Large Scale Data Warehouses on Grid: Oracle Database 10g and HP ProLiant Servers

Meikel Poess

Oracle Corporation 400 Oracle Parkway Redwood City, CA-94065 USA meikel.poess@oracle.com Raghunath Othayoth Nambiar

Hewlett-Packard Company 20555 Tomball Parkway Houston, TX-77070 USA raghunath.othayoth@hp.com

Abstract

Grid computing has the potential of drastically changing enterprise computing as we know it today. The main concept of Grid computing is to see computing as a utility. It should not matter where data resides, or what computer processes a task. This concept has been applied successfully to academic research. It also has many advantages for commercial data warehouse applications such as virtualization, flexible provisioning, reduced cost due to commodity hardware, high availability and high scale-out. In this paper we show how a large-scale, high performing and scalable Grid based data warehouse can be implemented using commodity hardware (industry standard x86based), Oracle Database 10G and Linux operating system. We further demonstrate this architecture in a recently published TPC-H benchmark.

1. Introduction

Grid Computing has the potential of drastically changing enterprise-computing as we know it today. The main concept of Grid computing is to see computing as a utility. It should not matter where data resides, or what computer processes a task. It should be possible to request information or computation and have it delivered – as much and whenever it is needed. This is analogous to the way electric utilities work, in that one does not know where the generator is, or how the electric grid is wired. One just asks for electricity and gets it. The goal is to make computing a utility - a commodity, and ubiquitous. This, however, is the

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 31st VLDB Conference, Trondheim, Norway, 2005 view of utility computing from a user's point of view. From an implementer point of view Grid is much more; it is about data virtualization, resource provisioning and availability.

Virtualization enables components on all levels, such as storage, processors and database servers to collaborate without creating rigidity and brittleness in the system. Rather than statically determining where a database physically locates its data or which exact server the database runs on, virtualization enables each component of the Grid to react to change more quickly and to adapt to component failures without compromising the entire system. Provisioning ensures that all those that need or request resources are getting what they need. Once resources are virtualized, they can be dynamically allocated for various tasks based on changing priorities. Both, hardware resources and data need to be allocated to databases and application servers. Most importantly, resources are not standing idle while tasks are not being serviced. High availability and scalability ensures that all the data and computation must always be there - just as a utility company must always provide electric power - even when systems are scaled out.

Research institutes have embraced the idea of Grid Computing for quite some time. Taking advantage of idle computing resources around the globe, academic Grids have been established to solve complex computational problems. For instance, SETI@home, a project at the University of Berkeley uses idle PCs on the Internet to mine radio telescope data for signs of extraterrestrial intelligence [1]. It uses a proprietary architecture consisting of a data server, located at U.C. Berkeley, which sends about 350KB of radio frequency data at a time, so called work units to subscribed computers. These computers run a client program as a background process, as a GUI application, or as a screensaver. Results are returned to the data server. In November 2004 IBM, along with representatives of the world's leading science organizations, launched World Community Grid. Similar to SETI@home it consolidates computational power of subscribed computers around the globe. In its first project it directs its computing power to research designed to identify the proteins that make up the Human Proteome and, in doing so, better understand the causes and potential cures for diseases like malaria and tuberculosis [8]. SETI@home and the World Community are Grid architectures that scavenge computing cycles of idle, heterogeneous, completely distributed set of servers. This concept can be applied very successfully to problems that are of common interest to the "world community". The data (radio waves, human protein) is not subject to any security concerns since it is readily available. Most importantly the outcome (discovery of extraterrestrial life or cures for diseases) is of everybody's interest. Lastly, the tasks sent to nodes in these Grids involve highly computational problems on a relatively small data set, 350KB in the SETI@home case, keeping computers easily busy for days with low network requirements. ObjectGlobe [9], similar to other projects like Jini [10], Mariposa [11], Garlic [12] or Amos [13] have studied distributed query processing. They constitute infrastructures that facilitate distributed processing of complex queries executing multiple operators from different sites while also dealing with security concerns. ObjectGlobe takes this concept further into an open environment. TPC-H like Data warehouse applications execute large joins and sort operations making it necessary to tightly couple the nodes of a distributed system such as the data warehouse Grid we are proposing in this paper. The above projects propose a framework for highly distributed systems rather than tightly connected systems.

Applying the above concepts directly to corporate data warehouses is appealing but difficult because data security concerns and network performance problems. Corporate data warehouses contain business intelligence that must not be made known to competitors. Also, the results are generally not in the public interest. Although network bandwidth is increasing, computing table joins of terabyte-sized tables between nodes connected via the Internet is not feasible. Businesses rely on having business intelligence delivered at a determined time. The characteristics of Grid are very appealing for today's corporate data centers. Companies must respond to accelerating business changes fueled by churning market demands, an increasingly dynamic global economy and constant technological innovations. Traditionally, data warehouses have been deployed as a series of islands of systems that cannot easily share resources. This development results in significantly underutilized IT systems and soaring costs. Grid based data warehouses can solve this dilemma.

This paper presents technologies from HP and Oracle that leads the way to large-scale data warehouse Grids that deliver high performance while providing flexible provisioning, high availability and re-source virtualization. As the prices for Linux run commodity hardware (industry standard x86-based) have dropped steadily and as performance and reliability of these systems have improved enormously, the industry is having a serious look at the industry standard server based Grid configurations for large-scale data warehouse applications. In addition Oracle's shared-disk architecture is ideal to the key features of a data warehouse Grids: virtualization, provisioning and availability.

The myth that shared-disk implementations have scaling issues for large-scale data warehouses is addressed by a 12 node published 1000GB scale factor TPC-H benchmark, which delivers performance comparable to equally sized Symmetrical Multi-Processing (SMP) systems. New technologies used in this benchmark, such as InfiniBand, HP's high performance storage arrays and Oracle Grid support, overcome these scalability issues. This benchmark demonstrates a milestone towards a true Grid. Even though this benchmark did not exercise all the aspects of Grid Computing, it addresses key features like scalability, performance and TCO.

The organization of this paper is as follows: In Section 2, the paper defines Oracle and HP's vision for a data warehouse Gird. It further explains how large-scale Grid based data warehouses can be built using industry standard hardware outlining the necessary hardware architectures and hardware features that are supported in HP, such as shared-disk storage, high performance interconnect and high performance computers. Section 3 outlines the features in Oracle Database 10g to implement a scalable shared-disk data warehouse Grid. Finally, Section 4 gives an in depth analysis of a published 1000GB TPC-H benchmark [1] with the emphasize on how the HP and Oracle's hardware and software features were used to achieve high performance.

2. Large-Scale Data Warehouse Grids

The goals of Grid Computing are closely aligned with capabilities and technologies that Oracle and HP have been developing for years. Oracle's latest RDBMS release, Oracle Database 10g provides substantial Grid computing technology. Figure 1 shows how Oracle and HP envision Grid Computing by orchestrating many small servers and storage subsystems into one virtual computer. There are three levels of abstraction: On the first level are server nodes, on the second level are database applications and on the third level are storage subsystems. This three level architecture, which allows for a very flexible Grid implemen-

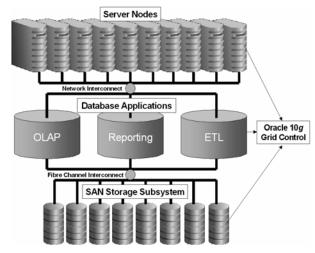


Figure 1: Large Scale Data Warehouse Grid

tation, requires a shared-disk implementation. Shared-disk architectures are characterized by all nodes having access to all data in the database. Hence, there is only one database - logically and physically. Contrary, in a shared-nothing architecture the database is divided into partitions. Each server node has exclusive access to a single partition. Only one server node may own and access a particular partition at a time.

Both architectures share many advantages. They are able to use commodity hardware to reduce Total Cost of Ownership (TCO). They can overcome natural limitations of SMP systems. SMP systems can be scaled-up by adding hardware such as processors, memory or disk arrays. However, they are limited by the resources available on a single server. Today's largest commercially available SMP systems (Sun Fire(TM) E25K Server) are limited to 144 CPUs (72 dual threaded processors) [14]. In contrast Grid systems can be scaled-out virtually without limits by adding additional server nodes.

In addition to the above advantages shared-disk architectures provide dynamic resource sharing between applications of one database and between databases and provisioning by virtualization of resources. Another advantage is the increased availability. Availability is very critical for today's large-scale data warehouses, which provide business-critical services and must therefore operate 24x7. Adding or removing additional systems (scale-out) can be done without interrupting the system as a whole.

2.1. Advantages of Grid Based Data Warehouses

2.1.1 Commodity Hardware

Performance and reliability of Linux run commodity hardware (industry standard x86-based) based systems have improved enormously while prices have dropped steadily. TPC-H results show that commodity hardware based systems can deliver the same performance as large SMP systems at half the price [3,4].

HP ProLiant servers provide features differentiating them from the competition. The number and variety of options and features available for HP ProLiant industry standard servers has grown rapidly and continues to grow today. Development of the ProLiant servers illustrates HP's consistent efforts to provide customers with the world's broadest industry standards based server portfolio and industry-leading innovation in areas such as management, availability, security and virtualization. For instance, the ProLiant DL585 servers used in the TPC-H benchmark are 4U rack-optimized 4-way servers created for large data center deployments requiring enterprise-class performance, uptime, and scalability plus easy management and expansion. They offer customers running 32-bit applications increased performance and memory addressability. While allowing IT organizations to protect their large x86 investments, they also provide a path to more powerful, 64-bit computing to meet evolving business needs for greater processing power and performance.

2.1.2 Scale-Out vs. Scale-Up

Scalability is the ability of a system to maintain desired performance levels as it grows. This is necessary when a system expands to accommodate increased data or a larger number of concurrent users. A system can be either scaled up or scaled out. Traditionally, data warehouse applications have been deployed on scale up (high-end Symmetrical Multi-Processing [SMP] systems) systems. In recent vears the industry has also observed another trend: scaleout configurations. This is fueled by drop in prices for commodity (industry standard x86-based) hardware and their improvement in performance and reliability. Scale-up is achieved by adding devices, typically processors, memory, disks, and network interface cards, to an existing database server. This is also referred to as vertical scaling or vertical growth. Multiple services or applications can be serviced by a single node, which reduces the total administration costs. Server capacity can be easily increased in a server with sufficient expansion capability by adding CPUs, memory and other components. Software licensing costs may be lower since the software is hosted on only one server. On the other hand, all services will be unavailable if the server is down. The availability of the server is limited and depends on server resources. If server load is maximized or the server fails, then the services may be discontinued until the server is replaced with a more capable or operational server. The scalability of the server is limited, and depends on server resources. If it is running at maximum capacity, the server cannot be scaled up. The only alternative would be to replace the existing server. Typically, the initial expense of scale-up server is higher than 4 CPU industry standard server used for scale out. This is due to the larger capability and often more complex architecture of large SMP servers.

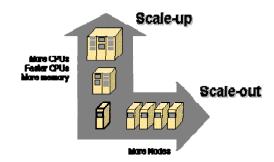


Figure 2: Scale-Up vs. Scale-Out

Scale-out is achieved by adding servers and distributing the existing database application across all of the servers. Scale out architecture has the potential for scalable and high-hosting highly available services, allowing capacity improvements without interruptions. Servers can be maintained and supported much more easily since services are not required to go down to add or remove a server for repair or replacement. There are linear and predictable costs associated with scaling out an application across multiple servers. On the other hand, depending on the initial per-

formance requirements, the cost per server and the cost of infrastructure may be higher than the implementation on one server. Software licenses may be higher when licenses are sold on a per-server basis. Also, management may increase for each server added to the array if appropriate best practices are not defined for the environment. This paper explores the potential of scale-out architecture to host scalable and highly available applications.

2.1.3 Provisioning

Provisioning means allocating resources where they are needed. With Oracle's shared-disk architecture resources can be virtualised. This enables enterprises to dynamically allocate resources for various enterprise tasks based on the changing business priorities reducing underutilized resources and decreasing overcapacities. Data warehouse applications can easily share compute and data resources by migrating between servers of the Grid to leverage available resources. Schedulers on the Grid track resource availability, and assign resources accordingly.

For instance in a business intelligence environment to increase the productivity of data analysts during the day, most resources should be allocated to OLAP queries against data marts, while resources for reporting queries against the data warehouse should be limited. Furthermore, let's assume that the data warehouse is synchronized twice a day with an operational data-store as part of an extraction transformation and load (ETL) task. Hence, during short periods it is imperative that resources are assigned to the ETL process. With a shared-disk based business intelligence Grid resources can be assigned to databases and applications without shutting down databases. In our example above, during the day most nodes of the Grid can be assigned the OLAP applications, while only a few are assigned to reporting queries. Twice a day, when the ETL application runs, nodes can be temporarily detached from the OLAP database and assigned to the ETL application. At night, when no OLAP activity exists most nodes can be assigned to answer reporting queries.

2.1.4 Availability

Another advantage of shared-disk Grid architectures is their increased availability. Availability is very critical for today's large-scale data warehouses, which provide business-critical services and must therefore operate 24x7. For instance, Amazon.com's online recommendation system is fed by an industry standard Grid data warehouse system. Shared-disk Grid increases availability by employing multiple symmetrical systems sharing the same data. Besides increased availability symmetrical systems lead to an increase in computing power. Nodes can fail without compromising the availability of the entire system. They can also be extended or shrunk without bringing down the entire system increasing system availability through system (hardware and software) upgrades.

2.2 Hardware requirements of a shared-disk Grid

2.3.3 Storage

Shared-disk cluster architecture is characterized by all nodes sharing all the data in the database, which necessities a Storage Area Network (SAN). HP SAN provides the data communication infrastructure and management features for most demanding scale-out clusters. In a SAN, server nodes can be added and removed while their data remains in the SAN. Multiple servers can access the same storage for more consistent and rapid processing. The storage itself can be easily increased, changed, or re-assigned. In a SAN, multiple compute servers and backup servers can access a common storage pool. Properly designed SAN storage is highly available, allowing many servers to access a common storage pool with the same degree of availability

A typical SAN is assembled from adapters, SAN switches, and storage enclosures. Fibre Channel host bus adapters are installed into hosts like any other PCI host bus adapters, SAN switches provide scalable storage connectivity of almost any size, storage enclosures place the array controller functionality close to the physical disk drives.

2.3.2 High Speed Inter-Node Communication

Data warehouse queries tend to involve complex, long running operations. In a shared-disk Grid configuration, they are broken down into smaller sub-queries, which are disseminated to participating nodes of the Grid. Upon completion, the results of these sub-queries are forwarded to other nodes for further processing or coalesced into the final result. Performance for these operations is directly correlated to the speed of the inter-node connection.

In a Grid configuration, the network connecting the independent computing nodes is called the Interconnect. In a data warehouse Grid, a low latency and high throughput Interconnect is critical for achieving high performance. During typical data warehouse queries, such as join operations and data load, it is important to effectively pass messages and to transfer data blocks between nodes. It is not uncommon for queries to require more than 100 MB/s throughput per node. The two most commonly available cluster interconnect technologies on industry standard hardware are Gigabit Ethernet (GigE) and InfiniBand (IB).

Cluster technology is now beginning to be adopted by mainstream customers and performance between nodes will largely determine the performance scalability. Although Ethernet technology is ubiquitous in IT organizations, a dedicated IB fabric not only significantly increases overall performance, but more than justifies the additional cost.

3. Grid Support within Oracle Database 10g

In Oracle Database 10g, users at separate nodes can access the same database while simultaneously sharing resources with all other nodes. It provides benefits from the increased transaction processing power and higher availability of multiple nodes. Also, scaling out with multiple nodes enables a cluster to overcome the limitations of a single node, such as memory, CPU, I/O-bandwidth etc. This enables the Grid configuration to supply much greater computing power. This is why Oracle 10g is the ideal platform for Grid based data warehouses.

In a Grid based data warehouse the execution of any particular operation must adapt to the resources available to it at the time it starts (Automatic Resource Allocation). For instance, if a SQL operation is the only operation running in a Grid, it should be given all resources. In contrast, if there are multiple SQL operations running, each should be given the same amount of resources without overloading the Grid. This is important in scale out situations: while the Grid grows new resources should be made available to operations without any user intervention, especially without changing the SQL text.

Furthermore, for a Grid system, it is important to support features like rapid recovery from failures, support for physical and logical standby databases, online maintenance operations, sophisticated diagnosis and repair of failure conditions and Transparent Application Failover. Oracle Database 10g provides these features.

3.1 Dynamic Resource Allocation

Under the term Dynamic Resource Allocation fall many features within the Oracle RDBMS. This paper, however, focuses on those that are applicable to Grid based data warehouse systems: Determining the optimal execution model for parallel query processing, choosing the optimal degree of parallelism, minimizing inter node processing and performing inter process communication efficiently.

Most data warehouse SQL operations are executed in parallel: that is, operations are divided into smaller portions necessary to run in parallel in multiple processes. This is called parallel execution or parallel processing, a feature most RDBMS implementations offer today. The next sections explain how Oracle Database 10g optimally performs parallel processing in Grid data warehouses by utilizing all available resources, dynamically choosing the degree of parallelism and minimizing inter process communication.

3.1.1 Parallel Processing in an Oracle Grid System

One of the most challenging tasks for Grid based data warehouse systems is to perform parallel processing efficiently. In Oracle database 10g, processes executing on a portion of the data are called parallel execution servers. One process, known as the parallel execution coordinator parses each SQL statement to optimize and parallelize its execution. After determining the execution plan of a statement, the parallel execution coordinator decides the parallelization method for each operation in the execution plan. The coordinator must choose whether an operation can be performed in parallel and, if so, how many parallel execution servers on which nodes to enlist. Then, it dispatches the execution of an operation to several parallel execution servers and coordinates the results from all of the server processes.

In detail, the parallel execution coordinator dynamically divides the work into units called granules. One granule at a time is sent to a single parallel execution server. It depends on the operation how granules are generated. For instance, for a table scan, granules are generated as ranges of physical blocks of the table to be scanned. The mapping of granules to execution servers is determined dynamically at execution time. When an execution server finishes with one granule, it gets another granule from the coordinator if there are any granules remaining. This continues until all granules are exhausted, that is, the operation is completed.

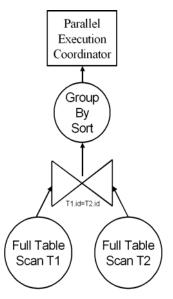
The number of parallel execution servers used for an operation is the degree of parallelism (DOP). The default DOP for any operation is set to twice the number of CPUs available: $DOP = 2 \times \#CPUs$

3.1.2 Inter- and Intra-Operation Parallelism

Data warehouse SQL statements usually perform complex, multiple table accesses, join and sort operations. The Oracle RDBMS applies a producer/consumer model to data operations. This means only two operations in a given execution tree need to be performed simultaneously. Each operation is given its own set of parallel execution servers. Therefore, both operations have parallelism. Parallelization of an individual operation where the same operation is performed on smaller sets of rows by parallel execution servers achieves what is termed intra-operation parallelism. When two operations run concurrently on different sets of parallel execution servers with data flowing from one operation into the other is called inter-operation parallelism. Consider the following simple join between two tables T1 and T2 aggregating on one column.

```
SELECT MAX(T_1.val) maximum FROM T_1, T_2 WHERE T_1.id= T_2.id GROUP BY maximum;
```

Assuming that the DOP for this operation is 64 this query is parallelized in the following fashion (see Figure 3 for the query's execution plan). Each parallel execution server set (S1 and S2) consists of 64 processes. S1 first scans the table T1 while S2 concurrently fetches rows



probes to finish the hash-join in parallel.

operation parallelism).

After S1 has finished scanning the entire table T1 it scans the table T2 in parallel. It sends its rows to parallel execution servers in S2, which then perform the

from S1 building the

hash table for the fol-

lowing hash join op-

(inter-

eration

Figure 3: Parallel Execution
Plan FP1

After S1 is done scanning the table T2 in parallel and sending the rows to S2, it switches to performing the GROUP BY in parallel. This is how two server sets run concurrently to achieve inter-operation parallelism across various operators in the query execution plan while achieving intra-operation parallelism in executing each operation in parallel.

Another important aspect of parallel execution is the repartitioning of rows while they are sent from parallel execution servers in one server set to another. For the query plan in Figure 2, after a server process in S1 scans a row of T1, which server process of S2 should it send it to? The partitioning of rows flowing up the query tree is decided by the operator into which the rows are flowing into. In this case, the partitioning of rows flowing up from S1 performing the parallel scan of T1 into S2 performing the parallel hash-join is done by hash partitioning on the join column value. That is, a server process scanning T1 computes a hash function of the value of the column T1.id to decide the number of the server process in S2 to send it to. Depending on the table partitioning Oracle can optimize this operation using partial or full partition wise joins minimizing inter process communication (see Section 3.4)

In a Grid environment each node hosts a subset of the parallel server processes. In a 16 node 4 CPU per node configuration, each node holds approximately 64 parallel servers. The entire system may hold 1024 parallel servers. The query coordinator prepares the execution depending on resource availability in the Grid. In addition to the usual query parsing, query plan generation and query parallelization steps, query optimization must also determine which DOP to choose and on which node to execute the query.

3.2 Dynamic Degree of Parallelism

In a Grid environment, the degree of parallelism cannot be static but must be adjusted according to resource availability (servers, memory, disk I/O etc.). Resource availability changes in two ways. At any given time in a Grid more or fewer systems can be available to a user for running SQL operations. This could be because resources are shared between multiple applications and/or different users in a Grid or because the system scaled out). Since the DOP is directly related to the number of CPUs, the DOP adjusts automatically as new nodes are added to the system. The DOP for an 8-node 4 CPU per node Grid is 64 (see Section 3.1). If the number of nodes doubles to 16 the DOP is automatically adjusted to 128. Similarly, if nodes are eliminated from the Grid, the DOP decreases proportionally. The second way system resources can be different depending on how many operations are actually running on the system. Usually, there are many users connected concurrently to a data warehouse system issuing queries at any given time. In order to optimize execution time for all users and to utilize all system resources. Oracle database 10g dynamic parallelism allows for adjusting the degree of parallelism before each operation starts (Adaptive Multiuser). When a system is overloaded and the input DOP is larger

than the default DOP, the Adaptive Multiuser Algorithm uses the default degree as input. The system then calculates a reduction factor that it applies to the input DOP. For instance, using a 16 node 4-CPU per node Grid, when the first user enters the system and it is idle, it will be granted a DOP of 128. The next user entering the system will be given a DOP of 64, the next 32, and so on. If the system settles into a steady state of, let's assume 16 concurrent users, all users will be given a DOP of 8, thus dividing the system evenly among all the parallel users.

3.3 Inter Process Communication (IPC)

Inter Process Communication (IPC) refers to sending data and control messages between parallel execution processors of a Grid. IPC is very high during data warehouse operations when join, sort and load operations are performed in parallel. Oracle uses a message-based protocol with its own flow control. Messages between processes on the same node are passed on using shared memory. Messages between processes on different nodes are sent using an operating system dependent IPC protocol. Oracle supports a variety of protocols and wires. With Linux, Oracle supports Ethernet and Infiniband. The protocols that can be run on Ethernet are TCP/IP and UDP/IP. Infiniband supports TPC/IP, UDP/IP and uDAPL. Performance of the different protocol/wire combination differs significantly.

Grid based data warehouses demand a high performance IPC. The amount of interconnect traffic depends on the operation and the number of nodes participating in the operation. Join and sort operations utilize IPC more than simple aggregations because of possible communication between parallel execution servers. The amount of interconnect traffic varies significantly depending on the distribution method. Partial Partition Wise Joins (see next section) in which only one side of the join is redistributed result in less interconnect traffic, while Full Partition Wise Joins (see next section) in which no side needs to be redistributed result in the least interconnect traffic.

The amount of interconnect traffic also depends on how many nodes participate in a join operation. The more nodes that participate in a join operation, the more data needs to be distributed to remote nodes. For instance, in a 4-node Grid cluster with 4 CPU on each node to maximize load performance with external tables the DOP is set to 32 on both the external and internal tables. This will result in 8 parallel server processes performing read operations from the external table on each node as well as 8 parallel server processes performing table creation statements on each node. On the other hand if there are 4 users on average on the systems issuing queries, it is very likely that each user's query runs locally on one node reducing the number of remote data distribution to almost zero.

3.4 Decreasing Inter Process Communication

Oracle database 10g provides functionality to minimize IPC traffic for large join operations, significantly improving performance and scalability for Grid based data ware-

house operations. The most important features are partition-wise joins and node locality. In the default case, parallel execution of a join operation by a set of parallel execution servers requires the redistribution of each table on the join column into disjoint subsets of rows. These disjoint subsets of rows are then joined pair-wise by a single parallel execution server. If at least one of the tables is partitioned on the join key Oracle can avoid redistributing partitions of this table.

3.4.1 Full and Partial Partition-Wise Join

Partition-wise joins minimize the amount of data exchanged between parallel execution servers. In a Grid data

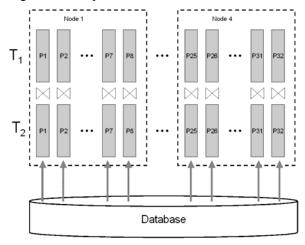


Figure 4: Pull Partition-Wise Join

warehouse this significantly reduces response time by limiting the data traffic over the interconnect (IPC), which is the key to achieving good scalability for massive join operations. Depending on the partitioning scheme of the tables to be joined partition-wise joins can be full or partial.

A full partition-wise join divides a large join into smaller joins between a pair of partitions from the two joined tables. For the optimizer to choose the full partition-wise join method both tables must be equipartitioned on their join keys. That is, they have to be partitioned on the same column with the same partitioning method. Parallel execution of a full partition-wise join is similar to its serial execution. Instead of joining one partition pair at a time, multiple partition pairs are joined in parallel by multiple parallel query servers. The number of partitions joined in parallel is determined by the DOP.

Figure 4 illustrates the parallel execution of a full partition-wise join between two tables T1, and T2 on 4 nodes. Both tables have the same degree of parallelism and the number of partitions, namely 32. As illustrated in the picture, each partition pair is read from the database and joined directly. There is no data redistribution necessary, thus minimizing IPC communication, especially across nodes. Defining more partitions than the degree of parallelism may improve load balancing and limit possible skew in the execution. If there are more partitions than parallel query servers, each time one query server completes the

join of one pair of partitions, it requests another pair to join. This process repeats until all pairs have been processed. This method enables the load to be balanced dynamically when the number of partition pairs is greater than the degree of parallelism. For example, 128 partitions with a degree of parallelism of 32.

Unlike full partition-wise joins, partial partition-wise joins can be applied if only one table is partitioned on the join key. Hence, partial partition-wise joins are more common than full partition-wise joins. To execute a partial partition-wise join, Oracle database 10g dynamically repartitions the other table based on the partitioning of the partitioned table. Once the other table is repartitioned, the exe-

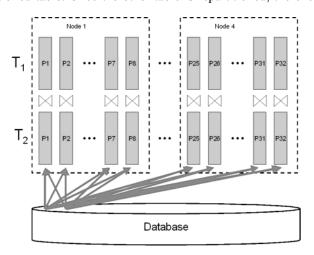


Figure 5: Partial Partition-Wise Join

cution is similar to a full partition-wise join. The redistribution operation involves exchanging rows between parallel execution servers. This operation leads to interconnect traffic in Grid environments since data needs to be repartitioned across node boundaries.

Figure 5 illustrates a partial partition-wise join. It uses the same example as in Figure 2, except that T2 is not partitioned. Before the join operation is executed the rows from T2 are dynamically redistributed on the join key as illustrated by the arrows from the database into P1...P32. For readability not all arrows are drawn. During this redistribution data is sent from one parallel server to another parallel server using IPC.

3.4.2 Node Locality

Another feature in Oracle to reduce IPC traffic is node locality. If the Adaptive Multiuser algorithm determined the DOP to be less or equal to double the number of CPUs per node, queries are run locally on one node. Which node gets to execute the operation is dynamically determined by the load on each system at that time. As resources are becoming available and the DOP is larger than double the number of CPUs per node, operations are executed on multiple nodes. However, Oracle Database 10g always tries to limit the number of nodes executing one operation. This is advantageous because it minimizes the interconnect require-

ments. For instance, if there are 16 users running on a 16-node 4 CPUs per node Grid, the DOP for each operation is 8. Hence, each operation runs on one node. If the number of user drops to 8, each operation is run on two systems.

3.4.3 Dynamic Partition of Grid Resources

Types of operations in a data warehouse range from long running periodic reporting queries over OLAP type queries to data maintenance operations etc. For some data warehouses it is possible to dedicate a specific time period to each of the above operations. In this case, the entire system either calculates reporting queries, ir it is available for on-line users or for data maintenance operations. ever, some data warehouses, especially globally operating systems, cannot afford to dedicate the entire system to specific tasks, but must run them concurrently. In Oracle Database 10g it is possible to dedicate a subset of the Grid to specific tasks. This can be done dynamically without shutting down the system. For instance, during peak hours the system can be available for OLAP users. During off hours, 20% of the system is dedicated to occasional user queries while 40% is dedicated to reporting queries and 40% is dedicated for data maintenance operations.

So far we have assumed that the entire Grid runs one database. It is also possible to run multiple databases on the same Grid. These databases, sharing the same disk subsystem, can take advantage of further features in Oracle Database 10g, such as Transportable Tablespaces. Transportable Tablespaces offers Grid users an extremely fast mechanism to move a subset of data from one Oracle database to another. It allows Oracle data files to be unplugged from a database, moved or copied to another location, and then plugged into another database. Unplugging or plugging a data file involves only reading or loading a small amount of metadata. Transportable Tablespaces also support simultaneous mounting of read-only tablespaces by two or more databases.

If data needs to be shared as it is created or changed, rather than occasionally shared in bulk, Oracle Streams can stream data between databases or nodes in a Grid. It provides a unique way for information sharing, combining message queuing, replication, events, data warehouse loading, notifications and publish/subscribe. Additionally:

- keep two or more copies in sync as updates are applied
- capture database changes,
- propagate database changes to subscribing nodes,
- apply changes,
- detect and resolve any conflicts.
- be used directly by applications as a message queuing feature, enabling communications between applications in the Grid.

4. Industry Standard TPC-H on a Shared-Disk Grid Configuration

TPC-H, the industry's leading decision support benchmark, exemplifies decision support systems that examine large volumes of data, execute queries with a high degree of complexity, and provides data used to answer critical business questions. It consists of a suite of business oriented ad-hoc queries and concurrent data modifications. The queries and the data populating the database have been chosen to have broad industry-wide relevance. Queries are run in two ways to simulate a single user and multi user environment: First queries are submitted by a single stream. In this test each query has all resources available, running massively in parallel (single-user or power test). Then, multiple streams are run concurrently, each running one guery at a time (multi-user or throughput test). Systems that want to excel in a TPC-H benchmark run have to prove that they can handle both the single and multi-user tests.

The TPC-H performance metric is called the TPC-H Composite Query-per-Hour Performance (QphH@Size), and reflects multiple aspects of the system's queries. TPC-H capability to process The Price/Performance metric is expressed as \$/OphH@Size. where \$ is the 3 year cost of ownership of the hardware and software components. In addition, there is a timed portion of database load time reported as secondary metric. TPC-H benchmark specification is available http://www.tpc.org/tpch/.

This section gives a detailed overview of the recently released 1000GB TPC-H benchmark by HP and Oracle. It demonstrates that:

- clustered ProLiant systems with AMD Opteron—x86 processors deliver performance comparable to large SMP systems,
- large-scale data warehouses can be successfully deployed using an industry standard hardware Grid configuration to deliver world record performance,
- the Linux operating system (Red Hat Enterprise Linux AS 3) handles the throughput and processing demands required to achieve the benchmark result,
- Oracle Database 10g delivers consistent, high performance query execution in Grid environments

This result builds on an earlier 3000GB TPC-H benchmark result on 8-node HP ProLiant cluster to demonstrate the commitment of HP and Oracle to this architecture. The benchmark proactively supports current and potential customers that are considering, industry standard hardware running Linux for data warehouses.

4.1 Benchmarked Configuration

The benchmarked configuration was a 12-node ProLiant DL585 cluster connected to a 14.7TB storage fabric SAN (see Figure 6), comprised of 12 HP StorageWorks SAN Switch 2/16s (SAN Switch 2/16) and 48 HP StorageWorks Modular Smart Array 1000 (MSA1000). Each HP ProLiant DL585 server contained six, dual-port HP StorageWorks Fiber Channel 2214DC Host Bus Adapters (HBAs). Each port connects to one of twelve SAN Switch 2/16s. Each SAN Switch 2/16s has four HP StorageWorks MSA1000s connected to it.

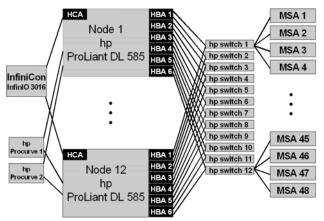


Figure 6: TPC-H Benchmark Configuration

Each server was configured with one InfiniCon Systems InfiniServ 7000 Host Channel Adapter (IB controller) connected to an InfiniCon Systems InfinIO 3016 switch (IB switch), used as cluster interconnect, as shown in Figure 6. Cluster interconnect protocol was UDP/IP over InfiniBand. The InfiniBand interconnect was chosen because it provides higher performance and lower latency than Gigabit Ethernet. The InfiniCon3016 switch is one of several roughly equivalent products to implement this approach. Each server has two on-board GigE NICs, each connected to a HP ProCurve 4148gl switches. One ProCurve 4148gl switch was used for cluster manger communication (cluster heart-beat) and the other was used for user connectivity.

4.2.1 Server

The ProLiant DL585 is an x86, 4-way server, based on AMD Opteron processor technology. Each ProLiant DL585 server was configured with four 2.2GHz/1MB AMD Opteron Model 848 processors (Figure 7). The design of the DL585 server is optimized for AMD Opteron 8000 series chipset and the AMD 800 series of Opteron microprocessors.

The AMD Opteron processor implements the x86 instruction set with a 64-bit memory space. The processor runs 32-bit x86 programs in native mode without application code changes and provides a 64-bit mode for running 64-bit applications. The processor provides program-controlled execution in either 32-bit or 64-bit mode. 32-bit applications can run on top of a 64-bit operating system.

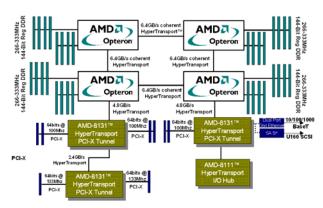


Figure 7: Detailed Server Configuration

The compatibility for 32-bit jobs is in the microcode, which imposes only a small penalty for the conversion to a fixed-length instruction set.

The ProLiant DL585 supports 64GB of memory. In the benchmarked configuration, 8GB of memory per server was selected because Oracle Database's System Global Space (SGA) does not require large amounts of memory in Decision Support Systems. For the same reason a 32-bit version of Oracle Database 10g was used.

In ProLiant DL585, HyperTransport links, a universal high speed chip-to-chip communications technology, is used to interconnect processors and memory. Applications that require high bandwidth and low latency access to system memory; for example, contain large sort operations common in data warehouse queries; benefit significantly from fast memory access. Oracle Database inter-process communications within a node, which is implemented using shared memory, also benefit from fast memory access. The ProLiant DL585 supports cache-coherent non-uniform memory access (ccNUMA). As each processor contains its own memory controller, when a processor accesses memory on its own local memory system, the latency is relatively low, especially when compared to a similar SMP system. If a processor accesses memory located on a different processor, then the latency will be higher. Many operating systems and databases servers can take advantage of ccNUMA.

HyperTransport links are also used to connect processor to I/O sub system. The ProLiant DL585 server contains eight 64-bit PCI-X slots. In the benchmarked configuration, six of them were used for fibre channel storage adapters, one for InfiniBand adapter and one unused. Typical business critical queries require large sequential scans and large multi-table joins,gain significantly from ProLiant DL585's high performance I/O sub-system.

4.2.2 Storage

The HP StorageWorks product set provides strong read performance suitable for DSS applications. It can be configured to support a large number of high-performance yet low-cost fiber-channel connections between the nodes and the arrays.

In the benchmarked configuration, the servers and storage arrays were interconnected using SAN Switch 2/16s. SAN Switch 2/16 is a 16 port fibre channel storage switch that offers 2Gb connectivity for an entry-level SAN, and the ability to grow to a large SAN infrastructure.

Features include redundant power supplies and cooling makes it ideal for supporting corporate infrastructure implementations. Each SAN Switch 2/16s had 12 servers and four MSA1000s connected to it. The benchmarked configuration can be extended to support higher throughput (adding more storage arrays) and larger scale out configurations (adding more servers) by using SAN switches with higher port counts or by interconnecting multiple SAN Switch 2/16s.

The storage array used in the benchmarked configuration was MSA1000, which is a 2Gb fibre channel storage system designed for entry-level to midrange SAN environments, provides low-cost, scalable and high performance storage. With the addition of two more drive enclosures, it can control up to 42 drives, at present allowing for a capacity of 12Terabytes. More spindles result in higher random throughput. In the benchmarked configuration four 36GB 15krpm Ultra320 disk drives per MSA1000 was choosen to achive the targeted performance and price/performance.

One of the key factors that helped in achieving the high I/O throughput was HP Drive Array Technology used in MSA 1000, that distributes data across a series of individual hard drives to unite these physical drives into one or more higher-performance logical arrays and volumes. Distributing the data allows for concurrent access from multiple drives in the array, yielding faster I/O rates with no overhead on the server. HP Drive Array Technology also supports fault-tolerant configurations that protect against data loss due to drive failures. Depending on the performance and application requirements, logical drives in the array can be set to a different level of fault tolerance. The RAID configuration methods supported by the MSA1000 Controller include:

- No fault tolerance/data striping (RAID 0)
- Drive mirroring and striping (RAID 10)
- Distributed data guarding (RAID 5)
- Advanced Data Guarding (RAID ADG)

Further data protection can be achieved by assigning one or more online spares to any fault-tolerant array. The Automatic Data Recovery feature rebuilds data onto a spare or replacement drive when another drive in the array fails in the background.

HP Drive Array technology for MSA 1000s supports various stripe sizes. Stripe size of the array refers to the size of the stripes written to each disk. Striping improves performance by splitting up files into small pieces and distributing them to multiple hard disks.

In the benchmarked configuration, each MSA1000 had two RAID0¹ (data striping) volumes of four 36GB 15krpm Ultra320 disk drives, hosting data files. Eight MSA1000s had two additional 36GB 15krpm Ultra320, RAID10 protected, hosting database redo log files. The RAID0 volumes were created using a stripe size of 256K. Stripe size of

256K was picked to match with the I/O request that Oracle database issues. Oracle database issues I/O requests in of the Oracle database chunks parameters multi block read count (number of blocks in one request) and db block size (size of a database block). To achieve best read performance the array IO/request (stripe size [256K)]* number of disks [4]=1MB) and Oracle database (multi block read count I/Orequest size db block size [16K]=1MB) were set to be the same.

The TPC-H database was partitioned efficiently to reduce overall amount of data required by the queries. The benchmarked configuration had 96 data volumes, two per MSA1000s. All tables of TPC-H schema, except nation and regions, were evenly distributed over these data volumes. Oracle was configured to use a DOP of 96.

MSA1000 Array Accelerator (battery backed cache) can increase performance in database configurations. It performs protected posted-write caching and read-ahead caching, allowing data to be accessed much faster than from disk storage. In protected posted-write caching, data is written to the cache memory on the Array Accelerator rather than directly to the drives. The read-ahead cache detects sequential accesses to the array, reads ahead data, and stores the data in the cache until the next read access arrives. If the data is of a sequential nature, the data can be loaded immediately into memory, avoiding the latency of a disk access. The MSA1000s in the benchmarked configuration was configured with 256MB of Array Accelerator Cache, set to 100% read ahead, because of read intensive nature of the queries.

In addition, certain Linux operating system and Oracle Database features were enabled to improve I/O throughput. Linux operating system and Oracle database engines were configured to support asynchronous I/O. Asynchronous I/O allows a process to submit an I/O request without waiting for it to complete. Because the system does not put the process to sleep while the I/O request is submitted and processed, the calling process is able to perform other tasks while the I/O proceeds. To further enhance I/O throughput, two of the HP StorageWorks 2214DC driver parameters were changed from their default. The gl2xintrdelaytimer the delay in milliseconds posted prior to generating an interrupt was set 0 (default is 3ms) resulting no interrupt migration. The ql2xmaxqdepth - number of outstanding requests per logical unit, which was reduced to 24 from default value of 32.

The cluster obtained actual performance of 7.2GB/Sec or about 600MB/Sec per node. Typical data warehouses often require lower I/O throughput than the "TPC-H" work load.

4.3 Query Scalability

Oracle Database 10g Grid shows very good query scalability. Figure 8 shows the elapsed time for the power and throughput runs of the 1000GB TPC-H benchmark. Bar 0 shows the elapsed time for the power run while bars 1 through 7 show the elapsed time for stream 1 through stream 7 of the throughput run. The elapsed time of the

¹ Production systems will usually use RAID10 or higher to ensure fault tolerance. This was not used in the test to reduce costs and improve performance.

power run is 3163s. The elapsed time of the throughput run varies between 17375s and 19618s. Hence, the ratio between the power run and the various streams varies between 5.5 and 6.2 indicating that the streams of the throughput run scale super-linearly. The reasons for the super-linear behaviour are that during a throughput run the system is utilized more efficiently due to multiple streams can overlap I/O with CPU. Secondly, queries in the throughput run take advantage of node locality. Instead of running across nodes, which might saturate the interconnect, queries run locally on one node, reducing IPC to a minimum. The good scalability between the power and throughput run shows that the query locality feature in Oracle database 10g works and significantly increases performance for multiple users issuing queries simultaneously.

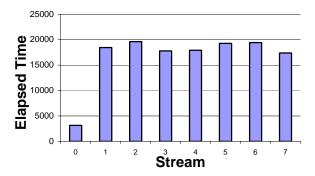


Figure 8: Power vs. Throughput Run

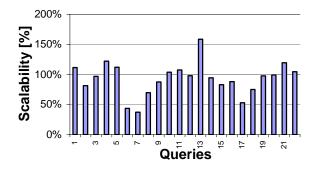


Figure 9: Query Scalability

Figure 9 shows query elapsed time scalability of all 22 TPC-H queries from a 8 to a 12 node Grid. Please note that the elapsed times presented in Figure 9 are not taken from any published TPC result. They were obtained during a lab exercise to obtain empirical data as evidence for Grid scalability. Furthermore, the numbers in the x-axis do not correspond directly to the TPC-H query numbers. The y-axis shows the scalability in percent. 100 percent indicates perfect scalability. Any scalability larger than 100% indicates super scalar behavior. As the chart shows some queries scale poorly while other queries show super scalar behavior. Overall the scalability is at about 90%.

5. Conclusion

In this paper we have shown that the commodity hardware (industry standard x86-based) and software components used in this benchmark are proven cost-effective building blocks for large scale DSS Grids. Enterprise customers can now deploy Oracle Databases 10g on HP ProLiant hardware with Red Hat Linux and obtain high performance at very low cost compared to an equivalent SMP scale-up configuration.

With the 1000GB TPC-H publication we have demonstrated that an InfiniBand-enabled Grid scales to 12 nodes. However, the architecture and database system support scale-out to many more nodes, including failover nodes. The MSA1000 storage fabric and ProLiant cluster size can be extended to support the higher throughput requirements of larger scale-out configurations by using SAN switches with higher port counts or by inter-connecting them. Each MSA1000 has the capacity of accommodating maximum of 42 disk drives and with 300GB disk drives so that the total storage can be expanded by over a factor of 87 from the configured system without adding extra arrays.

With its 48 AMD Opteron CPUs the Grid demonstrated exceptional performance. The MSA1000 based storage SAN enabled throughput of 7.2GB/sec, with additional throughput available as disk drives, nodes and faster SAN switches can be added to the cluster. More spindles will result in higher random throughput. In real-world terms, this means that customers can extend the warehouse and add datamarts within the same, centrally managed, storage environment. The InfiniBand technology compared with the ability to add I/O throughput means that the overall configuration can easily reach higher levels of performance.

Customers demand a reasonable Return on Investment (ROI) from their data warehouses and a low TCO. The industry standard AMD Opteron based ProLiant servers help to reduce overall solution cost. Nodes can be added inexpensively to improve performance or provide failover redundancy. Oracle database 10g can then seamlessly integrate the extra nodes into the cluster. The MSA1000 delivers a cost-effective storage subsystem with a high degree of parallelism, redundancy and performance.

HP and Oracle offer powerful Grid management tools. HP Systems Insight Manager software provides powerful monitoring and control of HP ProLiant Servers and storage. Oracle Enterprise Manager 10g with its Grid Control functionality enables all database instances to be managed in parallel. This reduces costs and ensures consistent management across the cluster.

Typically DSS at these sizes provide business-critical services and require 24x7 availability. The ProLiant servers and MSA storage offer strong availability and reliability capabilities. The MSA1000 supports many fault-tolerant capabilities and features including advanced RAID protection, online spares and automatic data recovery. ProLiant systems and StorageWorks SAN can be configured with redundant adapters, SAN switches and arrays. Along with Oracle Database 10g RAC node failover capability, the

benchmark configuration can be extended to better support mission-critical business intelligence applications.

6. Acknowledgement

The authors would like to thank Ray Glasstone, Bryon Georgson, Jay Workman, Anne Kohlstaedt, George Lumkin and Alex Morgan for their valuable input in writing this paper. Special thanks to Mike Nikolaiev, Tracey Stewart and Ray Glasstone for providing guidance and resources to write this paper.

7. REFERENCES

- [1] David P. Anderson, Jeff Cobb, Eric Korpela, Matt Lebofsky, Dan Werthimer: SETI@home: an experiment in public-resource computing. Commun. ACM 45(11): 56-61 (2002)
- [2] HP ProLiant DL585 Cluster 48P Benchmark Result ID: 104102501: http://www.tpc.org/tpch/results/tpch_result_detail.asp? id=104102501
- [3] PRIMEPOWER 2500 10309080: http://www.tpc.org/tpch/results/tpch_result_detail.asp? id=103090803
- [4] HP ProLiant DL740 Cluster 32P 3000GB TPC-H 1041030202 http://www.tpc.org/tpch/results/tpch_result_detail.asp? id=1041030202
- [5] Ralph Kimball, Kevin Strehlo: Why Decision Support Fails and How To Fix It. SIGMOD Record 24(3): 92-97 (1995)
- [6] Ralph Kimball: The Data Warehouse Toolkit: Practical Techniques for Building Dimensional Data Warehouses. John Wiley 1996
- [7] E. F. Codd: A Relational Model of Data for Large Shared Data Banks. Commun. ACM 13(6): 377-387 (1970)
- [8] World Community Grid: www.worldcommunitygrid.org and 1.ibm.com/grid/grid_press/pr_1116.shtml
- [9] R. Braumandl, M. Keidl, A. Kemper, D. Kossmann, A. Kreutz, S. Seltzsam, K. Stocker: ObjectGlobe, Ubiquitous query processing on the Internet. VLDB J. 10(1): 48-71 (2001)
- [10] J. Waldo. The Jini Architecture for Networkcentric Computing. Communications of the ACM, 42(7):76– 82, 1999.
- [11] M. Stonebraker, P. Aoki, W. Litwin, A. Pfeffer, A. Sah, J. Sidell, C. Staelin, and A. Yu. Mariposa: A wide-area distributed database system. The VLDB Journal, 5(1):48–63, January 1996.
- [12] L. Haas, D. Kossmann, E. Wimmers, and J. Yang. Optimizing queries across diverse data sources. In Proc.

- of the Conf. on Very Large Data Bases (VLDB), pages 276–285, Athens, Greece, August 1997.
- [13] Vanja Josifovski and Tore Risch. Integrating heterogenous overlapping databases through objectoriented transformations. In Proc. of the Conf. on Very Large Data Bases (VLDB), pages 435–446, Edinburgh, GB, September 1999.
- [14] Sun Microsystems, http://www.sun.com/servers/highend/sunfire e25k/index.xml