</li> <li>*L* document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)</li> <li>*O* document referring to an oral disclosure, use, exhibition or other means</li> <li>*P* document published prior to the international filing date but later than the priority date claimed</li> </ul> | <ul style="list-style-type: none"> <li>*T* later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention</li> <li>*X* document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone</li> <li>*Y* document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art</li> <li>*Z* document member of the same patent family</li> </ul> |
|---|--|

Date of the actual completion of the international search 19 FEBRUARY 2000	Date of mailing of the international search report <b>20 MAR 2000</b>
---	--

Name and mailing address of the ISA/US Commissioner of Patents and Trademarks Box PCT Washington, D.C. 20231 Facsimile No. (703) 305-3230	Authorized officer  UYEN LE Telephone No. (703) 305-4134
---	---