



HIGH AVAILABILITY CLUSTER
MULTIPROCESSING
BEST PRACTICES

March 2005

Contents

Overview	2
Designing High Availability	2
Cluster Components	3
References	9
About the Author	9

Overview

The IBM High Availability Cluster Multiprocessing (HACMP™) product supports a wide variety of configurations, and provides the cluster administrator with a great deal of flexibility. With this flexibility comes the responsibility to make wise choices: there are many cluster configurations that are workable in the sense that the cluster will pass verification and come on line, but which are not optimum in terms of providing availability. This document discusses the choices that the cluster designer can make, and suggests the alternatives that make for the highest level of availability.

Designing High Availability

HACMP is clustering for high availability. Most fundamentally, it provides high availability by using the redundant hardware in the cluster to bring a spare on line in the event of a hardware or software failure. From this, it reasonably follows that every cluster element that is normally in use when an application is running should have a back up - a spare element of the same type. This is often expressed by saying that there should be no "single point of failure" - no hardware or software element for which there is no backup in the cluster.

While the principle of "no single point of failure" is generally accepted, it is sometimes deliberately or inadvertently violated. It is inadvertently violated when the cluster designer does not appreciate the consequences of the failure of a specific component. It is deliberately violated when the cluster designer chooses not to put redundant hardware in the cluster. The most common instance of this is when cluster nodes are chosen that do not have enough I/O slots to support redundant adapters. This choice is often made to reduce the price of a cluster, and is generally a false economy: the resulting cluster is still more expensive than a single node, but has no better availability.

For more information or to comment on this document, please email:

hafeedbk@us.ibm.com

A cluster should be carefully planned so that every cluster element has a backup. Best practice is that either the paper or on-line planning worksheets be used to do this planning, and saved as part of the on-going documentation of the system.

Cluster Components

Here are the recommended practices for important cluster components.

Nodes

HACMP supports clusters of up to 32 nodes, with any combination of active and standby nodes. While it is possible to have all nodes in the cluster running applications (a configuration referred to as "mutual takeover"), the most reliable and available clusters have at least one standby node - one node that is normally not running any applications, but it available to take them over in the event of a failure on an active node.

Additionally, it is important to pay attention to environmental considerations. Nodes should not have a common power supply - which may happen if they are placed in a single rack. Similarly, building a cluster of nodes that are actually logical partitions (LPARs) in a single processor is useful as a test cluster, but should not be considered for availability of production applications.

Nodes should be chosen that have sufficient I/O slots to install redundant network and disk adapters. That is, twice as many slots as would be required for single node operation. This naturally suggests that processors with small numbers of slots should be avoided. Use of nodes without redundant adapters should not be considered best practice. Blades are an outstanding example of this. And, just as every cluster resource should have a backup, the root volume group in each node should be mirrored, or be on a RAID device.

Nodes should also be chosen so that when the production applications are run at peak load, there is still sufficient CPU cycles and I/O bandwidth to allow HACMP to operate. The production application should be carefully benchmarked (preferable) or modeled (if benchmarking is not feasible) and nodes chosen so that they will not exceed 85% busy, even under the heaviest expected load.

Note that the takeover node should be sized to accommodate all possible workloads: if there is a single standby backing up multiple primaries, the takeover must be capable of servicing multiple workloads. On hardware that supports dynamic LPAR operations, HACMP can be configured to allocate processors and memory to a takeover node before applications are started. However, these resources must actually be available, or acquirable through Capacity Upgrade on Demand. The worst case situation – e.g., all the applications on a single node – must be understood and planned for.

Networks

Networks and adapters are the routes by which HACMP transfers heart beats between nodes. This is an area where more is better; the most reliable clusters have three or more networks to ensure that no simple combinations of hardware or software failures prevent heart beats from flowing.

HACMP strongly recommends that there be at least one non-IP network connecting a node to at least one other node. For clusters with more than two nodes, the most reliable configuration provides two non-IP networks on each node. The distance limitations on non-IP links – particularly RS-232 – has often made this requirement difficult to meet. For such clusters, HACMP V5.1 disk heart beating should be strongly considered. Disk heart beating allows the easy creation of multiple non-IP networks without requiring additional hardware. A cluster without at least one non-IP heart beat path from each node should not be considered as “best practice”, it is barely viable.

(The purpose of the non-IP heart beat link is often misunderstood. The requirement comes from the following: HACMP heart beats on IP networks are sent as UDP datagrams. This means that if a node or network is congested, the heart beats can be discarded. If there were only IP networks, and if this congestion went on long enough, the node would be seen as having failed, and HACMP would initiate takeover. Since the node is still alive, HACMP takeover can cause both nodes to have the same IP address, and can cause the nodes to both try to own and access the shared disks. This situation is sometimes referred to as “split brain”. Data corruption is all but inevitable in this circumstance.)

An installation will often find that it must access a particular node in an HACMP cluster, for purposes such as running reports or diagnostics. To support this, the best practice is to define a node

alias for each cluster node. This has the advantage that HACMP will keep that IP address available despite individual adapter failures (provided there are spare adapters on that network). Experience shows that there is some temptation to use a boot or standby address for this purpose; that temptation must be resisted. Such use of a boot or standby adapter – say, as target of a telnet operation – will interfere with the HACMP use of that adapter, and conceivably cause takeover to fail.

Adapters

Each network defined to HACMP should have at least two adapters per node. While it is possible to build a cluster with fewer, the reaction to adapter failures is more severe: the resource group must be moved to another node. AIX 5L™ V5.2 provides a Network Interface Backup (NIB) facility that can be used to provide particularly fast responses to adapter failures. This must be set up with some care in an HACMP cluster; the appropriate documentation should be consulted. When done properly, this provides the highest level of availability against adapter failure.

Many IBM @server® pSeries® processors contain built-in Ethernet adapters. If the nodes are physically close together, it is possible to use the built-in Ethernet adapters on two nodes and a "cross-over" Ethernet cable (sometimes referred to as a "data transfer" cable) to build an inexpensive Ethernet network between two nodes for heart beating. Note that this is not a substitute for a non-IP network.

Some adapters provide multiple ports. One port on such an adapter should not be used to back up another port on that adapter, since the adapter card itself is a common point of failure. The same thing is true of the built-in Ethernet adapters in most @server p5 and pSeries processors and currently available blades: the ports have a common adapter. When the built-in Ethernet adapter can be used, best practice is to provide an additional adapter in the node, with the two backing up each other.

Applications

The most important part of making an application run well in an HACMP cluster is understanding the application's dependencies. That includes both the resources that HACMP directly manipulates - such as IP addresses, volume groups and file systems - and those that it does not - such as configuration information. The latter is often a source of problems in clusters: if the configuration information is not kept on a shared volume group, it is easy to forget

to update it on all cluster nodes when it changes. This can prevent the application from starting or working correctly, when it is run on a backup node. Hence, best practice in this area is to keep all application configuration information on a shared volume group in the resource group for the application.

The above recommendation may prove infeasible for some applications and installations. In that case, the HACMP V5.2 File Collections facility should be used to keep the relevant configuration information in sync across the cluster.

HACMP provides the ability to monitor an application. This can either check for process death, or run a user-supplied monitor method. The latter is particularly useful when the application provides some form of transaction processing - a monitor can run a null transaction to ensure that the application is functional. Best practice for applications is to have both process death and user-supplied application monitors in place.

HACMP supplies a number of tools and utilities to help in customization efforts like pre- and post- event scripts. Care should be taken to use only those for which HACMP also supplies a man page – those are the only ones for which upwards compatibility is guaranteed.

Even the most carefully planned and configured cluster will have problems if it is not well maintained. A large part of best practice for an HACMP cluster is associated with maintaining the initial working state of the cluster through hardware and software changes.

Testing

Simplistic as it may seem, the most important thing about testing is to actually do it.

A cluster should be thoroughly tested prior to initial production (and once cverify runs without errors or warnings). This means that every cluster node and every interface that HACMP uses should be brought down and up again, to validate that HACMP responds as expected. Best practice would be to perform the same level of testing after each change to the cluster. HACMP V5.2 provides a cluster test tool that can be run on a cluster before it is put into production. This will verify that the applications are brought back on line after node, network and adapter failures. The test tool should be run as part of any comprehensive cluster test effort.

Additionally, regular testing should be planned. It's a common safety recommendation that home smoke detectors be tested twice a year - the switch to and from daylight savings time being well-known points. Similarly, if the enterprise can afford to schedule it, node fallover and fallback tests should be scheduled bi-annually. These tests will at least indicate whether any problems have crept in, and allow for correction before the cluster fails in production.

On a more regular basis, `clverify` should be run. Not only errors but also warning messages should be taken quite seriously, and fixed at the first opportunity. (HACMP V5.2 automatically runs `clverify` daily. Installations at that level should make a practice of checking the logs daily, and reacting to any warnings or errors.)

Maintenance

Prior to any change to a cluster node, take an HACMP snapshot. If the change involves installing an HACMP, AIX 5L or other software fix, also take a `mksysb` backup. On successful completion of the change, use SMIT to display the cluster configuration, print out and save the `smit.log` file. The Online Planning Worksheets facility can also be used to generate a report of the cluster configuration.

Enterprises that have a number of identical or nearly identical clusters should, as best practice, maintain a test cluster identical to the production ones. All changes to applications, cluster configuration, or software should be first thoroughly tested on the test cluster prior to being put on the production clusters. The HACMP V5.2 cluster test tool can be used to at least partially automate this effort.

Software maintenance or upgrades (AIX 5L, HACMP or application) should be first applied to a standby node (after taking the above-mentioned backups). Once that node has been rebooted (which should be done whether or not the maintenance specific instructions call for it), the resource group for the application should be moved to that node. In the event of immediate difficulties, the application can then be moved back to the original production node. If the application runs without obvious problems, it should be allowed to continue on the standby node for some period of time before the production node is upgraded. Note that while HACMP will work with mixed levels of AIX 5L or HACMP in the cluster, the goal should be to have all nodes at exactly the same levels of AIX 5L, HACMP and application software. Additionally, HACMP prevents changes to the cluster configuration when mixed levels of HACMP are present.

Change control is vitally important in an HACMP cluster. In some organizations, databases, networks and clusters are administered by separate individuals or groups. When any group plans maintenance on a cluster node, it should be planned and coordinated amongst all the parties. All should be aware of the changes being made to avoid introducing problems. Organizational policy must preclude “unilateral” changes to a cluster node. Additionally, change control in an HACMP cluster needs to have a goal of having all cluster nodes at the same level. It is insufficient to upgrade just the node running the application.

To this end, the best practice is to use the HACMP C-SPOC facility for any change to shared volume groups. If the installation uses AIX 5L password control on the cluster nodes (as opposed to NIS or LDAP), C-SPOC should also be used for any changes to users and groups. HACMP will then ensure that the change is properly reflected to all cluster nodes.

Monitoring

HACMP provides a rich set of facilities for monitoring a cluster. The actual facilities used may well be set by enterprise policy (e.g., Tivoli® is used to monitor all enterprise systems). In the event that there is no such policy, clstat should be used. That is, it should be running all the time to allow easy determination of cluster state. Furthermore, HACMP can invoke notification methods (such as a program to send a message or e-mail) on cluster events, or even send a page. Best practice is to have notification of some form in place for all cluster events associated with hardware or software failures.

IBM's HACMP product was first shipped in 1991 and is now in its 14th release, with over 60,000 HACMP clusters in production world wide. It is generally recognized as a robust, mature high-availability product. For more information about HACMP, contact your IBM Representative, or visit:

<http://www.ibm.com/servers/eserver/pseries/ha/>

References

IBM Web Pages:

HACMP for AIX 5L Version 5.2

http://www.ibm.com/servers/aix/products/ibmsw/high_avail_network/hacmp.html

IBM Learning Services Classes:

pSeries Logical Partitioning (LPAR) for AIX 5L, course code Q1370

HACMP System Administration I: Planning and Implementation, course code Q1554

HACMP System Administration II: Maintenance and Migration, course code Q1557

HACMP System Administration III: Problem Determination and Recovery, course code Q1559

HACMP System Administration II: Master Class, course code Q1556

About the Author

Tom Weaver

Tom Weaver has been the HACMP product architect since 1994. In that role, he has been responsible for defining the content and direction of the HACMP product. He regards it as the pinnacle of a career replete with opportunities to improving the reliability, availability and serviceability of IBM products.



© IBM Corporation 2005

IBM Corporation
Systems and Technology Group
Route 100
Somers, New York 10589

Produced in the United States of America
March 2005
All Rights Reserved

This document was developed for products and/or services offered in the United States. IBM may not offer the products, features, or services discussed in this document in other countries.

The information may be subject to change without notice. Consult your local IBM business contact for information on the products, features and services available in your area.

All statements regarding IBM future directions and intent are subject to change or withdrawal without notice and represent goals and objectives only. IBM, the IBM logo, the e-business logo, @server, AIX 5L, HACMP, pSeries, Tivoli are trademarks or registered trademarks of International Business Machines Corporation in the United States or other countries or both. A full list of U.S. trademarks owned by IBM may be found at:
<http://www.ibm.com/legal/copytrade.shtml>.

Other company, product, and service names may be trademarks or service marks of others.

Information concerning non-IBM products was obtained from the suppliers of these products or other public sources. Questions on the capabilities of the non-IBM products should be addressed with those suppliers.

The IBM home page on the Internet can be found at:
<http://www.ibm.com>.

The IBM @server p5 and pSeries home page on the Internet can be found at:
<http://www.ibm.com/servers/eserver/pseries>.