

Using Geographic Clustering For Disaster Recovery

By Thomas Casey and Herb Linnell

The GeoHA software, developed by CLAM Associates, provides a disaster-tolerant computing environment for critical data and applications. GeoHA works with IBM's HACMP software to provide automated disaster recovery. In this article we examine how the GeoHA software provides automated disaster recovery. We review the core high availability services provided by the HACMP software, then describe how the GeoHA software extends this high availability to encompass two sites. We conclude the article by discussing the benefits provided by the GeoHA software.

The '90s have been very kind to Michael Jordan and the Chicago Bulls. They have not been as kind to companies who failed to safeguard mission-critical computing resources from disasters. In this decade, a seemingly unending series of disasters—both natural and man-made—have wreaked billions of dollars worth of destruction. Figure 1 lists some of the major disasters that have struck the United States alone since 1990. Add the worldwide disasters during this time, and the total losses are staggering.

And a disaster does not have to be as catastrophic as an earthquake or hurricane. It can be as mundane as a water main break that washes out a data center. Various types of disasters can threaten the availability of critical computing resources, such as the following:

- ◆ Natural disasters—floods, earthquakes, hurricanes, tornadoes, blizzards, wildfires
- ◆ Infrastructure disasters—power failures, fires, water main breaks

- ◆ Operational disasters—hardware faults, viruses, human error

The bottom line is that anything that causes downtime hurts a business, both in terms of lost revenue and lost opportunity. Clearly, prudent companies must implement a disaster recovery solution that ensures critical data and applications remain continuously available, no matter what. For AIX users, the Geographic High Availability (GeoHA) software is that solution.

The GeoHA software, developed by CLAM Associates, provides a disaster-tolerant computing environment for critical data and applications. GeoHA works with IBM's High Availability Cluster Multiprocessing (HACMP) software to provide automated disaster recovery. HACMP uses loosely coupled clustering technology to



Thomas Casey



Herb Linnell

Major U.S. Disasters of the 1990s

| Disaster | Year | Estimated Losses |
|------------------------|------|------------------|
| Hurricane Andrew | '92 | \$24 billion |
| Los Angeles earthquake | '94 | \$20 billion |
| Midwest floods | '92 | \$12 billion |
| Southeast drought | '94 | \$3 billion |
| Hurricane Opal | '95 | \$2 billion |
| Los Angeles riots | '92 | \$2 billion |
| California wildfires | '91 | \$1.7 billion |
| Hurricane Iniki | '92 | \$1.6 billion |
| California floods | '95 | \$.7 billion |

Source: Contingency Planning Research

Figure 1. Major U.S. disasters

prevent individual components, including processors, networks, and disks, from being single points of failure within a cluster.

GeoHA extends HACMP to encompass two physically separate data centers, or sites. GeoHA, by using two sites, prevents an individual site from being a single point of failure within the cluster. Each site maintains an updated copy of essential data and can run key applications, ensuring that mission-critical computing resources remain continuously available at a geographically separate location should a failure or disaster disable an entire site. HACMP prevents a failure in the computer room from disabling a business. GeoHA prevents a site failure from disabling a business.

Figure 2 shows the increasing levels of protection provided by a stand-alone system, an HACMP cluster, an HACMP cluster with RAID, and a GeoHA geographic cluster.

In this article we examine how the GeoHA software provides automated disaster recovery. We begin by reviewing the core high availability services provided by the HACMP software. Next, we describe how the GeoHA software extends this high availability to encompass two sites. We conclude the article by discussing the benefits provided by the GeoHA software.

HACMP

A highly available HACMP cluster is a group of independent processors networked together, sharing critical resources, that cooperate to provide application services to clients. A cluster manager agent, which runs on each processor, is the central control mechanism for providing high availability.

The cluster manager monitors local hardware and software subsystems, and tracks the availability of other processors in the cluster. When a monitored component fails, the cluster manager detects the loss and shifts that component's workload to a designated alternate in the cluster. The HACMP software prevents individual components—including processors, networks, and disks—from being "single points of failure" within a cluster. For example, the cluster manager monitors network adapters by sending keep-alive packets every half second over the service adapters. The cluster manager uses the presence or absence of keep-alive activity as a sign of the adapter's health. If there is no keep-alive activity for five seconds, the cluster manager instructs a standby adapter to take over the IP address of the service adapter.

Once detected, swapping adapters take about three seconds.

Since becoming generally available in 1992, HACMP has been installed in thousands of sites. A mature product, HACMP provides extensive processor, disk subsystem, and network support.

GeoHA

The GeoHA software integrates into HACMP's well-established high availability framework. GeoHA uses the standard HACMP protocol to provide high availability services across a geography. For example, GeoHA uses HACMP daemons and scripts to monitor cluster components and to drive the failover and recovery process. The benefit to users is familiarity. Since GeoHA builds upon HACMP, the learning curve is significantly reduced. Figure 3 shows an example of a GeoHA geographic cluster.

Now let's take a look at this geographic cluster to see how it differs from a standard HACMP cluster. GeoHA adds three components to provide disaster tolerance. These added components are as follows:

- ◆ A second site
- ◆ Geographic networks connecting the sites
- ◆ Geographic mirrors that shadow data entered at one site to the second site

| Levels of Availability | | |
|--|--|---|
| System Type | Level of Resilience | Configuration Increment |
| Stand-alone | Vulnerable to operator error, software defects, hardware failure, single-building mishaps, metropolitan-area mishaps | |
| HACMP cluster | Survives operator error, most administrative assault, processor failure | Cluster hardware and software |
| HACMP cluster with RAID | Survives operator error, most administrative assault, protected hardware failure with RAID | Cluster hardware and software, RAID controllers, 25% disk increment |
| GeoHA geographic cluster with mirrored disks | Survives operator error, most administrative assault, protected hardware failure, plus single-building and metropolitan-area mishaps | Cluster hardware and software, GeoHA software, high-speed point-to-point geographic link, 100% disk increment |

Source: Gartner Group

Figure 2. Levels of availability

A GeoHA Geographic Cluster

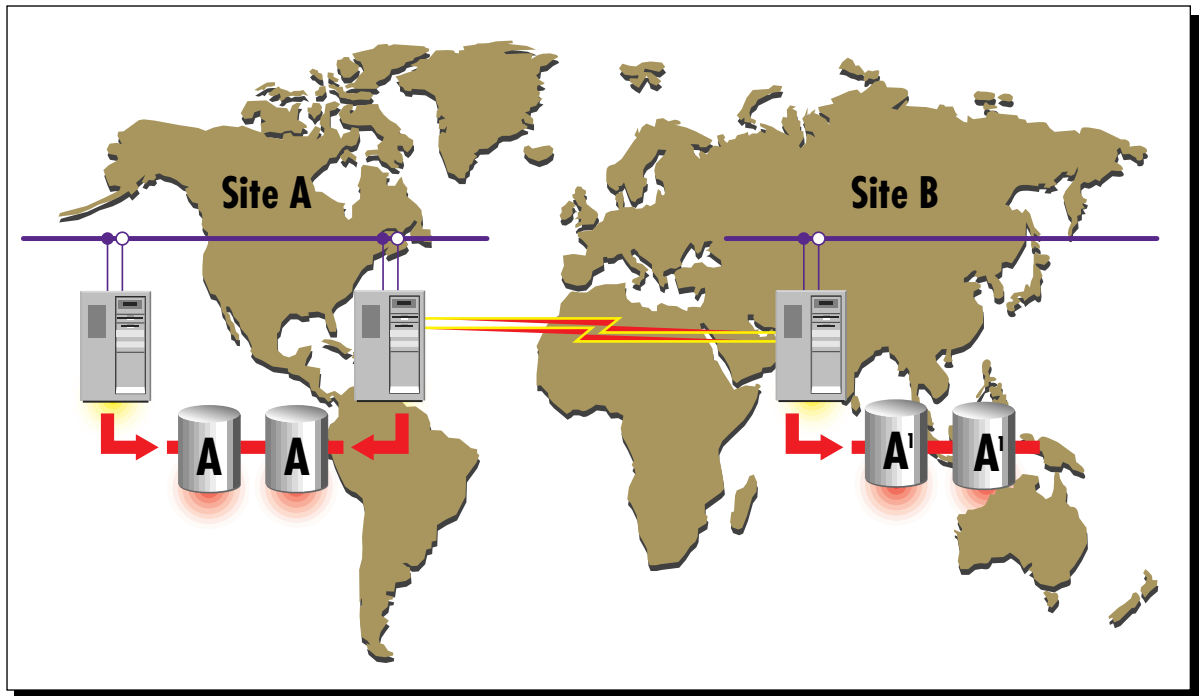


Figure 3. A GeoHA geographic cluster

A *site* is a data center that runs the GeoHA software. A GeoHA geographic cluster includes two sites. The sites should be separated by enough distance to prevent the same disaster from disabling both sites. An individual site can have from one to seven nodes; the geographic cluster can have a total of eight nodes. It is important to note that, together, the two sites form a single HACMP cluster.

The sites in this example are connected by multiple geographic networks. A *geographic network* is a TCP/IP network dedicated to mirroring data between the two sites within the GeoHA cluster. It is also used for HACMP heartbeat traffic. Clients do not have access to the geographic networks. Ideally, a GeoHA cluster should have a minimum of two geographic networks that follow diverse routes to connect the sites. This prevents the same disaster from disabling all the network connections between sites.

Each site has geographic mirrors. Mirroring the data across sites ensures that, even if a site fails, an up-to-date copy of the data is available on the other side. During failover, this copy can be brought online so that operations can

continue. The GeoHA software supports all HACMP cluster configurations:

- ◆ Hot standby configuration where a remote site backs up a local site
- ◆ Mutual takeover configuration where both sites are active and each can take over for the other if it should fail
- ◆ Concurrent access configuration where a single disk image is shared by the nodes across the geography (supported only in synchronous mode)

Providing Disaster Tolerance

Now, having seen what a GeoHA cluster looks like, let's take a look at how the GeoHA software can continue to provide mission-critical data and application services even if an entire site is disabled. To provide this protection, GeoHA must do the following:

- ◆ Mirror data at both sites
- ◆ Monitor both sites to detect failures

- ◆ When a site fails, track data written to the local site so that it can be propagated to the remote site during reintegration
- ◆ Reintegrate the failed site with an accurate copy of the data

Mirroring Data Between Sites

Essential to providing disaster tolerance is having up-to-date copies of critical data at both sites. GeoHA provides a geographic mirroring device that distributes disk I/O transactions to both sites in the geographic cluster, and a geographic messaging system that ensures the data is delivered. Data written at a local site is transmitted across a geographic network to the remote site, so that both sites have an updated copy of the data. Then, if one site becomes inaccessible, the data can be accessed at the remaining site, allowing operations to continue. The key point to understand is that the geographic mirroring system does not simply copy data from point A to point B. The geographic mirroring system must also:

- ◆ Maintain the integrity of the data
- ◆ Track data changes made to the primary site that have not yet been duplicated to the remote site
- ◆ Deliver the data quickly and efficiently across the geographic network

The geographic mirroring device is a pseudo-device driver layered above the AIX Logical Volume Manager (LVM) or a disk device. The geographic mirroring device has two components: a local geographic mirror device that sends data across the geography to the remote site, and a remote geographic mirror device that receives the data. In this schema, the local site initiates a transaction while the remote site receives disk block transmissions from the local site. The roles of local and remote change, depending on the site that initiates the transaction.

Synchronous and Asynchronous Geographic Mirroring

Data may be mirrored synchronously or asynchronously. With synchronous mirroring, the data is actually written first at the remote site, then at the local site. No more input is handled until both writes are complete. Synchronous

mirroring provides the highest data integrity between sites; the local and remote sites have exactly the same data at all times. Note, however, that application performance at the local site may be slowed somewhat by synchronous mirroring.

With asynchronous mirroring, the data is first written to the local site, and then sent to the remote site. Input can continue at the local site while the previously entered data is mirrored to the remote site. Asynchronous mirroring allows the remote site to lag behind the writes to the local disk. While this may improve performance, it is possible to lose the data queued at the remote site if a disaster strikes.

Ensuring Data Delivery

Geographic messaging is a kernel Remote Procedure Call (RPC) protocol that guarantees the delivery of geographic mirroring data should a network fail, provided that at least one alternate, functioning network exists in the GeoHA environment. Geographic messaging automatically re-routes requests in the event that a request fails to reach a destination node.

Monitoring Both Sites

Earlier, we said that geographic clustering was an extension of highly available clustering. A geographic cluster has additional software logic that allows the cluster manager to extend its coverage to encompass a remote site. In highly available clustering, the cluster manager on a processor exchanges “heartbeats” with the cluster managers on the other nodes in the cluster so that they can detect a change in the health of a particular processor. In geographic clustering, the cluster managers at a site exchange heartbeats not only with the processors at the local site, but also with the cluster managers at the remote site. In this way they can determine when there has been a change in the status of the peer site.

Responding to Failures

The purpose of the GeoHA software is to keep key data and application services available even if a disaster destroys one of the sites in the cluster. Now let’s take a look at what happens when something goes wrong. First we will look at what happens when a specific component fails, then we will look at what happens when a disaster disables an entire site.

**Essential
to providing
disaster tolerance
is having
up-to-date copies
of critical data
at both sites.**

Handling Failures within a Site

A *local failure* is the failure of a specific system component within a specific site within the geographic cluster. The component could be a processor, a disk or disk adapter, a local area network, or local area network adapter. The basic facilities of highly available clustering are used to handle local failures.

Each system component (including networks, adapters, and disks) that is a potential “single point of failure” has an automatic replacement designated for it. The cluster manager monitors the status of each designated component, detects the failure of one of these components, and redistributes that component’s workload to a backup component to maintain the availability of cluster resources. In this way, component failures are handled within site boundaries, and are transparent to the other site in the geographic cluster.

Handling Site Failures

Unlike the failure of a system component, which was contained within a specific site, a disaster requires that the viable site take over for the site that has been disabled. The viable site must be able to do the following:

- ◆ Detect the failure
- ◆ Track data written to the local site so that it can be propagated to the remote site during reintegration
- ◆ Continue to provide mission-critical data and application services
- ◆ Reintegrate the failed site when it comes back online

These capabilities are what distinguishes geographic clustering from remote mirroring or data replication. Geographic clustering software is able to respond intelligently to changes in the status of the cluster so that it can maintain high availability. No operator intervention is required to restore services.

Earlier we said that the cluster managers exchange heartbeats across the geography. When the cluster managers at a local site cannot communicate over any of the geographic networks connecting the two sites, they assume the remote site may have failed. The cluster managers at the local site must then determine if the failure is a

global network failure or a site failure, and take the appropriate action to preserve data integrity between the sites.

Site Isolation

Site isolation, or global network failure, occurs if all the networks used to transmit geographic mirroring data are disabled. Despite global network failure, the GeoHA cluster may still be sending keep-alives over a serial network to determine that the remote site is functioning. When site isolation occurs, the site configured as the dominant site continues functioning; the other site will bring itself down gracefully to preserve data integrity between the sites. When the networks are functioning again, the non-dominant site may rejoin the cluster; the GeoHA software then synchronizes the data.

Site Failure

The failure of all nodes at a site is a site failure. The HACMP software on the viable site initiates the takeover of the resources of the failed site. In a concurrent access configuration, no transfer of ownership is necessary. In a cascading configuration, any resources defined for takeover will be taken over by the nodes on the viable site. The nodes at the takeover site mark the last geographic mirroring transaction completed, and then keep track of all data entered after communication with the other site stopped.

Suppose, for example, you have two sites, Toronto and Vancouver, with two nodes at each site. Toronto has nodes alpha and beta; Vancouver has nodes gamma and delta. The resources owned by node alpha at Toronto should be configured to cascade to node beta for a local failure and to node gamma (Vancouver) if node beta is not available. Similarly, the resources owned by node beta will cascade to node alpha or to node delta if node alpha is not available.

Reintegrating a Failed Site

When a site reintegrates, the updates that occurred during the failure need to be resynchronized between the two sites. This process is handled by the geographic mirroring component. When the first remote node sends the message that it is ready to rejoin the cluster, the geographic mirroring device on the functioning site delays the regular HACMP reintegration process until resynchronization completes.

Geographic clustering software is able to respond intelligently to changes in the status of the cluster so that it can maintain high availability.

The nodes at the viable site continue to process data while HACMP is bringing up-to-date the node that is rejoining. When the remote node successfully rejoins the cluster, all of its configured HACMP applications are then available. Once the first node is up, the site can continue operations. The time for synchronization depends on several factors:

- ◆ Number of geographic mirrors involved in the process
- ◆ Size of the geographic mirroring devices
- ◆ Amount of data that has been updated since the last outage
- ◆ Network bandwidth available
- ◆ Network traffic

Benefits of the GeoHA Software

The GeoHA software provides the following benefits.

Reliable data integrity and consistency. GeoHA's geographic mirroring and geographic messaging components ensure that if a site fails, the failed site's data is consistent with the surviving site's data. When the failed site reintegrates into the cluster, GeoHA updates that site with the current data from the operable site, once again ensuring data consistency.

Automatic failure detection. The GeoHA software automatically detects a site failure and initiates the recovery process, including notification that a site has failed.

Automatic or manual fallover. When a site fails, you can manually execute a process to transfer control of data and applications to the operable site. Or you can use HACMP's recovery scripts to automate this process. Using HACMP results in less costly downtime while eliminating the possibility of inducing system failure during recovery.

Fast recovery. The GeoHA software provides fast recovery of mission-critical data and applications at the operable site. The GeoHA software can reduce the time needed to restore computing resources to minutes, although individual times can vary depending on the amount of resources that need to be shifted to

the surviving site and the amount of application recovery processing required.

Flexible configurations. The GeoHA software supports a wide range of configurations, allowing you to configure the disaster recovery solution unique to your needs. Configurations can range from an online backup machine turned on nightly to receive an updated copy of a database to a concurrent access configuration, where both sides have simultaneous access to the same database.

Geographically separate sites. The GeoHA software imposes no constraints on the distance between the sites. This allows you to separate the sites by enough distance to ensure that the same disaster does not render both sites inoperable.

Flexible, scalable architecture. You can scale a GeoHA geographic cluster simply by adding memory and I/O controllers to individual processors, by swapping in more powerful processors, or by adding more processors. A GeoHA cluster can support up to eight nodes.

Based on industry-leading technology. The GeoHA platform is built from proven components—cost-effective, high-performance RISC System/6000 servers, and the industrial-strength AIX operating system.

Filesystem and database independent. The geographic mirroring device behaves the same as the disk devices it supports. Because the mirroring is transparent, applications configured to use geographic mirroring do not have to be modified in any way. In this sense GeoHA is a “generic” solution that works with any file-system or database management system.



Thomas Casey, CLAM Associates, Inc., 101 Main Street, Cambridge, MA 02142. Internet: tom@clam.com. Mr. Casey is the manager of CLAM's technical writing group. He has a BS from Trinity College and an MS from Emerson College, both in Boston.

Herb Linnell, CLAM Associates, Inc. 101 Main Street, Cambridge MA 02142. Internet: herb@clam.com. Mr. Linnell is a senior member of CLAM's product development staff. He has a BSCE from the University of New Hampshire in Durham and a MSCE from Northeastern University in Boston.

GeoHA is a “generic” solution that works with any filesystem or database management system.