

Oracle Database 10g Architecture on Windows

An Oracle Technical White Paper
December 2003

Oracle Database 10g Architecture on Windows

Executive Overview.....	3
Introduction.....	3
Oracle Database 10g Architecture on Windows.....	4
Thread Model.....	4
Services.....	5
Scalability Enhancements.....	6
4GB RAM Tuning (4GT) Support.....	6
Large User Populations.....	6
VLM Support.....	6
Large Page Support – New for Oracle Database 10g.....	7
Fiber Support – New for Oracle Database 10g.....	7
Affinity and Priority Settings.....	8
NUMA Support – New for Oracle Database 10g.....	8
64-bit Support.....	9
File I/O Enhancements.....	10
Cluster File System.....	10
64-bit File I/O.....	10
Raw File Support.....	11
Conclusion.....	11

Oracle Database 10g Architecture on Windows

EXECUTIVE OVERVIEW

With the introduction of Oracle Database 10g for Windows, Oracle once again provides the enterprise scalability, reliability, and high performance that customers require. Oracle Database 10g provides enterprise-class data solutions through tight integration with the advanced features of the Windows operating system and the underlying hardware. By using a native, thread-based Windows service model, Oracle Database 10g ensures high performance and scalability. Additional performance improvements have been made with the introduction of Large Page support, a fiber model, and NUMA support in Oracle Database 10g. Oracle can provide enterprise-class performance through the use of large and raw file support, large memory support, and grid computing via clustering. Performance and scalability enhancements are made available with the release of the 64-bit Oracle database on Windows Server 2003. This paper discusses how the Oracle database architecture has been designed to take full advantage of advanced Windows operating system features and the underlying hardware.

INTRODUCTION

With the introduction of Windows Server 2003, Oracle Database 10g has become the leading database for the Windows platform. From the outset, Oracle's goal has always been to provide the highest performing and most tightly integrated database on Windows and, as a result, Oracle invested early to move its market-leading UNIX database technology to the Windows platform. In 1993, Oracle was the first company to provide a relational database for Windows NT.

Initially, Oracle's development efforts were concentrated on improving the performance and optimizing the architecture of the database on Windows. Oracle7 on Windows NT was re-designed to take advantage of several features unique to the Windows platform including native thread support and integration with some of the Windows NT administrative tools such as Performance Monitor and the Event Viewer.

However, Oracle Database 10g on Windows has evolved from the basic level of operating system integration to utilize some of the more advanced services in the Windows platform including the Itanium-based 64-bit version of Windows Server 2003. As always, Oracle is continuing to innovate and leverage new Windows

technologies. This white paper discusses the architecture of Oracle Database 10g on Windows in detail.

ORACLE DATABASE 10g ARCHITECTURE ON WINDOWS

Oracle Database 10g has the same features and functionality on Windows as on UNIX.

However, underneath the covers, significant work has been done to take advantage of Windows-specific operating system features to improve performance, reliability, and stability.

When running on Windows, Oracle Database 10g contains the same features and functionality as it does on the various UNIX platforms that Oracle supports. However, the interface between Oracle Database 10g and the operating system has been substantially modified to take advantage of the unique services provided by Windows. As a result, Oracle Database 10g on Windows is not a straightforward port of the UNIX code base. Significant engineering work has been done to make sure that Oracle Database 10g exploits Windows to the fullest and also to guarantee that Oracle Database 10g is a stable, reliable, and high performing system upon which to build applications.

Thread Model

The architecture of Oracle on Windows is based on threads, rather than processes. Threads provide faster context switches; a much simpler SGA allocation routine that does not require the use of shared memory; faster spawning of new connections; and decreased overall memory usage.

Compared to Oracle Database 10g on UNIX, the most significant architectural change in Oracle Database 10g on Windows is the conversion from a process-based server to a thread-based server. On UNIX, Oracle uses processes to implement background tasks such as database writer (DBW0), log writer (LGWR), dispatchers, shared servers and the like. In addition, each dedicated connection made to the database causes another operating system process to be spawned on behalf of that session. On Windows, however, all of these processes are implemented as threads inside a single, large process. What this means is that for each Oracle database instance, there is only one process running on Windows for the Oracle database server itself. Inside that process will be many running threads with each thread corresponding directly to a process in the UNIX architecture. So, if there were 100 Oracle processes running on UNIX for a particular instance, that same workload would be handled by 100 threads in one process on Windows.

Operationally, client applications connecting to the database are unaffected by this change in database architecture. Every effort has been made to ensure that the database operates in the same way on Windows as it does on other platforms, even though the internal process architecture has been converted to a thread-based approach.

The original motivation to move to a thread-based architecture had to do with performance issues with the first release of Windows NT when dealing with files shared among processes. Simply converting to a thread-based architecture and modifying no other code dramatically increased performance as the particular operating system bottleneck was avoided. No doubt that the original motivation for the change is no longer present; however, the thread architecture for Oracle remains since it has been proven to be a very stable, maintainable one. In addition, there are other benefits that arise out of the thread architecture. These include faster operating system context switches among threads (as opposed to processes); a much simpler SGA allocation routine which does not require the use of shared

memory; faster spawning of new connections since threads are more quickly created than processes; decreased memory usage since threads share more data structures than processes do; and finally, a perception that a thread-based model is somehow more “Windows-like” than a process-based one.

Internally, the code to implement the thread model is compact and very isolated from the main body of Oracle code. Fewer than 20 modules provide the entire infrastructure needed to implement the thread model. In addition, robustness has been added to the architecture through the use of exception handlers and also through routines used to track and de-allocate resources. Both of these additions help allow for 24x7 operation with no downtime due to resource leaks or an ill-behaved program.

Services

In addition to being thread-based, Oracle Database 10g is also not a typical Windows process. It is a Windows *service*, which is basically a background process that’s registered with the operating system, started by Windows at boot time, and which runs under a particular security context. The conversion of Oracle into a service was necessary to allow the database to come up automatically upon system reboot, since services require no user interaction to start. When the Oracle database service starts, there are none of the typical Oracle threads running in the process. Instead, the process basically waits for an initial connection and startup request from SQL*Plus, which will cause a foreground thread to start and which will eventually cause the creation of the background threads and of the SGA. When the database is shutdown, all the threads that were created will terminate, but the process itself will continue to run and will wait for the next connection request and startup command. In addition to the Oracle database service, further support was added which allows the automatic spawning of SQL*Plus to start up and open the database for use by clients. Finally, the Oracle Net Listener is a service since it too needs to be running before users can connect to the database. Again, all of this is basically an implementation detail that does not affect how clients connect to or otherwise use the database, although this is very relevant for administrators of the database on Windows.

The Oracle database runs as a Windows service, which is a background process that can be started by Windows when booting up.

The Oracle database on Windows supports accessing large amounts of memory through a variety of means, including 4GB RAM Tuning, Very Large Memory, and Address Windowing Extensions. Because Oracle can use the maximum possible memory, 64GB, on 32-bit Windows 2000 and Windows Server 2003, users experience better scalability and throughput.

Over the years, Oracle has consistently built its database to serve large user populations. Oracle Real Application Clusters 10g increases capacity for user connections and at the same time increases throughput.

Scalability Enhancements

One of the key goals of the Oracle Database 10g product on Windows is to fully exploit any technologies that can help increase scalability, throughput, and database capacity. The following section describes a few of these technologies, how they affect Oracle, and the benefits that can be derived from them.

4GB RAM Tuning (4GT) Support

Windows 2000 Server (Advanced and Datacenter editions) and Windows Server 2003 (Enterprise and Datacenter editions) include a feature called 4GB RAM Tuning (4GT). This feature allows memory-intensive applications running on Windows to access up to 3GB of memory as opposed to the standard 2GB that is allowed in other editions of Windows. The obvious benefit to Oracle Database 10g is that 50% more memory becomes available for database use, which can increase SGA sizes or connection counts. All Oracle database server releases since version 7.3.4 have supported this feature with no modifications necessary to a standard Oracle installation. The only configuration change required is to ensure that the /3GB flag is used in Windows's boot.ini file.

Large User Populations

One area in which much activity has been undertaken is an effort to support large numbers of connected database users on Windows. As far back as Oracle7 version 7.2, there have been customers in production with over 1000 concurrent connections to a single database instance on Windows NT. As time has progressed, that number has increased to a point where well over 2000 users can connect concurrently to a single database instance in production environments. When using the Oracle shared server architecture, which limits the number of threads running in the Oracle database process, over 10,000 simultaneous connections have been accomplished to a single database instance. In addition, network multiplexing and connection pooling features can also allow a large configuration to achieve more connected users to a single database instance. Finally, Oracle Real Application Clusters can be used to again increase connection counts dramatically by allowing multiple server machines access to the same database files, thereby increasing capacity for user connections and at the same time increasing throughput as well.

VLM Support

One of the key Windows 2000-specific additions originally introduced in Oracle8i was support for Very Large Memory (VLM) configurations. Oracle Database 10g enhances this support and allows the database on Windows to break through the 3GB address space limit normally imposed by 32-bit Windows 2000 and Windows Server 2003. Specifically, a single database instance can now access up to 64GB of database buffers when running on a machine and an O/S that support that much physical memory. In addition, this support in Oracle Database 10g is very tightly integrated with the database buffer cache code inside the database kernel, thereby

allowing very efficient use of the large amounts of RAM available for database buffers. By configuring a database with a large amount of buffers, disk I/O can be diminished since more data is cached in memory. This leads to a corresponding increase in throughput and performance.

Under the covers, Oracle Database 10g on Windows takes advantage of the Address Windowing Extensions (AWE), which are built into all Windows 2000 and Windows Server 2003. The AWE are a set of API calls which allow applications to access more than the traditional 3GB of RAM normally accessible to Windows 2000 and Windows Server 2003 applications. The AWE interface takes advantage of the Intel Xeon architecture and provides a fast map/unmap interface to all memory in a machine.

The AWE calls allow a large increase in database buffer usage up to 64GB of buffers total. This support is purely an in-memory change with no changes or modifications made to the database files themselves.

Large Page Support – New for Oracle Database 10g

Large Page support is a new feature of Oracle Database 10g, which provides a performance boost for memory-intensive database instances on both 32-bit and 64-bit Windows Server 2003. By taking advantage of newly introduced operating system support, Oracle Database 10g can now make more efficient use of processor memory addressing resources. Specifically, when Large Page support is enabled, the CPUs in the system will be able to more quickly access the Oracle database buffers in RAM. Instead of addressing the buffers in either 4KB (on 32-bit) or 8KB (on 64-bit) increments, the CPUs are told to use 4MB or 16MB page sizes when addressing the database buffers. To enable this new feature, the registry variable ORA_LPENABLE should be set to 1 in the Oracle key of the Windows registry. This feature is particularly useful when the Oracle buffer cache is several gigabytes in size. Smaller-sized configurations will still see a gain when using Large Pages, but it will not be as great as when the database is accessing large amounts of memory.

Large Page support, new for Oracle Database 10g, boosts performance for memory-intensive database applications, especially in cases when the buffer cache is several gigabytes in size.

Fiber Support – New for Oracle Database 10g

Another new performance enhancing feature added to Oracle Database 10g on Windows is support for fibers as the basis for the Oracle database, as opposed to threads. Fibers are a Windows concept much like threads, except that fibers are user-scheduled instead of operating system scheduled. What this means is that when the fiber support is turned in Oracle Database 10g, it is not the operating system that determines which Oracle code is run when, but rather it is the database itself that is doing the scheduling of the fibers. By having the database schedule the fibers instead of Windows, they can be more efficiently scheduled based on the current database state. In addition, switching between fibers is a cheaper operation from a CPU perspective than switching between threads, in much the same way that threads are cheaper than processes. Certain configurations are not supported

Fibers, new for Oracle Database 10g, provide faster context switching than threads and are database-scheduled. Thus, they improve overall database performance and throughput.

when running with fibers, but for most user workloads, fibers can be a way to increase performance and throughput when running large-scale applications. Fiber support is turned on and off with a configuration file, while the default configuration remains a thread-based database.

Affinity and Priority Settings

The Oracle Database 10g supports the modification of both priority and affinity settings for the database process and individual threads in that process when running on Windows.

By modifying the value of the ORACLE_PRIORITY registry setting, a database administrator can assign different Windows priorities to the individual background threads and also to the foreground threads as a whole. Likewise, the priority of the entire Oracle process can also be modified. In certain circumstances, this may improve performance slightly for some applications. For instance, if an application generates a great deal of log file activity, the priority of the LGWR thread can be increased to better handle the load put upon it. Likewise, if replication is heavily used, those threads that refresh data to and from remote databases can have their priority bumped up as well.

Much like the ORACLE_PRIORITY setting, the ORACLE_AFFINITY registry setting allows a database administrator to assign the entire Oracle process or individual threads in that process to particular CPUs or groups of CPUs in the system. Again, in certain cases, this can help performance. For instance, pinning DBW0 to a single CPU such that it does not migrate from one CPU to another can in some cases provide a slight performance improvement. Also, if there are other applications running on the system, using ORACLE_AFFINITY can be a way to keep Oracle confined to a subset of the available CPUs in order to give the other applications time to run.

Both ORACLE_PRIORITY and ORACLE_AFFINITY are described in more detail in the Windows-specific documentation that accompanies Oracle Database 10g on Windows.

NUMA Support – New for Oracle Database 10g

With the addition of Non-Uniform Memory Access (NUMA) support in Windows Server 2003, Oracle Database 10g can now better exploit high-end NUMA hardware in which a server is comprised of several computing “nodes”. Since each node in a NUMA machine accesses different parts of physical RAM at different speeds, it is essential that the database can determine the topology of a NUMA machine and adjust its scheduling, memory allocations, and internal operations accordingly. In particular, when running on a NUMA machine on Windows Server 2003, Oracle Database 10g automatically sets the ORACLE_AFFINITY setting to an appropriate default value at startup to maximize resource utilization on the machine. In addition, the memory allocations made by the database when allocating SGA and PGA memory are made in a NUMA-aware fashion such that

Database administrators can assign CPU affinities and priorities to specific Oracle threads to improve their performance.

Oracle Database 10g can automatically detect NUMA hardware and optimize itself by efficiently utilizing NUMA node affinities.

the memory in the machine is accessed as efficiently as possible from all the various nodes in the server. Finally, the number of database writer threads (or fibers) is configured such that there is one per node, again as a performance-enhancing operation.

64-bit Support

The next major step for the Oracle database architecture on Windows has been achieved with the move to 64-bit Itanium, which greatly improves scalability. Because the Oracle database has already been ported to other 64-bit platforms, the move to 64-bit Windows results in a stable, high performing database from Oracle.

The next leap in Oracle database performance and scalability on Windows has been achieved with the 64-bit Oracle database on Intel Itanium-based machines and 64-bit Windows Server 2003. Since being the first to make a developer's release of a database publicly available on 64-bit Windows, Oracle has continued to lead the way toward 64-bit Windows computing by also releasing a production version of the Oracle Database on the same day that 64-bit Windows Server 2003 was launched. The development teams at Oracle have been working closely with the vendors of these technologies to guarantee that the Oracle database on Windows works optimally on the 64-bit hardware and operating system.

As with other Oracle 64-bit ports to different UNIX variants, a 64-bit port of the Oracle database to Windows will be able to handle more connections, allocate much more memory, and provide much better throughput than the 32-bit version of the database on Windows. Oracle's performance and scalability greatly benefit from the larger caches and memory available on Itanium systems. There is no longer a 4GB memory limitation as on 32-bit systems, making 64-bit Oracle perfect for large transaction processing or business intelligence applications. Moreover, Oracle benefits from the improved parallelism, scheduling, and lower number of branch mis-predictions available in Itanium. All these performance enhancements are transparently available in the Oracle database; thus, they require no code changes to be made.

One of the major transparent performance improvements employed by Oracle is profile-guided optimization (PGO). With Intel's Electron compiler, Oracle has designed its database to perform optimally for typical customer workloads. By using the workloads during compilation, a feedback loop is provided to the compiler, which then can analyze the most heavily and lightly used code paths. Based on that information, the compiler can arrange the code paths to be more efficient when run on Itanium. Just by using PGO with no other changes, Oracle has seen approximately a 15%-25% improvement in performance.

The migration path from 32-bit to 64-bit Oracle is very straightforward. There is no need to recreate databases, nor is a full export and import required. All that is needed will be to copy the current datafiles to the new system, install the 64-bit version of Oracle, start the database as normal, and run a few SQL scripts to update the data dictionary.

From an architectural perspective, the current, proven thread-based architecture is used for the 64-bit port. As a result, creating the new 64-bit Oracle software basically entailed re-compiling, re-linking, re-testing and re-releasing the new

version. Very little new code was written during the move to 64-bit since the underlying operating system APIs are substantially the same. In addition, since the Oracle database has already been ported to other 64-bit ports, moving to 64-bit is a straightforward process that will produce a quality, stable product in a very short period of time.

Oracle's close working relationship with Intel has also greatly helped in making the 64-bit Oracle database for Windows a stable, high-performing platform. Specifically, by working closely with Intel compiler development teams, Oracle has been able to significantly increase database throughput through the use of Intel's compiler technology for Itanium. Compared to Oracle9i, Oracle Database 10g is now more fully optimized for Itanium hardware than it had been in the past.

File I/O Enhancements

One other area in which much work has been done in the Oracle database code concerns support for cluster files, large files, and raw files. The Oracle cluster file system is an integral part of Oracle Database 10g that makes administration and installation of Oracle clusters easier on Windows. In an effort to guarantee that all features of Windows are fully exploited by Oracle, the database supports 64-bit file I/O to allow the use of files larger than 4GB in size. In addition, physical and logical raw files are supported as data, log, and control files in order to enable Oracle Real Application Clusters on Windows and also for those cases where performance needs to be maximized.

Cluster File System

With Oracle Database 10g, Real Application Cluster (RAC) manageability has been greatly improved through the introduction of an Oracle cluster file system (CFS). The Oracle CFS was created for use with RAC specifically. Oracle RAC executables are installed on either the CFS or on raw files. In the latter case, at least one database instance runs on each node of the cluster. In a single Oracle home install with CFS, the database will exist on the shared storage, generally a storage array. The Oracle software will be accessible by all nodes in the cluster, but controlled by none. All CFS machines have equal access to all the data and can process any transaction. In this way, RAC with CFS ensures full database software redundancy for Windows clusters while simplifying installation and administration.

64-bit File I/O

Internally, all Oracle Database 10g file I/O routines support 64-bit file offsets, meaning that there are no 2GB or 4GB file size limitations when it comes to data, log, or control files as is the case on some other platforms. In fact, the limitations that are in place are generic Oracle limitations across all ports. These limits include 4 million database blocks per file, 16KB maximum block size, and 64K files per database. If these values are multiplied, the maximum file size for a database file

The Oracle database on Windows supports a cluster file system, easing manageability. 64-bit file I/O support permits file sizes beyond 4GB. Raw files, or unformatted disk partitions, are supported to provide some performance gain over using a file system.

on Windows is calculated to be 64GB while the maximum total database size supported (with 16KB database blocks) is 4 petabytes.

Raw File Support

Like UNIX, Windows supports the concept of raw files, which are basically unformatted disk partitions that can be used as one large file. Raw files have the benefit of no file system overhead, since they are unformatted partitions, and as a result, using raw files for database or log files can produce a slight performance gain. However, the downside to using raw files is manageability since standard Windows commands do not support manipulating or backing up raw files. As a result, raw files are generally used only by very high-end installations and by Oracle Real Application Clusters unless the CFS is used.

To use a raw file, all Oracle needs to be told is the filename specifying which drive letter or partition to use for the file. For instance, the filename `\\.\PhysicalDrive3` tells Oracle to use the 3rd physical drive as a physical raw file as part of the database. In addition, a file such as `\\.\log_file_1` is an example of a raw file that has been assigned an alias for ease of understanding. Aliases can be assigned with the Oracle Object Link Manager (OLM). OLM provides an easy to use GUI interface and maintains the links across the cluster and reboots. When specifying raw filenames to Oracle, care must be taken to choose the right partition number or drive letter, as Oracle will simply overwrite anything on the drive specified when it adds the file to the database, even if it's already an NTFS or FAT formatted drive.

To Oracle, raw files are really no different from other Oracle database files. They are treated in the same way by Oracle and can be backed up and restored via Recovery Manager as any other file can be.

CONCLUSION

In summary, Oracle's database on Windows has evolved from a port of its UNIX database server to a well-integrated native application that takes full advantage of the services and features of the Windows operating system and underlying hardware. Oracle continues to improve the performance, scalability, and capability of its database server on Windows, while at the same time producing a stable, highly functional platform on which to build applications. Oracle is fully committed to providing the highest performing, most well integrated database on the Windows platform on both the 32-bit and the 64-bit versions of Windows. For further information on Oracle's Windows products, please visit:

<http://otn.oracle.com/windows>

<http://www.oracle.com/windows>



Oracle Database 10g Architecture on Windows
December 2003
Author: David Colello
Contributing Authors: Alex Keh

Oracle Corporation
World Headquarters
500 Oracle Parkway
Redwood Shores, CA 94065
U.S.A.

Worldwide Inquiries:
Phone: +1.650.506.7000
Fax: +1.650.506.7200
www.oracle.com

Oracle Corporation provides the software
that powers the internet.

Oracle is a registered trademark of Oracle Corporation. Various
product and service names referenced herein may be trademarks
of Oracle Corporation. All other product and service names
mentioned may be trademarks of their respective owners.

Copyright © 2003 Oracle Corporation
All rights reserved.