

An Oracle White Paper  
September 2010

# Oracle Clusterware 11g Release 2

Introduction .....	1
Introducing the Oracle Grid Infrastructure.....	2
Easier Grid Installation using Oracle Grid Infrastructure .....	4
Typical and Advanced Installation .....	4
Prerequisite Checks, Secure Shell Setup, and FixUp Scripts.....	4
Oracle Clusterware Files stored in Oracle ASM.....	5
Easier Grid Management using Oracle Clusterware .....	6
Grid Naming Service and Automatic-Virtual IP Assignment.....	7
Policy-based and Role-separated Cluster Management .....	7
Clusterized (cluster-aware) Commands.....	9
EM-based Resource and Cluster Management .....	9
Improved Availability – Tuning “Under The Hood” .....	11
Advanced Availability.....	11
Fencing Flexibility and Third Party Cluster Solution Support.....	11
Redundant Interconnect Usage .....	12
Reboot-less node fencing in Oracle Clusterware 11g Release 2 .	13
Cluster Time Synchronization and Network Management.....	14
Cluster Health Monitor: Integrated with Oracle Grid Infrastructure	14
Managing Oracle RAC Databases using Oracle Clusterware .....	15
Managing Any Kind of Application using Oracle Clusterware.....	15
Conclusion .....	16

## Introduction

Oracle Clusterware is portable cluster software that allows clustering of independent servers so that they cooperate as a single system. Oracle Clusterware was first released with Oracle Database 10g Release 1 as the required cluster technology for Oracle Real Application Clusters (RAC). Oracle Clusterware is an independent cluster infrastructure, which is fully integrated with Oracle RAC, capable of protecting any kind of application in a failover cluster.

Oracle Clusterware 11g Release 2 sets a milestone in the development of Oracle's cluster solution. Oracle Clusterware, combined with Oracle Automatic Storage Management (ASM), has become Oracle's Grid Infrastructure software. The integration of these two technologies sets a new level in the management of grid applications.

Oracle Grid Infrastructure introduces a new server pool concept allowing the partitioning of the grid into groups of servers. "Role-separated Management" can be used by organizations, in which cluster, storage, and database management functions are strictly separated. Cluster-aware commands and an Enterprise Manager based cluster and resource management simplify grid management regardless of size. Further enhancements in Oracle ASM, like the new ASM cluster file system or the new dynamic volume manager, complete Oracle's new Grid Infrastructure solution.

The new features of Oracle Clusterware 11g Release 2 discussed in this paper show that Oracle Clusterware does not only continue to be the foundation for Oracle RAC, but has evolved to be the basis for any clustered environment. Oracle Clusterware 11g Release 2 provides a superior level of availability and scalability for any application, eliminating the need for other third party cluster solutions.

## Introducing the Oracle Grid Infrastructure

Oracle Clusterware is a technology transforming a server farm into a cluster. A cluster is defined as a group of independent, but connected servers, cooperating as a single system. Oracle Clusterware is the intelligence in this system that provides the cooperation.

Oracle Clusterware was introduced with Oracle Database 10g Release 1 as the underlying clustering software required for running Oracle Real Application Clusters (RAC). As part of the Oracle RAC stack, Oracle Clusterware is also used by the clustered version of Oracle ASM and is tightly integrated into the Oracle RAC cluster stack.

Oracle Clusterware is a complete clustering solution that can be used outside of Oracle RAC. In these environments, Oracle Clusterware serves as a failover cluster solution, protecting any kind of application.

In both environments, Oracle Clusterware is capable of managing resources, processes, and applications in the cluster as well as for maintaining node membership and ensuring fencing.

With the Oracle Grid Infrastructure, Oracle integrated Oracle ASM, the proven storage management solution for the Oracle Database, and Oracle Clusterware in one software bundle. Oracle has thereby combined two of its strongest products for cluster environments to form a universal grid foundation.

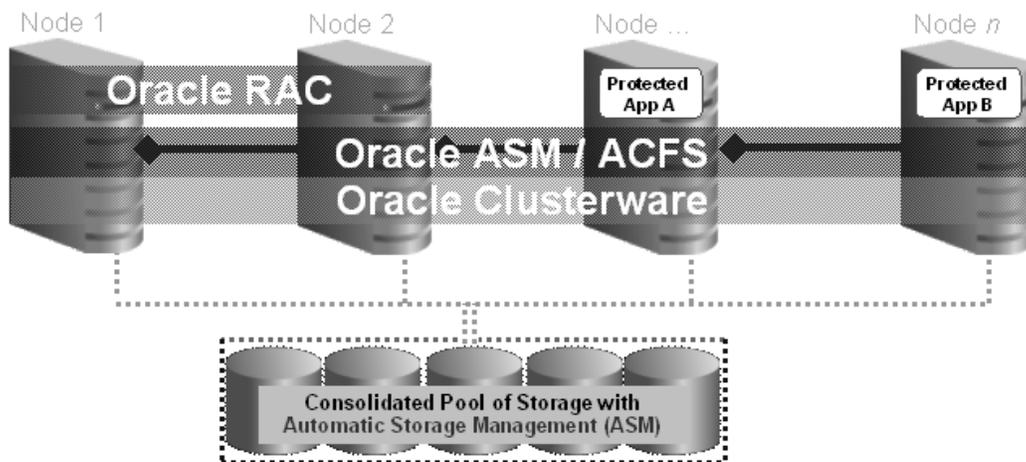


Figure 1 - Oracle Clusterware Configuration Overview

Enhancements to both products, like the Oracle ASM-based cluster file system (ACFS), an Enterprise Manager based graphical user interface for managing the entire cluster, or features like server pools and Grid Plug and Play (GPnP) enable the dynamic management of a grid regardless of size, type, or the number of deployed applications.

The Oracle Grid Infrastructure is not only the foundation for Oracle RAC. It is the basis for any clustered environment, providing a superior level of availability and scalability for any kind of application.

With the addition of a robust storage management solution, it eliminates the need for any third-party management offering and reduces costs in a grid environment as described in figure 2:

1. Capital Expenditures (CAPEX) can be reduced using Oracle’s Grid Infrastructure as a consolidation platform for consolidating all kind of applications in a cluster-based grid.
2. Operation Expenditures (OPEX) are reduced using Oracle Clusterware to lower Data Center costs based on a consolidated environment. In addition, administrative costs are reduced, managing more applications on less servers using only one management interface.

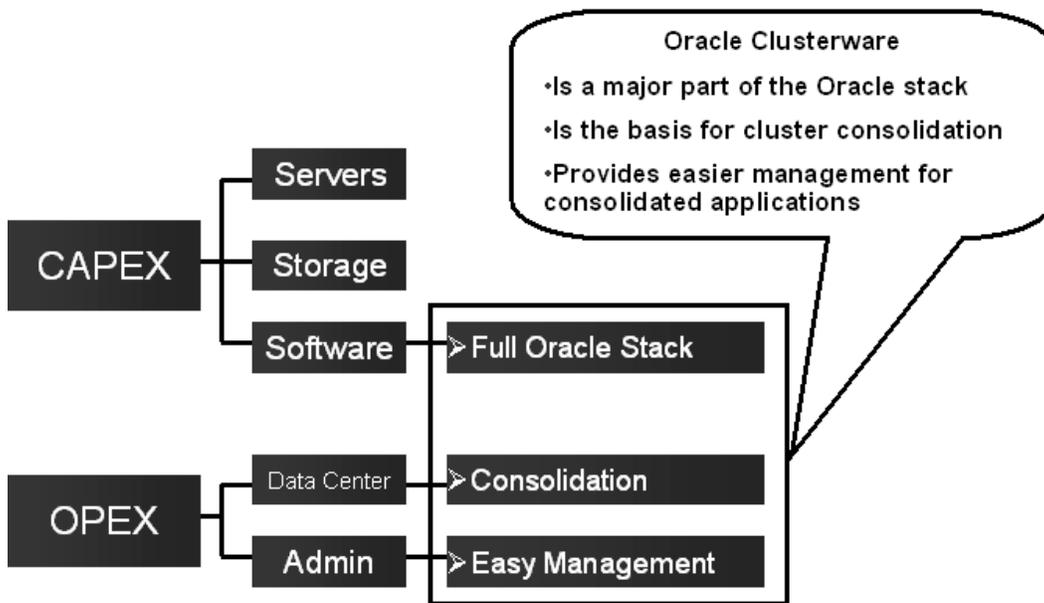


Figure 2: Saving costs with Oracle Clusterware

## Easier Grid Installation using Oracle Grid Infrastructure

The Oracle Grid Infrastructure installation has been simplified and improved in the following ways to provide a general-purpose grid infrastructure:

1. **Validation** – any input made during the interactive installation using Oracle Universal Installer (OUI) is validated immediately to prevent failures later.
2. **Automation** – the new installer has been optimized to complete certain steps automatically and to anticipate a typical configuration.
3. For **mass deployments**, Oracle Grid Infrastructure provides a new “Software Only Installation”, which separates the software deployment and the configuration. Once the software has been installed, the configuration can be performed at a later time.
4. For **single instance environments**, Oracle Grid Infrastructure can be installed in a “standalone server” configuration, providing Oracle ASM and a functionality to restart Single Instance Databases, called Oracle Restart.

### Typical and Advanced Installation

Oracle Grid Infrastructure can be installed in either a “Typical Installation” or an “Advanced Installation”. The typical installation provides a fast way to install Oracle Grid Infrastructure with recommended defaults to simplify the installation, providing fewer options to customize the environment. If more flexibility is required, the Advanced Installation should be used, providing more customization options for more complex architectures.

### Prerequisite Checks, Secure Shell Setup, and FixUp Scripts

Building a grid based on independent servers requires certain prerequisites to be met on each server before cluster software can successfully operate on those machines, ensuring a secure communication within the cluster and therefore stable cluster operation. The Oracle Universal Installer (OUI) for Oracle Database 11g Release 2 now performs all prerequisite checks for an Oracle Grid Infrastructure installation using the Cluster Verification Utility (CVU).

If one or more of the prerequisite checks fail, because the server configuration does not meet certain requirements to install the Oracle Grid Infrastructure, the OUI will list the requirements that were not met and provide guidance for their rectification.

The ability of the OUI to check remote nodes relies on a proper setup of Secure Shell (SSH). The OUI will check whether the nodes in the cluster can be reached via SSH and it will not allow proceeding with the installation until SSH has been properly set up.

Furthermore, the new Oracle Universal Installer offers the option of an automatic SSH configuration across the cluster. If SSH was already (partially) set up on the machines, OUI offers to check the setup only, check and correct, if required, or set up SSH from scratch.

Once SSH has been set up and all checks have been performed, OUI lists the details of the checks that failed. For some checks – depending on the nature of the check – the new OUI is now able to provide a fix by means of a FixUp script as shown in figure 4.

The FixUp script must then be run as root on all nodes that require fixing. Alternatively, the configuration can be changed manually. And even though, it would not be recommended, failed checks can be ignored and fixed at a later time.

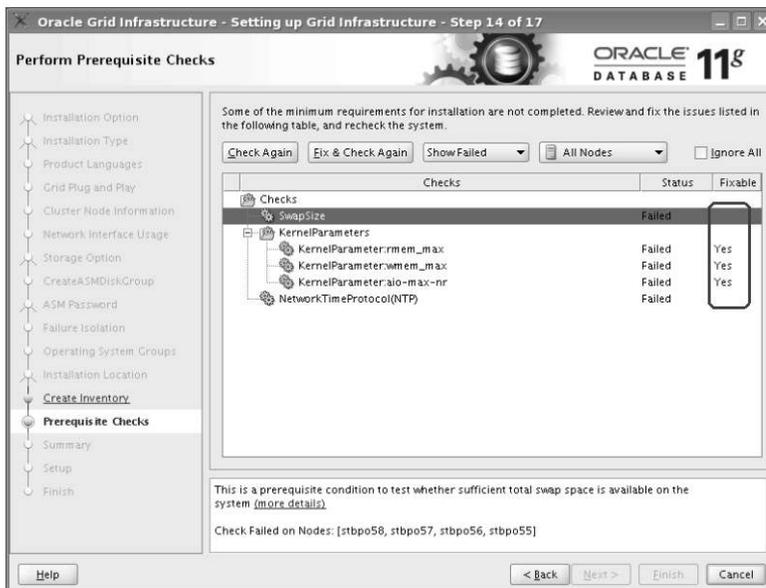


Figure 3: OUI Prerequisite Check Result Page and FixUp Scripts

## Oracle Clusterware Files stored in Oracle ASM

Oracle Automatic Storage Management (ASM) is now integrated with Oracle Clusterware in the Oracle Grid Infrastructure. Oracle ASM with Oracle Database 11g Release 2 provides a complete storage management solution. Part of this solution is the ability to store the Oracle Clusterware files, namely the Oracle Cluster Registry (OCR) and the Voting Files (VF – also called Voting Disks), eliminating the need for a third party storage solution.

Both file types are stored in Oracle ASM disk groups as other database related files are and therefore utilize the ASM disk group configuration with respect to redundancy. This means that a normal redundancy disk group will hold a 2-way-mirrored OCR. A failure of one disk in the disk group will not prevent access to the OCR. In case of high-redundancy disk group (3-way-mirrored), two independent disks can fail without impacting the access to the OCR. For external redundancy, no protection is provided by Oracle.

Only one OCR per disk group can be used in order to protect from physical disk failures. Oracle recommends using two independent disk groups for the OCR to ensure uninterrupted cluster operation in case the OCR is subject to corruptions rather than disk failures.

The Voting Files are managed in a similar way. They follow the Oracle ASM disk group with respect to redundancy, but are not managed as normal ASM files in the disk group. Instead, each voting disk is placed on a specific disk in the disk group. The disk and the location of the Voting Files on the disks are stored internally within Oracle Clusterware.

## Easier Grid Management using Oracle Clusterware

The management of Oracle Clusterware and the Oracle Grid Infrastructure has been improved to make it easy to manage any kind of application in a grid regardless of size. For the management of Oracle RAC databases, specialized management tools have been incorporated to make the management of Oracle RAC databases in the grid even easier.

The management improvements simplify the management in at least one of these areas:

1. Manage a grid regardless of size – up to hundreds of servers in the cluster
2. Make Oracle Clusterware independent of third party solutions
3. Easy maintenance – reduce the number of (manual) administrative tasks
4. Easy management of third party applications (using one management tool only)

In order to avoid unnecessary upfront configuration steps for third party applications and to avoid re-configuration steps later, most of the new management features are ideally configured during the initial installation.

Managing Oracle Clusterware files in Oracle ASM is one example for such a feature. While this feature eases the installation by providing an easy way to store the Oracle Clusterware files in a shared storage infrastructure, it provides additional benefits due to the fact that those files now virtually do not require any continuous management.

Other examples include the Grid Naming Service together with the Automatic-Virtual IP Assignment and the Oracle Cluster Time Synchronization Service, which are both typically activated during installation, making it easier to install Oracle Grid Infrastructure. These features show their real benefits on a day-to-day basis when managing the grid.

## Grid Naming Service and Automatic-Virtual IP Assignment

Dynamic management of large grid environments having the ability to add or remove nodes on demand requires a dynamic naming scheme. Static naming, like the name resolution of a Virtual Internet Protocol (VIP) address to a server name in the Domain Name Service (DNS), is problematic in dynamic environments, since any change with respect to the number of nodes in the cluster would require respective updates.

The new Grid Naming Service (GNS) solves this problem. Based on Dynamic Host Configuration Protocol (DHCP) assigned VIP addresses (the DHCP server assumed here is not part of the Oracle Grid Infrastructure software) names can now be dynamically resolved via GNS. Having GNS and Auto-VIP assignment in place, means that whenever a server is added to or removed from the cluster, no VIP and name assignment needs to be performed. The cluster manages itself in these regards.

The use of GNS and the new, dynamic name resolution is optional and is typically enabled during installation (or as a post-installation maintenance task). Before activating the GNS, a domain must be delegated in the corporate DNS to the GNS and the GNS IP, which are the only static entries in the DNS. When clients request access to a node in the delegated domain, they will then transparently find their ways into the cluster.

## Policy-based and Role-separated Cluster Management

Using the Oracle Grid Naming Service (GNS) solves the problem of using fixed assigned names in an otherwise dynamic environment. However, even when using GNS, applications still need to be dynamically assigned to run on certain servers or, more generally, hardware resources in the cluster, when highly dynamic environments are used. The new server pool concept in Oracle Clusterware solves this problem.

Server Pools partition the cluster. They create logical divisions in the cluster, combining individual servers into groups of servers, typically running similar workload. Server pools therefore enable dynamic capacity assignment when required, but also ensure isolation where necessary (one may think about server pools as “dedicated groups of servers in the cluster”).

Server pools do not determine the placement of applications or resources in the cluster. Policies need to be used to run the application on the right number of servers. What is required is that an application can be run on a minimum amount of servers in order to meet its workload requirements. Server pools are therefore defined using three attributes besides their name:

1. Min – specifies the “minimum” number of servers that should run in the server pool
2. Max – states the “maximum” number of servers that can run in the server pool.
3. Imp – “importance” specifies the relative importance between server pools. This parameter is of relevance at the time of the server assignment to server pools or when servers need to be re-shuffled in the cluster due to failures.

Because applications are now defined to run in server pools rather than on named servers, using these attributes on all server pools defines a policy for the assignment of applications to hardware resources in the cluster.

In this context, it should be noted that managing the cluster “in the old way” – having named servers and assigning applications and resources to distinct, individual servers – is still possible. This is to support backward compatibility.

Furthermore, using named servers may be required when heterogeneous servers are used in the grid. The assumption is typically that only homogeneous servers are used. If servers of fairly different capacity (CPU / memory) are used in the same cluster, named servers should be used to reflect these disparities when assigning applications to server pools.

Role-separated management can be used on top of server pools and is ideally enabled when the cluster is used as a shared Grid Infrastructure, serving as a foundation for more than one organizational group in the company.

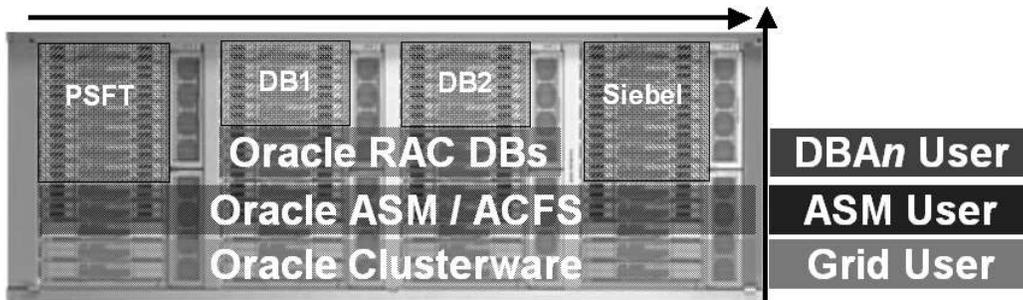
Companies that have implemented a strict separation of duty between their IT organizations can now reflect this separation, even though these groups will concurrently operate on a shared Grid Infrastructure. Role-separated management ensures that different (organizational) groups will not fight over resources by restricting each group to manage and operate only within those (hardware) resources that have been assigned to the them.

In Oracle Clusterware 11g Release 2, role-separated management comes in two modes, as illustrated in figure 4. The vertical mode, as in previous Oracle releases, uses Operating System user/group credentials when implementing the Oracle RAC stack. For Oracle Clusterware 11g Release 2, a dedicated group assignment is used.

The “horizontal mode” is new in Oracle Clusterware 11g Release 2 and enables a clear separation between various organizational groups sharing the same Grid Infrastructure. The new approach is based on server pools and basically restricts certain Operating System (OS) users to modify server pools that are assigned to another Operating System user for management.

The Operating System user owning the installation (that was used for the installation of the Oracle Clusterware software) and the master user (administrator or root, depending on the OS) are also administrators in the cluster.

In Oracle Clusterware 11g Release 2, privileges to operate on server pools are assigned using internal Access Control Lists (ACL). Integration with directory services like the Oracle Internet directory is planned for later releases. The default installation would not foresee a horizontal separation of duty.



**Figure 4: Role-separated Management**

### Clusterized (cluster-aware) Commands

Managing large grids requires performing management tasks on remote nodes (remote operations). Managing a grid with hundreds of nodes is only possible, if such operations can be invoked once, but operate on all nodes or certain groups of nodes in the cluster. Oracle Clusterware provides these cluster-aware command line based commands in terms of clusterized commands. These commands, for example, allow checking the state of the cluster by querying each individual server in the cluster using one command only.

### EM-based Resource and Cluster Management

Oracle Enterprise Manager (EM) extends the idea of a centralized cluster management and brings it to a new level. Not only can Oracle EM be used to manage the cluster as a whole, it can also be used to manage all applications in the cluster using one graphical user interface.

In the first release of Oracle Grid Infrastructure, only Oracle Enterprise Manager Database Control can be used to manage the cluster and all its resources – applications and databases – as a whole. This functionality will be adapted in Oracle Enterprise Manager Grid Control as soon as the new release of Oracle Enterprise Manager Grid Control is available.

Using Oracle EM Database Control, however, requires at least one cluster-aware database to be installed in the cluster. This database does not need to be running to perform some resource-related operations in the cluster. Nevertheless, full access to all cluster operations will only be possible, if the database that is used by Oracle Enterprise Manager Database Control is online.

ORACLE Enterprise Manager 11g Database Control

Cluster Database

Cluster: cluster7 >

Manage Resources

Page Refreshed Jun 4, 2009 11:27:29 AM PDT Refresh

Oracle Grid Infrastructure provides High availability framework to protect any application that is registered with Grid Infrastructure. You can Create, Administer and Monitor these Resources using this interface.

Resources 23 ( 19 4 )  
(Including Internal Oracle Resources)

Search  Go Advanced Search

Show Oracle Resources

View Edit Delete Start Stop Relocate Add Resource Add Application VIP

Select All | Select None | Show All Details | Hide All Details

Select	Details	Name	Cardinality	Current State	Target State	Running Hosts	Resource Type	Owner
<input type="checkbox"/>	Show	ApacheVIP	1	↑	↑	stbpo57	app appvip type	root
<input type="checkbox"/>	Show	MyApache	1	↑	↑	stbpo57	cluster_resource	root
<input type="checkbox"/>	Show	myclock	1	↑	↑	stbpo57	cluster_resource	oracle
<input type="checkbox"/>	Show	ora.DATA.dg	Runs on all servers	↑	↑	stbpo55, stbpo56, stbpo57, stbpo58	ora.diskgroup type	oracle
<input type="checkbox"/>	Show	ora.LISTENER.lsnr	Runs on all servers	↑	↑	stbpo55, stbpo56, stbpo57, stbpo58	ora.listener type	oracle
<input type="checkbox"/>	Show	ora.LISTENER_SCAN1.lsnr	1	↑	↑	stbpo56	ora.scan_listener type	oracle
<input type="checkbox"/>	Show	ora.LISTENER_SCAN2.lsnr	1	↑	↑	stbpo58	ora.scan_listener type	oracle
<input type="checkbox"/>	Show	ora.LISTENER_SCAN3.lsnr	1	↓	↓	n/a	ora.scan_listener type	oracle
<input type="checkbox"/>	Show	ora.asm	Runs on all servers	↑	↑	stbpo55, stbpo56, stbpo57, stbpo58	ora.asm type	oracle
<input type="checkbox"/>	Show	ora.eons	Runs on all servers	↑	↑	stbpo55, stbpo56, stbpo57, stbpo58	ora.eons type	oracle
<input type="checkbox"/>	Show	ora.gsd	Runs on all servers	↓	↓	n/a	ora.gsd type	oracle
<input type="checkbox"/>	Show	ora.net1.network	Runs on all servers	↑	↑	stbpo55, stbpo56, stbpo57, stbpo58	ora.network type	root
<input type="checkbox"/>	Show	ora.oc4j	1	↑	↑	stbpo55	ora.oc4j type	oracle
<input type="checkbox"/>	Show	ora.ons	Runs on all servers	↑	↑	stbpo55, stbpo56, stbpo57, stbpo58	ora.ons type	oracle
<input type="checkbox"/>	Show	ora.racftb.db	Runs on server pool(s) ora.RACpool	↓	↓	n/a	ora.database type	oracle
<input type="checkbox"/>	Show	ora.registry.acfs	Runs on all servers	↑	↑	stbpo55, stbpo56, stbpo57, stbpo58	ora.registry.acfs type	root
<input type="checkbox"/>	Show	ora.scan1.vip	1	↑	↑	stbpo56	ora.scan_vip type	root
<input type="checkbox"/>	Show	ora.scan2.vip	1	↑	↑	stbpo58	ora.scan_vip type	root
<input type="checkbox"/>	Show	ora.scan3.vip	1	↓	↓	n/a	ora.scan_vip type	root
<input type="checkbox"/>	Show	ora.stbpo55.vip	1	↑	↑	stbpo55	ora.cluster_vip_net1 type	root
<input type="checkbox"/>	Show	ora.stbpo56.vip	1	↑	↑	stbpo56	ora.cluster_vip_net1 type	root
<input type="checkbox"/>	Show	ora.stbpo57.vip	1	↑	↑	stbpo57	ora.cluster_vip_net1 type	root
<input type="checkbox"/>	Show	ora.stbpo58.vip	1	↑	↑	stbpo58	ora.cluster_vip_net1 type	root

View Edit Delete Start Stop Relocate

Figure 5: Resource Overview in Oracle Enterprise Manager

## Improved Availability – Tuning “Under The Hood”

Oracle Clusterware 11g Release 2 has been improved to provide better availability. For example, a new agent-based monitoring system is used for monitoring all resources. These memory resident agents allow more frequent checks using fewer resources. More frequent checks means faster detection of failures and a faster recovery time. In case of the Oracle listener, the average failure detection time was reduced from 5 minutes to 30 seconds, while the check interval was reduced from every 10 minutes to 1 minute.

In addition, Oracle Clusterware also reduces planned downtime required for software maintenance. Oracle Clusterware 11g Release 2 uses “Zero Downtime Patching for Oracle Clusterware”, “Out-of-Place Upgrade” and “Software Only Install.” to perform all software maintenance operations with a minimum amount of downtime.

### Advanced Availability

Using a new, more advanced High Availability framework, Oracle Clusterware 11g Release 2 is now able to provide a much higher level of availability for any application managed in the cluster. The agent-based monitoring system used for monitoring all deployed resources is only one example; the extended dependency model is another one.

Today’s applications are not only based on simply one component or process anymore. Complex applications, like Oracle Siebel or Oracle Peoplesoft, as well as SAP applications, are based on many components and processes.

To provide High Availability for the whole application means more than one process needs to be considered and monitored constantly. Using a more flexible dependency model, these and other real world configurations are now easily modeled in Oracle Clusterware 11g Release 2.

For application developers, Oracle Clusterware provides cluster API’s to interact with Oracle Clusterware on a programmatic basis. These APIs have now been incorporated into the Agent Development Framework, allowing the development of Oracle Clusterware agents in a much simpler and more efficient way. Alternatively, the default script-agent provided by Oracle Clusterware 11g Release 2 can be used to protect any kind of application as in previous releases.

### Fencing Flexibility and Third Party Cluster Solution Support

Traditionally, Oracle Clusterware uses a STONITH (Shoot The Other Node In The Head) comparable fencing algorithm to ensure data integrity in cases, in which cluster integrity is endangered and split-brain scenarios need to be prevented. For Oracle Clusterware this means that a local process enforces the removal of one or more nodes from the cluster (fencing).

In addition to this traditional fencing approach, Oracle Clusterware now supports a new fencing mechanism based on remote node-termination. The concept uses an external mechanism capable of restarting a problem node without cooperation either from Oracle Clusterware or from the operating system running on that node. To provide this capability, Oracle Clusterware 11g Release 2 supports the Intelligent Management Platform Interface specification (IPMI), a standard management protocol.

In order to use IPMI and to be able to remotely fence a server in the cluster, the server must be equipped with a Baseboard Management Controller (BMC), which supports IPMI over a local area network (LAN). Once this hardware is in place in every server of the cluster, IPMI can be activated either during the installation of the Oracle Grid Infrastructure or after the installation in course of a post-installation management task using CRSCCTL.

Oracle Clusterware continues to support third party cluster solutions under Oracle Clusterware. For certified solutions (certified solutions can be found in My Oracle Support / Certify) Oracle Clusterware will integrate with the third party cluster solution in a way that node membership decisions are deferred to the third party cluster solution. Only, if a decision is not made within a certain amount of time, Oracle Clusterware will perform corrective actions, using one of the fencing mechanisms described.

Maintaining a third party cluster solution under Oracle Clusterware increases the complexity of the cluster stack and makes the cluster management more difficult. Oracle therefore recommends avoiding having more than one cluster solution on the same system. For Oracle RAC environments it is worth noticing that Oracle Clusterware is mandatory and provides all required functionality. No other third party solution should therefore be required.

## Redundant Interconnect Usage

While in previous releases bonding, trunking, teaming, or similar technology was required to make use of redundant network connections between the nodes to be used as redundant, dedicated, private communication channels or “interconnect”, Oracle Clusterware now provides an integrated solution to ensure “Redundant Interconnect Usage”. This functionality is available starting with Oracle Database 11g Release 2, Patch Set One (11.2.0.2).

The Redundant Interconnect Usage feature does not operate on the network interfaces directly. Instead, it is based on a multiple-listening-endpoint architecture, in which a highly available virtual IP (the HAIP) is assigned to each private network (up to a total number of 4 interfaces).

By default, Oracle Real Application Clusters (RAC) software uses all of the HAIP addresses for private network communication, providing load balancing across the set of interfaces identified as the private network. If a private interconnect interface fails or becomes non-communicative, then Oracle Clusterware transparently moves the corresponding HAIP address to one of the remaining functional interfaces.

Oracle RAC Databases, Oracle Automatic Storage Management (clustered ASM), and Oracle Clusterware components such as CSS, OCR, CRS, CTSS, and EVM components employ Redundant Interconnect Usage starting with Oracle Database 11g Release 2, Patch Set One (11.2.0.2). Non-Oracle software and Oracle software not listed above, however, will not be able to benefit from this feature.

## Reboot-less node fencing in Oracle Clusterware 11g Release 2

As mentioned, Oracle Clusterware uses a STONITH (Shoot The Other Node In The Head) comparable fencing algorithm to ensure data integrity in cases, in which cluster integrity is endangered and split-brain scenarios need to be prevented. In case of Oracle Clusterware, this means that a local process enforces the removal of one or more nodes from the cluster (fencing).

Until Oracle Clusterware 11g Release 2, Patch Set One (11.2.0.2) the fencing of a node was performed by a “fast reboot” of the respective server. A “fast reboot” in this context summarizes a shutdown and restart procedure that does not wait for any IO to finish or for file systems to synchronize on shutdown. With Oracle Clusterware 11g Release 2, Patch Set One (11.2.0.2) this mechanism has been changed in order to prevent such a reboot as much as possible.

Already with Oracle Clusterware 11g Release 2 this algorithm was improved so that failures of certain, Oracle RAC-required subcomponents in the cluster do not necessarily cause an immediate fencing (reboot) of a node. Instead, an attempt is made to clean up the failure within the cluster and to restart the failed subcomponent. Only, if a cleanup of the failed component appears to be unsuccessful, a node reboot is performed in order to force a cleanup.

With Oracle Clusterware 11g Release 2, Patch Set One (11.2.0.2) further improvements were made so that Oracle Clusterware will try to prevent a split-brain without rebooting the node. It thereby implements a standing requirement from those customers, who were requesting to preserve the node and to prevent a reboot, since the node runs applications not managed by Oracle Clusterware, which would otherwise be forcibly shut down by the reboot of a node.

With the new algorithm and when a decision is made to evict a node from the cluster, Oracle Clusterware will first attempt to shutdown all resources on the machine that was chosen to be the subject of an eviction. Especially IO generating processes are killed and it is ensured that those processes are completely stopped before continuing. If, for some reason, not all resources can be stopped or IO generating processes cannot be stopped completely, Oracle Clusterware will still perform a reboot or use IPMI to forcibly evict the node from the cluster.

If all resources can be stopped and all IO generating processes can be killed, Oracle Clusterware will shut itself down on the respective node, but will attempt to restart after the stack has been stopped. The restart is initiated by the Oracle High Availability Services Daemon, which has been introduced with Oracle Clusterware 11g Release 2.

## Cluster Time Synchronization and Network Management

Time synchronization between cluster nodes is crucial. While a deviating time between the servers in a cluster does not necessarily lead to instability, asynchronous times can make it harder to manage the cluster as a whole. One reason is that timestamps are written using the local node time. Log analysis can be impacted severely if the times in a cluster deviate significantly.

A central time server in the data center, accessed by NTP, is typically used to synchronize the server times to prevent deviating times between the cluster nodes. It is best to avoid sudden time adjustments on individual nodes, which can lead to node evictions when performed too abruptly.

To make the Oracle Grid Infrastructure independent from (failures of) external resources, the new Oracle Cluster Time Synchronization Service Daemon (OCTSSD) can be used alternatively to synchronize the time between the servers in one cluster.

The Oracle CTSS daemon is always installed and will always be running, but is configured in accordance to the configuration found on the system. If NTP is installed on the system, CTSS is started in an Observer Mode, not synchronizing the time. Only if NTP is not present on any server of the cluster, CTSS will be activated in active mode, synchronizing the time in the cluster, using one server as the reference server.

More flexibility has also been added regarding network management. Technically, the main enhancement is a network resource managed by Oracle Clusterware. As a local cluster resource, it constantly monitors the network on each server. When a network outage is detected, dependent resource like VIPs managed by Oracle Clusterware are informed and failed over to another node, if required.

Oracle Clusterware maintains one network resource per subnet in the cluster. Multiple subnet support is a new feature that facilitates the consolidation of applications and databases in the grid infrastructure. Multiple subnet support enables independent access of applications and databases using different subnets in the cluster, which then appears as an independent environment to both the database and application clients.

## Cluster Health Monitor: Integrated with Oracle Grid Infrastructure

The Cluster Health Monitor (CHM), formerly known as Instantaneous Problem Detector for Clusters or IPD/OS, is designed to detect and analyze operating system (OS) and cluster resource related degradation and failures in order to bring more explanatory power to many issues that occur in clusters where Oracle Clusterware and Oracle RAC are running.

It tracks the OS resource consumption at each node, process, and device level continuously. It collects and analyzes the cluster-wide data. In real time mode, when thresholds are hit, an alert is shown to the operator. For root cause analysis, historical data can be replayed to understand what was happening at the time of failure.

## Managing Oracle RAC Databases using Oracle Clusterware

Oracle RAC is tightly integrated with Oracle Clusterware. This means that Oracle Clusterware is required whenever Oracle RAC databases are used. No other cluster software can replace Oracle Clusterware in the Oracle stack or is required to run with an Oracle RAC database.

Tools like SRVCTL, which are dedicated to the management of Oracle pre-configured resources in the cluster and therefore ease typical administrative tasks for Oracle RAC databases, show the strong integration of both products.

On a lower level, for example, when it comes to fencing and instance evictions, Oracle RAC is considered a special application in the cluster. Other examples for the integration of Oracle Clusterware 11g Release 2 with the Oracle RAC database stack include internal agents, providing additional state information for Oracle ASM and Oracle RAC instances, and Clusterware-managed services as well as the Single Client Access Name (SCAN), which enables client connections to any database in the cluster using a single cluster alias.

## Managing Any Kind of Application using Oracle Clusterware

All kinds of applications and processes can be monitored and managed by Oracle Clusterware as resources in the cluster. Resources can either run as so called local resources on every node in the cluster, providing infrastructure components. Alternatively, they can be defined as cluster resources, in which case they manage an application in a failover manner.

Non-Oracle applications must be managed using the CRSCCTL tool. Oracle Enterprise Manager (EM) database control, however, can be used to manage all types of resources in the cluster.

Oracle Clusterware 11g Release 2 makes it easier than ever to register an application so it can be managed using Oracle Clusterware. All that is required are four steps, including a failover test in the cluster:

1. Create an Application Specific Action Script or Individual Agent
2. Create an Application VIP to access the Application
3. Configure and Register the Application with Oracle Clusterware
4. Check Start / Stop of the Application & Finalize

Oracle Enterprise Manager supports each of these steps for an easy, centralized deployment. Figure 6 for example shows the dialog that supports adding a resource to the cluster. This includes the action script creation and the registration of the resource.

Other dialogs support the Application-VIP creation or the relocation of a resource between servers in the cluster. Most of the cluster related management operations can be performed even, when the database assigned to the Oracle Enterprise Manager Database Control is down.

ORACLE Enterprise Manager 11g Database Control Cluster Database Help

**Add Resource** (Cancel) (Submit)

General Parameters Advanced Settings Dependencies

Name: MyApache

Resource Type: cluster\_resource (View) (Add)

Description: Failover resource for the Apache WebServer

Start the resource after creation

Placement

The following parameters define where the resource would be placed.

Placement:  Anywhere in the Cluster where this resource instance can be started  
 Restrict or Favor the placement of resource to the Server Pools: AppsPool  
 Restrict or Favor the placement of resource to the Hosts:

Placement Policy:  Restricted  
 Favored

Cardinality:  Specify number: 1  
 Set to size of Server Pool(s) on which the resource is running

Degree: 1

Active Placement:  Re-evaluate resource's placement during addition or restart of a cluster node

Action Program

Action Program defines the way to start, stop and check the status of a resource. Action Program could be an executable (Agent File) and / or a script (Action Script) that the Oracle Clusterware can invoke. Action Program should accept 'start', 'stop' or 'check' as argument to perform respective operations. User can implement all these operations using Agent File alone or Action Script alone or using a combination of both (some operations in Agent File and some in Action Script). If both implement the same operation, Agent File operation would override the Action Script operation.

Action Program: Use Action Script

Action Script Name: /myshared/scripts/myapache.scr (Create New Action Script)

Overwrite if already exists (on any node of the cluster)

General Parameters Advanced Settings Dependencies

**Figure 6: Adding a resource to the cluster using Oracle Enterprise Manager**

## Conclusion

Managing large grids, independent of the number of deployed applications can appear to be a challenge. Oracle Grid Infrastructure overcomes these challenges. Using Server Pools, Grid Plug And Play, as well as standardized and integrated management tools, Oracle Clusterware now manages all applications in a grid as if they were running on a single system, providing improved availability, scalability, and flexibility



Oracle Clusterware 11g Release 2  
September 2010  
Author: Markus Michalewicz  
Contributing Authors:

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  
oracle.com



| Oracle is committed to developing practices and products that help protect the environment

Copyright © 2009, Oracle and/or its affiliates. All rights reserved. This document is provided for information purposes only and the contents hereof are subject to change without notice. This document is not warranted to be error-free, nor subject to any other warranties or conditions, whether expressed orally or implied in law, including implied warranties and conditions of merchantability or fitness for a particular purpose. We specifically disclaim any liability with respect to this document and no contractual obligations are formed either directly or indirectly by this document. This document may not be reproduced or transmitted in any form or by any means, electronic or mechanical, for any purpose, without our prior written permission.

Oracle is a registered trademark of Oracle Corporation and/or its affiliates. Other names may be trademarks of their respective owners.