# RS/6000 and AIX clusters:  past, present and future

*July, 2000*

**Contents**

## Introduction

This document is divided into three sections.The first assumes knowledge of IBM systems or clustering, the latter two assume some familiarity with IBM's cluster offerings.

Ask the users of a computer system what they expect from it and you'll probably be told that as well as providing the required services, they want it to be running when they need to use it and that they want it to perform well. Ask the management responsible for the same system what they want, and they'd probably add to the list, the requirements that the system be easy to run and that the lifetime of the system be as long as possible, given the rate of change in technology.These four requirements, that the systembe available, scalable, easy to manage and that the investment in it be protected are common across just about all computing platforms and all of these are under increasing pressure.

As industries globalize and as competitive needs are pushing companies to do business round the clock, availability is becoming more critical. Computers are not less reliable than they were, rather it is that the business requirements mean that they must stay running for longer. Many companies canno longer function without their IT systems. And it's not just large enterprises or the systems that run the business that fall into this category. Small companies, department servers and even workstations now have availability requirements far in excess of where they would have been even a few years ago. When key systems are unavailable, businesses are unable to function andose money as a direct result. Keeping systems running is important today and will become even more so in the future.

Software packages are rapidly adding more and more function between releases. In the past, when many applications were developed in-house, the pace of change was slower and the requirements placed on the computer systems by the applications grew at a much slower pace, one linked closely to the development cycle. As companies have moved away from in-house software and now utilize many more off-the-shelf packages, they are at much greater risk that their computer systems will not be able to cope with the next version. One has only to look at the PC arena to see this happening with a vengeance. The difference in system requirements between the first and second versions of a particular software package might require a doubling in both processor speed and memory capacity! In a multi-user or client-server environment, this might translate into being unable to support all of your existing user base on your current platform. So, the requirement to upgrade the hardware is likely to occur on a more frequent basis.

Which brings us neatly to investment protection Companies want to know that the systems they purchase have a lifetime longer than the next software release. They need to know that they can expand capacities and change the capabilities of the systems as the business needs change, rather than being forced to replace the equipment they have just bought. When they do decide to purchase a new system, they need to know that it will seamlessly integrate andnteroperate with their existing equipment such that they can still continue to gain business benefit from their investmentBut it is not just hardware that comes into play here. It is much more expensive to replace the skills and experience of computer staff with years of experience in running or managing a computer system. This can be lost or made obsolete overnight if the system changes dramatically. The time taken to retrain and the ensuing impacts on the business probably has a higher potential price than the hardware or software.

Changing business requirements force computer systems to change. As new systems are added or new functions are added, the environment becomes more complex. In this day and age though, companies are often forced to use less and less skilled people to manage their computer environments. Other companies are new startups, just getting going in IT and so have not yet acquired the skills they need yet. There is a requirement, often one paid lip service to, that computer systems have to be easy to use. If the environment cannot be configured, managed and run easily, then it is almost impossible to get the best out of it. And if you can't get the best out of the system, then the business benefits of having it, rapidly start to evaporate.

Addressing these four key requirements, availability, scalability, investment protection and ease-of-use are the key to having computer systems that are of real, measurable benefit to the business. If you miss out on one, you're going to be at a disadvantage, unable maybe to be as responsive as your competitors or losing

out if you go after a new opportunity.  Miss out on two or more and any benefit you might hope to derive from these systems is likely to be temporary.  Meeting these requirements and meeting them well has a direct effect on the performance of the system and consequently, the performance of the business.

This document shows how to meet or exceed these requirements to deliver systems that are of real benefit to the business.  It discusses the requirements in more detail, explains how IBM AIX clusters meet these requirements today and how they will continue to do so into the future.
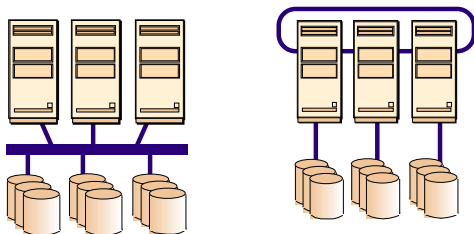

## Why a cluster?

Given these requirements, a cluster of systems is an ideal solution for customers.  A cluster is a group of computers, usually referred to as nodes, that are interconnected to provide a single computing resource.

Clusters offer much higher availability than a single system, allowing freedom from both planned and unplanned outages.  To address the planned outage component, workloads may be migrated to another cluster node whilst administrative change occurs, thus allowing changes to be made whilst still maintaining a useful service.  When the maintenance activities are complete, the workload may be moved back to its original location if desired. To handle unplanned outages, redundant components are used.  Redundancy in the cluster hardware, whether redundant servers, networks, storage, or adapters allow work to continue transparently when one or more hardware or software components fail. A service continues to be provided by simply switching to use other cluster components.  If a node fails, service may be moved to another node.  If a disk fails, service continues from another disk containing the same data, and so on.  A cluster acts as a single, continuously available system in this respect.  Availability of system resources is one of the greatest advantages of clustering.

The cluster model also provides investment protection because clusters offer both horizontal and vertical scalability.  Horizontal scalability is the ability to add more nodes to the cluster to provide more functionality.  These nodes may be relatively small and / or cheap, offering much cheaper upgradeability than might be possible upgrading a single large system, the synergistic effects making the sum offered by the cluster significantly greater than the sum of the component systems.  For those who need larger individual cluster components, vertical scalability is the ability to upgrade individual nodes to higher specifications.  Either of these enhancements may be made without requiring the service to be taken off-line.  Finally, by reusing existing computers in clusters, their useful lifetime may be dramatically extended.

Gathering individual computers together into a cluster can simplify their management too.  Cluster software allows the individual components to be treated as a single system such that changes may be made from a single point, the cluster software taking care of making the appropriate changes on the individual cluster nodes.

What is important here is that this cluster software, whilst providing an interface that insulates the user from the complexity of the cluster, does not isolate them.  It is important to be able to work with the individual components as well as with the cluster as a whole.  Clusters provide a means for easier management of your computing resources, whilst not restricting the potential capabilities of the systems.

Many people think that a cluster is overkill for their requirements or is only for large systems. This could not be further from the truth.  If the system is important to the business then it is a potential candidate to be clustered.  Whether it is a mission-critical enterprise server, or a server for a small work group, clustering

**A collection of 'whole' computers...**



**... used as a single computing resource.**

*Figure 1:  Shared-disk and shared-nothing clusters*

provides availability and scalability.  The smallest work group servers can have enterprise-class availability characteristics via clustering.

The ability to build clusters with small, cheap servers or by reusing servers that are surplus to requirements elsewhere within the business can allow dramatic improvements in availability and scalability throughout the company.

# Part I.  Business requirements

## Availability

Availability is probably the most important property your computer systems have.  If a system is not available to run the workloads it is intended to run and to perform the tasks that are vital to your business, then it doesn't matter how fast it is or how much memory it has, what you have is  an expensive pile of metal and plastic and little else.  For many companies, this is an accepted fact of life.  These companies have been running mission-critical computing environments for years.  For others, whether they are small or new businesses, availability is only now being seen to be critical to their success and competitive edge.

Take for example, the fiercely competitive arena of e-business, if your web site is not available, or you are unable to process an order because your backend server is down, you might, if you're lucky, get a second chance.  You certainly won't get a third.  You will have lost a new potential customer or damaged the relationship you have with an existing one.  It is not just e-businesses that need this level of availability though.  All businesses, large or small need increased availability, just to stay competitive these days.  Running 7x24, seven days a week for twenty four hours a day, used to be only required of the largest of companies, nowadays, most, if not all companies have to, or try to, run their computer systems this way.

Whether your customers are external to the business or users of your own computer systems within it, computer systems outage is a major problem for an increasing number of computing users these days.  These customers are insisting upon greater and greater systems availability.

As more and more companies become totally dependent on their computer systems being available around the clock, so it is critical the these systems stay up and running.  Any period of computer system outage can be directly translated into lost revenue for the business.  Typical outages cost an average company in the order of 10,000 US dollars per minute[1] and in many companies, this amount may be considerably more.

This requirement for continuous access to computer systems spans the entire spectrum of both customers and computer systems.  From the largest corporation to the smallest web startup.  From central-site servers to work group, and even desktop type systems, the need for computer availability has never been greater.

## Outage and causes of outage

Many industry surveys have been carried out to determine the causes of system outage, those periods of time when the system is not available to perform useful work.  The numbers vary between surveys, but the culprits invariably remain the same.

It is important however, to put things into perspective.  The major portion of the outage experienced by a computer system will be planned.  According to IBM surveys, over ninety percent of the outage of a system falls into this category.  This may be time required to perform administrative tasks such as taking backups, apply



Figure 2: Causes of outage

software and hardware maintenance or add new software or hardware components.  The need to remain at the leading edge of your business segment will drive continuous hardware and software change, but no business today can afford the periods of downtime that a few years ago were commonplace.  Increasing globalization and running systems and services around the clock to gain a competitive edge, is reducing or removing those periods of downtime that are required to perform normal systems maintenance.  Nowadays, users often need to keep the system operational whilst this maintenance is being performed.
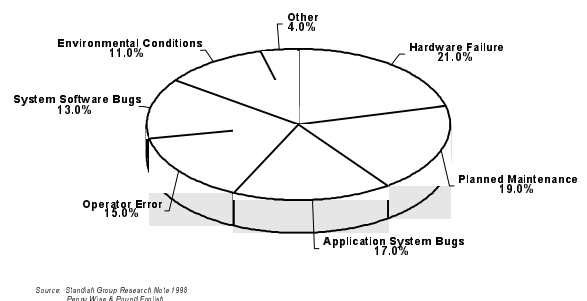
The remaining ten percent of outage is that which is unplanned. This is typically broken down still further into outage caused by software failures, user errors and hardware or environmental problems. Hardware is getting more and more reliable, reaching levels of reliability today that were unheard of, even a few years ago. Components with moving parts, those which are the most likely to fail, are now often achieving a Mean Time Between Failure(MTBF) in excess of one million hours. IBM has been leading in developing such components. Components such as circuit boards and microprocessors have MTBF numbers one or two orders of magnitude higher than these. A typical UNIX server without any additional availability features such are standard on an IBM server is a complex combination of these components and might typically be expected to be available for 99% of a year as a result.

Just as hardware is getting more reliable, software is, on average, becoming less reliable. Almost without exception, the main culprit here is application software[2]. Operating system facilities are exploited by many more users and consequently have been better tested in the real world than most applications. As applications offer more and more features to soak up the rapid growth in system resources, new code needs to be added. More code means more likelihood for bugs, leading in turn to an increased risk of an application software failure. This doesn't mean to say that operating systems are totally without fault, rather that they contribute relatively few failures to the overall software failure number.

User skills are also measurably reducing over time. Much of this, is brought about through the 'point-and-click' mentality of the PC. Whilst the PC has much to commend it, in terms of providing increased ease of use through devices such as 'wizards' and other GUIs, these do tend to isolate the user from the system. When a failure occurs, very often the very GUI that made the activity simple, now acts as a barrier, preventing diagnosis of the problem. Many users, who have never performed systems management outside of a GUI environment are lost when it gets to this point. Often, when placed into a more complex systems environment, their untrained actions cause more harm than good. UNIX systems are famous for their ability to do exactly what you ask them to regardless of whether you actually wanted this done.

Beyond the single site environment, the major causes of outage are those resulting from site failures, such as power outages, fires and floods or from larger scale disasters such as earthquakes, hurricanes, tornadoes etc. Just as over the past few years, companies have been looking at increasing the availability of their servers, having now hardened these to a greater or lesser extent, they are now finding that site failures or disasters pose a very real threat to their competitive edge. Disaster survivability is now something that all companies have to take seriously when considering their availability strategy.

If we look at the causes of outage, it is fairly easy to come up with a description of the sort of system we need, in order to avoid them. For starters, such a system must have the ability to be able to perform administrative tasks, including dramatic changes such as the addition of new hardware or the application of software maintenance, whilst still being able to provide a service to the end-users. If this alone can be addressed, then the vast majority of system outage can be avoided. But we still haven't dealt with the unplanned outage. There are still things that will cause the system to be unavailable. To handle these requires a system which can survive hardware and software failures. Finally, you need a system which is both easy to use, yet doesn't prevent the system administrator from getting at the underlying system when needed. For those environments with relatively unskilled users, having the ability to automate complex functions such that they can be performed error-free by the system when the need arises, is a key benefit too.

Merely having such a shopping list is not enough, any investment in a computing resource needs to be protected from obsolescence. Given the rate at which technology changes, and hence the rate at which it becomes obsolete it is important to ensure that a system will remain useful to the business for as long a time as possible. Better still, that the systems can be grown as the business needs grow. However, growing a system must be possible without increasing complexity.
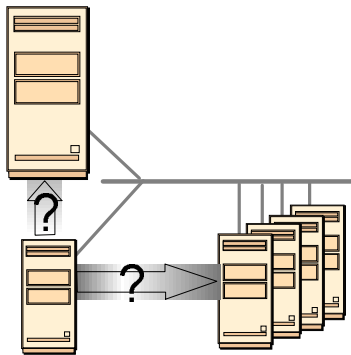
## Scalability

A computer that performs adequately today is not necessarily going to perform adequately in the future. Any resource in a system is finite so sooner or later, it is likely that as workloads increase or as the system is called upon to perform tasks other than those it was originally intended to, resources will get constrained and the performance will be impacted.

Most computer systems offer some means of extending their capabilities, whether this means adding more memory, additional processors to a multiprocessing system or adding further I/O capabilities, be these adapters or disks. This is vertical scalability, the ability to grow a system within the unit itself. Vertical scalability provides some means of putting off the day when the system finally runs out of performance.

But what if you have already bought into the system at the high end of it's capabilities and already have as much memory or as many CPUs as are supported by the system itself. What options are open to you now? What if the bottleneck in the system is not caused by something that you can simply add more of to alleviate? What if the workloads that you are trying to run are competing for resources or are making the system run in an unbalanced fashion?

Horizontal scalability provides the ability to add more computers to the system and to distribute work across them. The whole system thus scales well beyond the limits of a single system. The limit to this form of scalability is then brought about by the applications you decide to run. The simplest way to exploit multiple machines is to move different unrelated workloads onto different systems. In such a system, you might move your office applications onto one cluster node, your personnel applications to another and the production database to a third. However, there are often

*Figure 3: Horizontal and vertical scalability choices*

requirements for interoperability between applications. The scalability here is usually limited by the data that they use. For different applications that require access to the same data, some means of providing data access from multiple machines will be required. This might be via a database hosted on one cluster node and queried from another via a network, via a parallel database, where multiple cluster nodes participate in providing database services, or simply by giving multiple cluster nodes access to shared file systems or raw disk, the technique chosen based upon what the applications require.

Even though this gives considerable relief from the pressures on system performance, it is still not allowing exploitation of multiple systems against a single workload. To do this is perfectly possible, but again we are limited by the capabilities of the application. In this case the degree to which the application is parallelized. Parallel applications run on multiple systems, each instance or component communicating with the others to coordinate the activity of the whole. This allows multiple computers to process a single workload in a very much reduced time compared to a single system and allows previously impossible tasks to be undertaken by exploiting parallelism to it's ultimate limits. These so-called massively parallel or 'grand-challenge' systems are almost always processing scientific or technical workloads, though their use in solving commercial problems, in environments such as data mining is becoming increasingly common.

## Ease-of-use

Modern computer systems are complex things, there is no escaping this fact though much of this complexity is hidden from the user unless they decide to go looking for it. Usability really breaks down into two major areas, that of applications and systems management.

Application ease of use is really down to the application vendor.  Some have made a good job of this, others less so.  For many applications though, the usability is either common across all platforms or takes the look and feel of the underlying operating system.  This either assists or thwarts the application user depending upon their background and experience.   In most environments, the choice of application is a given that little can be done about.  Consequently, there is little that can be done in terms of application usability.

Where ease-of-use has less commonality is in the area of systems management.  Each computer system vendor provides their own unique facilities as well as many ones that are common across multiple platforms.

The management of the system can really be broken down into four main areas.  These are:

- Installation and configuration
- Change management
- Ongoing maintenance
- Problem management

Each of these have their own usability issues and requirements.

The initial installation and configuration of the system should be quick and simple.  It being important to get the system up and performing useful work as quickly as possible.  A system that requires multiple attempts to get working properly will not be generating revenue as quickly as one that works first time.  Similarly, one that takes two months to install will not be as beneficial to the business as one functioning properly in two weeks.

Whether the installation is performed by the customer, the system vendor or a services organization, changes in the future are likely to be made by the customer.  If the customer has less skill than the installer, then there is a potential risk that the system manager may cause more problems than are solved, user error being a significant cause of outage[3].  Changes need to be able to be made, simply and effectively with the minimum risk to the running environment.  Given that you are likely to be clustering for availability, ongoing maintenance activities and day-to-day systems management has to be simple enough to not cause problems for less skilled administrators.  Systems that allow automation of systems management activities or provide facilities to proactively rather than reactively manage a system are preferable too in that once set up, the system can either run 'lights-out' or manage many day-to-day operations itself, this being a good way to minimize risk due to operational errors.

When problems arise, a systems management interface must help, not hinder the identification and fixing of them.  Good tools also allow complex situations to be managed or 'learnt' such that if they happen again in the future, problem resolution is faster.
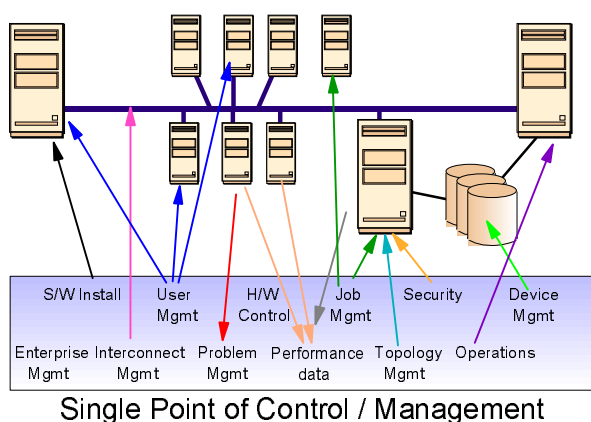
Usability is a combination of tools and familiarity.  Without wishing to get involved in the usual heated debate that seems to break out whenever the merits of one interface vs. another are discussed, let us simply say that what one user prefers to use, whether it is a graphical user interface (GUI), the command line or some other interface, is not necessarily the choice of another.  What is much more important is ensuring that the systems management facilities may be used from each of these different environments.  Whilst it may appear to some users that graphically representing systems management data makes it easier to understand, let it not be forgotten that before the days of GUIs, people managed



Single Point of Control / Management

Figure 4:  Usability improves with a single point of management

systems perfectly well, some might say better, via a 'green screen' interface. It is important to provide a choice of interface rather than forcing a user to take a particular route, that way they are more likely to be comfortable and use the interface rather than avoiding it.

Remember also, that an interface may be used in different ways. If a system provides only a GUI, then managing this system via say a dial-up line and a modem may not be possible. Similarly, if the system is configured with an ASCII terminal as a console, a GUI will not run in this environment too. Again, choice of management interface has to be balanced against the access methods within which it will be used.

In providing interfaces, it is also important to maintain consistency to ensure that skills are not lost or need to be relearnt as different versions of products are implemented. Whilst it is quite reasonable to expect that a new function might add options to a configuration menu or appear as a new icon in a GUI environment, these should be implemented in a common look-and-feel to the existing management environment. Again, familiarity encourages use.

Interfaces that are common across multiple platforms or allow the management of multiple platforms assist in maintaining skills. Similarly low-level tools that can be integrated into a higher-level systems management environment mask much of the apparent complexity. The increasing use of systems management tools which offer 'web-like' or 'PC-like' interfaces, whether web-based or not, is becoming an increasingly common means of achieving management economies. A user trained to manage PCs can contribute significantly more to management of more complex environments if the look-and-feel of the tools used are similar or interrelated. The same can also be said of management through 'higher level' tools such as Tivoli which provide an encompassing framework within which individual management activities may be driven.

It should also be remembered that usability is not only about making things easier. It also encompasses giving access to information. When a problem arises, a GUI is typically the wrong tool to use. Here, it is important to provide access to information in the most appropriate form, whether this be logfiles or trace information to allow a skilled administrator or service personnel to identify the cause of a problem and take remedial actions. Any systems management interface must be able to be bypassed as the need arises to allow access to the underlying system. Insulation not isolation is the key thing to remember here.

## Part II.  AIX clustering today

### Why IBM?

A typical cluster, regardless of it's end-use, is built of a number of hardware components.  In addition to the cluster nodes themselves, these include disk storage and networking hardware.  The cluster is completed with software components to make this collection of hardware function as a cluster.

Not unsurprisingly, the 'best' clusters are usually those which combine the 'best' components, both hardware and software.  The term 'best' is used here advisedly as what is suitable for one customer environment is not necessarily suitable for another.  The key to this is flexibility.  Given a set of requirements, be they for availability, scalability or whatever, it is important to be able to use those components which most closely meet them.

When it comes to AIX® clustering, the node will be an RS/6000® server.  This family of systems has been available for over ten years, with over one million systems sold,  resulting in the most proven UNIX hardware platform in the industry.  Proven by real customer use as well as offering leadership performance.  All RS/6000 servers are symmetric multiprocessor (SMP) capable, though offering uniprocessor entry points in many models.  Capable of running a wide range of applications, both 32 bit and 64 bit, in both commercial and technical environments the RS/6000 is an ideal cluster building block.

With the most powerful model in the family, being some 120 times more powerful than the entry RS/6000, scalability should be of little concern.  Customers can chose those cluster nodes that meet their performance requirements, clustering them in any combination.  In addition, to this horizontal scalability, all RS/6000 models can be upgraded with extra CPUs providing vertical scalability within the same node.  IBM SMPs are well balanced offering excellent near-linear scalability as additional processors and memory are added unlike many other SMP platforms which scale poorly.  Industry standard benchmarks confirm this when 24-way RS/6000 servers are easily capable of outperforming much larger competitive offerings when running real customer workloads.

But performance isn't everything.  As we have seen, it is not much use having a high performance system if it isn't available.  The base RS/6000 server builds upon years of IBM experience in delivering systems for running businesses and uses many unique technologies to enhance the availability of the nodes themselves.  If the components you build the cluster with are themselves exceptionally reliable, the cluster is even more so.

A typical RS/6000 server provides a wealth of hardware availability features.  These include concurrent maintenance and component reassignment, allowing many components to be replaced without the requirement to stop operational use.  This is further enhanced by unique facilities such as Predictive Failure Analysis, whereby components that are about to fail can notify the user.  Action can then be taken to repair or replace the component before a critical failure occurs.  Combining this with the capabilities of the Service Processors in most RS/6000 servers allows management of the system, even when powered down as well as remote management.

Software is needed to run any computer and the operating system is key to delivering a platform that is going to deliver benefit to the business.  For the RS/6000, the operating system is AIX.  The same AIX operating system runs on all models of the RS/6000 unlike some vendors who offer different operating systems on different models in their processor ranges.  AIX is a second generation UNIX operating system, designed for running businesses.  Rated the 'Top UNIX operating system'[5] and the 'Operating System with the Lowest Total Cost of Ownership'[6], AIX is packed with proven facilities.  Whilst pretty much any vendor offering clusters might claim to offer an operating system which meets the availability, scalability and manageability goals, IBM not only can prove that AIX excels here, they set the standards against which all other clustering offerings are judged.

When it comes to providing a system that minimizes outage due to planned outage, it is important to be able to do as much as possible within a single system without requiring a system to be taken down.  AIX

facilities such as the Logical Volume Manager and Journaled file system both enhance availability and remove the requirements for downtime to perform disk administration tasks. Moreover, these come bundled as part of the operating system as are the RAID, striping and mirroring functions built in to protect against disk failures. Many other vendors only offer these as additionally chargeable products. The ability to add and remove devices and subsystems on the fly without requiring a reboot is also key, as are the wide range of diagnostic and fault isolation capabilities. This range of features allows the vast majority of system administration tasks to be performed on a running system.

AIX offers backward binary compatibility between hardware platforms and operating systems. The very earliest RS/6000 servers can still use the latest AIX operating system. The large number of applications written to use AIX can be moved to more powerful systems as businesses grow without the need for recompilation or modification. This is not only an example of software scalability, but also an excellent investment protection strategy for companies to follow. Having scalability is one thing, but efficiently utilizing these system resources is another. To address these issues AIX provides a Workload Manager to ensure that applications which have given resource requirements can be sure of obtaining them regardless of other activities performed on the system.

As well as being scalable and available, the RS/6000 and AIX platform offers a wealth of systems administration and management tools. Being a UNIX system, the command line interface is ever present and can be used either directly or via standard UNIX scripting languages. AIX differs from many modern UNIXs in providing single commands that can perform very complex configuration and management tasks, other systems requiring multiple commands or activities to achieve the same end result. For the majority of day to day tasks, most users prefer to use the other management interfaces, the System Management Interface Tool (SMIT) and the graphical Web Systems Management (WebSM) being the most common. As with any environment, users may chose those tools and facilities most suitable to the task and may mix and match accordingly. Management of AIX systems remotely is possible via WebSM. Alternatively, AIX systems management may be performed via enterprise management tools such as Tivoli offering all of the operational management capabilities expected of environments designed to run and manage entire organizations.

### IBM cluster value proposition

Why should you choose an IBM cluster solution?

Quite simply, experience. When you chose an IBM cluster solution, you are choosing the only solutions built upon over forty years of experience in delivering available, scalable, mission-critical systems to completely meet your business needs. You are buying leadership technology that is proven to work, every time, time after time. Proven by the largest installed customer base, in the largest range of countries, industries and business sizes. We've been doing this for longer than anyone else, so you get a real tried-and-tested solution with an unrivaled track record.

Our experience gives you a better solution. IBM clusters build upon the proven strengths of RS/6000 and AIX with availability, scalability and systems management features built-in, not bolt-on. We understand what it takes to keep a business running and know exactly the facilities you need and how best to deliver them to you. You want a single product, not a set of building blocks or a product you must replace as your needs change. A single product, capable of fully meeting your needs today and in the future. Extendable as your business changes whilst keeping the same proven technology you've invested in and are familiar with. Moreover, IBM owns it's cluster products. It controls the speed and direction of the offerings to be as responsive as possible. We are not at the mercy of other companies when it comes to meeting customer demands. Whether you need availability for a single department or disaster recovery for the entire enterprise. Whether you need a loosely-coupled cluster of workstations or a worldbeating massively parallel computer, IBM solutions are always the ones that lead, that the competition try to beat, offering the most function at the lowest cost.

IBM's experience is not just technology, it's also about people.  People who understand your business and know that you want to be making money, not running systems.  People who know what it takes to design and build you a superior system that works first time, without any hassle.  People who have decades of experience in developing solutions for keeping businesses running.  People who know that however good the technology, without processes and management disciplines, you will not realise the full potential of an IT solution.  All of this is backed by the truly worldwide service and support you need when running a mission-critical system.

All this experience is yours when you invest in an IBM cluster solution.  Your business will reap the benefits sooner and quite simply will stay up and running, doing what it was designed to do, for longer.
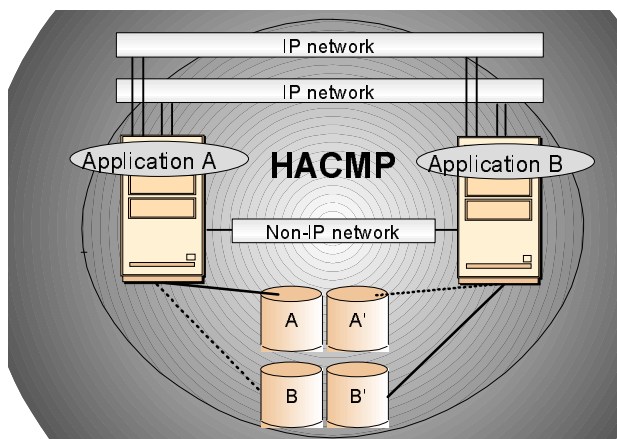
## IBM's AIX clustering products

IBM has been delivering clusters of RS/6000 servers to meet our customer's computing needs for almost ten years, far longer than any of our competitors have had similar offerings.  Given that customers typically cluster systems to increase availability, increase performance and simplify systems management, different solutions have been developed in each of these areas.  This section outlines the main IBM offerings in these areas.

### Clustering for availability

When it comes to availability, the cornerstone of IBM's AIX clustering strategy, today, as in the past is HACMP (High Availability Cluster Multi-Processing). This, industry-leading clustering solution has always been the one that our competitors have been aspiring to beat.  It has consistently delivered the capabilities that our customers demand for high availability.  More importantly, it is a truly proven solution with almost 40,000 licenses having been installed in over 16,000 customers worldwide.

The configuration of a typical two-node HACMP cluster is shown below.  Two RS/6000 servers are connected to each other with a variety of networking connections and they are both connected to a set of 'shared' protected disks.  Every entity in the cluster is regarded as a resource.  These resources include, applications, IP addresses, disks, file systems and so on.  If a node in an HACMP cluster fails, the resources it 'owned' and hence the services it was providing are automatically relocated to a backup node within the cluster.  The resources associated with an application, the disks and file systems where it's data is located, the application itself and the IP address used to connect to it are all grouped together to simplify management.  This resource group is the entity moved around the cluster



Figure 5:  A simple two node HACMP cluster

Even with a simple two node cluster, there are many possibilities.  A cluster such as this might be configured in an 'Idle Standby' configuration.  Here, one node is active, running the applications.  The backup node is idle.  In this case, when the failed node is restored, the workload is migrated back to it.  If  you don't want this to happen, a variant of this is a 'Rotating Standby' configuration.  Here, the rejoining node acts as the backup to the existing live server.  For those wishing to make more effective use of their computing resources, the backup may be running a non-critical workload which is stopped when the production node fails, a so called 'SimpleFallover'

13

configuration.  Finally, both nodes may be running production workloads in a 'Mutual Takeover' configuration, each acting as a backup to the other.

As the number of cluster nodes increases, so do the possibilities.  With up to thirty two nodes able to be configured in a cluster, extremely cost effective configurations, such as one machine acting as a backup to thirty one others are also possible.

Configuration flexibility is a key HACMP differentiator. Any RS/6000 server can exist in a cluster with any other RS/6000 servers.  Just because the RS/6000 is participating in a cluster doesn't mean that you have to forego any of the advanced availability features of the servers themselves.  Features such as service processors, hot-pluggable components, redundant power and cooling all go to make RS/6000 servers, the most reliable servers with which to build your cluster.

HACMP also supports a full range of networking and disk technologies.  A cluster is likely to contain one of more IP-based connections between the cluster nodes, including Ethernet, Token Ring, FDDI, ATM or the high speed switches, such are used in the RS/6000 SP.  HACMP also supports non-IP based connections such as RS232, Target Mode SCSI or Target Mode SSA.  Both IP and non-IP connections are used to ensure that there is no single point of failure in the networking configuration.  Heartbeat messages are sent between the cluster nodes over all of the network paths defined to HACMP.  These heartbeat flows are used to monitor the 'health' of each node, the networks to which they are attached and the network adapters that are used by them.

As with networking, HACMP supports a wide range of disk technologies and products.  These are typically SCSI, SSA or FC-AL based.  This storage is connected to the cluster nodes and is usually protected by the use of mirroring or RAID techniques.

One of HACMP's key strengths has always been its configuration flexibility, the ability to mix and match any combination of server, disk or networking products to meet a customer's exact requirements rather than forcing a customer to adopt technologies that do not easily fit into their existing environments.  By choosing components that exactly meet a set of customer requirements, combining these in a configuration that has no single point of failure and using HACMP to manage them, the best possible hardware and software can be allowed to do what they were bought to do, to run the business.

Statistically, site failures and disasters are only slightly less likely than hardware failures, when it comes to causes of computer outage[4].



Building outage
4.0%

Corrupt data
2.0%

Storm
4.0%

Fire
9.0%

Power Outage
28.0%

Earthquake
11.0%

Hardware Failure
15.0%

Water Damage
27.0%

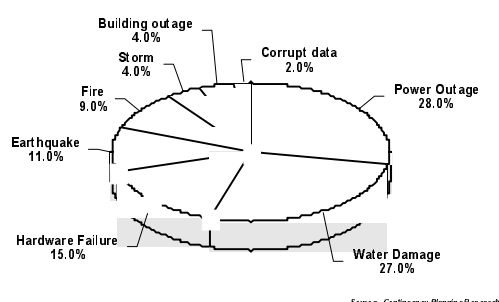Source:  Contingency Planning Research

*Figure 6:  Outage caused by site failures and disasters*

When deciding upon a high availability strategy, it is important to ensure that your choice of solution is also capable of being extended to support disaster recovery capabilities.  You do not want to add additional (and potentially incompatible) products to a working highly available environment, rather you want a seamless integration between your disaster recovery and availability products.

To support disaster recovery within the AIX environment, IBM has two main offerings today which are part of the HACMP family of products; the Geographic Remote Mirror (GeoRM) and the High Availability Geographic Cluster (HAGEO).  GeoRM provides the capability to mirror disk data between two RS/6000 servers connected by an IP network.  It supports a many-to-one model allowing multiple RS/6000 servers to be backed up on a single central server.  Mirroring may be synchronous to ensure maximum data availability or asynchronous for maximum performance at the cost of not having all your data available after a disaster.  Again, customer choice is key, mirroring being performed at a logical volume level.  GeoRM is particularly suitable for keeping a copy of your critical data off-site.  If a failure occurs, manual intervention is required to get access to the data on the remote site.
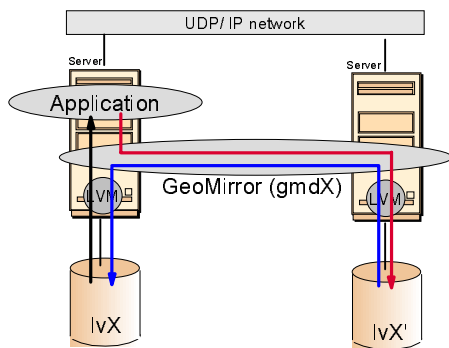
14

Figure 7:  A typical GeoRM configuration

For those customers who want automatic failure detection and failover, the HAGEO products is a better choice.  HAGEO is an extension to HACMP which builds upon the GeoRM disk mirroring functionality and, in addition, provides the necessary 'hooks' into HACMP to provide automatic failure detection and recovery and to support the automatic reintegration of a failed site.
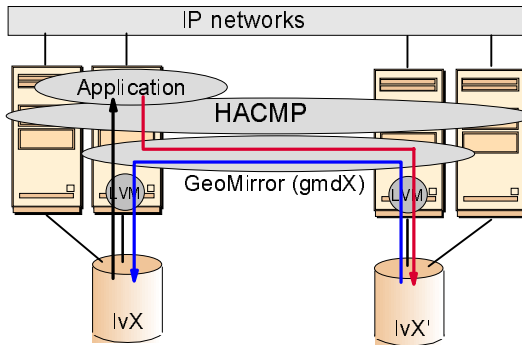


Figure 8:  A typical HAGEO cluster

Clustering for performance

Just as HACMP can be used to provide increased availability for a cluster, it can also be used to provide increased scalability and performance.  The Concurrent Resource Manager component, the Cluster Multi-Processing bit, of HACMP allows up to eight nodes to run the same application workload against a parallel shared disk database in a 'Concurrent Access' configuration.  This takes advantage of the AIX Concurrent Logical Volume Manager which allows an AIX Volume Group to be concurrently accessed by multiple servers.  In an environment like this, it is important to provide a mechanism to coordinate access to the shared disk resources.  HACMP provides it's own lock manager to do this, or alternatively a customer may choose to use a lock manager provided by an application vendor.  Whilst providing scalability in this fashion, a customer doesn't have to compromise on availability either.  This configuration also provides extremely fast failover in the event of a problem arising.  HACMP can also be used with shared-nothing parallel databases in configurations with up to 32 RS/6000 servers in a single HACMP cluster to provide still greater scalability for those applications which can exploit this.

When it comes to absolute scalability however, the prize must go to the RS/6000 SP, the system used as the basis of many of the most powerful computers in the world.  Scalability here can easily be an order of magnitude or more over a typical HACMP cluster,  SP systems consisting of hundreds of nodes being fairly commonplace.  The combination of the SP hardware and it's cluster management software, the Parallel System Support Program (PSSP), give it unrivaled cluster capabilities for both commercial and scientific customers.

Whereas the cluster nodes themselves are little different from their RS/6000 cousins, offering many of the same facilities, albeit in a slightly modified package, and indeed RS/6000 S80 servers themselves may be
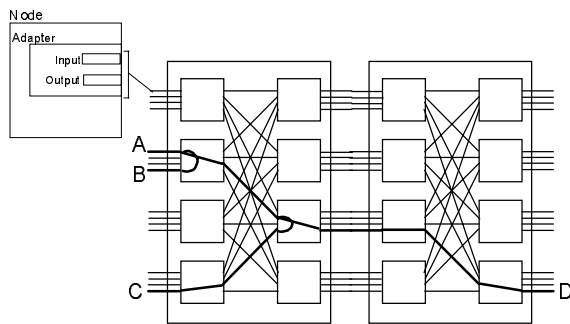
15

*Figure 9:  Example high speed switch interconnect*

included within an RS/6000 SP complex, the main hardware difference between the SP and an RS/6000 cluster comes from the high-speed switches used as cluster interconnects.

These high-speed switches are a characteristic of many SP systems. The switch provides the capability to exchange information between cluster nodes at a very high data rate, with a very low message latency. Each node connects to a switch through one or two adapters physically contained within the cluster node. Each adapter in turn provides one or two ports to the switch.  The number of adapters, adapter ports, switches, and the number of connections between switches are dependent on the number of nodes in the complex, the expected fabric utilization, the bandwidth needed, and on the amount of redundancy desired.  The latest SP switch technology adapters can sustain up to approximately 1000MB/s of communication bandwidth. The switch architecture is based on message passing.  PSSP supports a multithreaded, standards-compliant Message Passing Interface (MPI) via the IBM Parallel Environment for AIX, as well as maintaining its single-threaded MPI support. In addition, PSSP includes a Low-level Application Programming Interface for the SP Switch as well as providing an IP interface that makes the switch look like any other IP network.

Whilst vertical scalability, the ability to increase the processing capabilities of a single cluster node is important, when it comes to the largest of requirements, horizontal scalability, the ability to add more and more cluster nodes is more important.  Whilst many competitors claim to be able to cluster many systems together, the real key to scalability is to make efficient use of these extra nodes.  This is whether the SP's software really comes to the fore.

The Virtual Shared Disk (VSD) is a piece of software which provides the ability for multiple cluster nodes to access raw logical volumes as if the disk were attached locally to the node.  This allows scalability far beyond that which is possible using the Concurrent Logical Volume Manager allowing truly massive databases to be built, yet not compromising on availability.  The Recoverable VSD component provides transparent recovery of VSDs from failures.
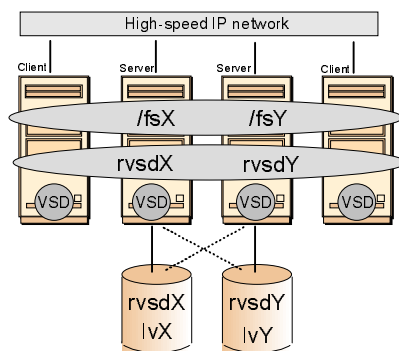
For those customer environments which need access to a cluster file system instead, the SP provides this to via it's General Purpose Parallel File System (GPFS).  GPFS is built on top of the RVSD technology to provide simultaneous access to file systems from up to 128 cluster nodes.



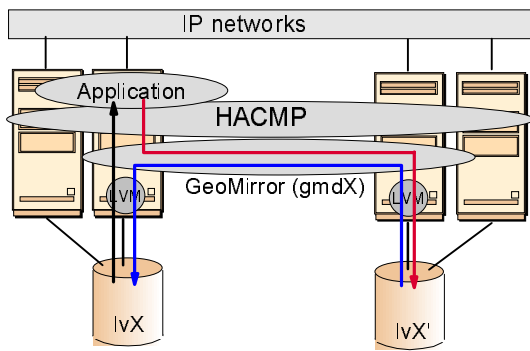*Figure 10:  Recoverable Virtual Shared Disk*

Figure 11: General Purpose Parallel File System

Outside of the commercial arena, the RS/6000 SP also runs a wide range of IBM and industry-standard parallel environments as well as providing parallelized versions of compilers and libraries such as the Engineering and Scientific Subroutine Library (ESSL) and Optimization Subroutine Library (OSL) allowing a cluster of nodes to be used as a single compute-intensive resource. The LoadLeveler product provides a powerful mechanism for managing workloads. Jobs can be submitted to a central pool for LoadLeveler to farm out to the nodes in the cluster based upon their characteristics or special requirements. LoadLeveler supports both serial and parallel jobs.

The collection of performance data is one of the key requirements of any performance environment. The Performance Toolbox Parallel Extensions (PTPE) function of PSSP collects and provides performance data for SP hardware and software through enhancements to the AIX Performance Toolbox product. This not only allows performance monitoring of unique SP subsystems, such as VSDs and the SP Switch, but also organizes the SP into a set of performance reporting groups with coordinating managers, and distributes the burden of monitoring nodes throughout the SP system. This eliminates the need for dedicated performance monitoring nodes. PTPE also provides average performance statistics for the SP system, rather than monitoring every data point on every node. This helps reduce the computational effort required for run-time monitoring of SP performance.

Clustering for ease of management

The PSSP software is not only a set of tools to improve the scalability of the SP system, it also provides a large suite of tools to make the management of a large cluster simpler.

The heart of the system is a single point of management; a Control Workstation from which the system administrator or operator can perform all local and remote administrative functions. This workstation provides an object-oriented graphical interface called Perspectives allowing simpler monitoring and control of hardware. This capability includes powering on and off nodes, changing settings in the ROS of a node, monitoring the state of the LEDs or LCDs of the nodes, and providing remote access to the service processor. In addition, Perspectives is also used for creating and monitoring events, managing storage and



SP Control Workstation

configuring systems. This is a single consolidated graphical user interface providing a common launch pad for PSSP system management applications through direct manipulation of system objects represented as icons. The interface is tightly integrated with the problem management infrastructure allowing users to easily take actions as required. It is also highly scalable for large systems, and can be easily customized to accommodate varying environments. The control workstation also provides a single point of installation and maintenance for software on all nodes in an SP system. Availability of the Control Workstation may be enhanced through the use of the High Availability Control WorkStation (HACWS), which itself is built upon HACMP.
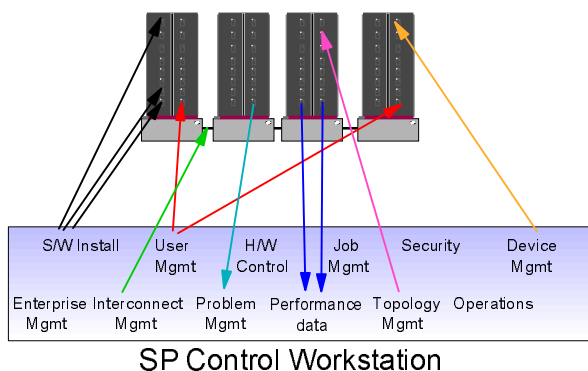
Figure 12: The Control Workstation concept

PSSP provides a full suite of system management applications for the RS/6000 SP and SP-attached RS/6000 servers in addition to the many tools provided as part of AIX. As well as allowing administrative and operational support of the system, many additional unique functions have been added to manage and fully exploit the capabilities of the SP to the utmost degree. These tools enable system administrators and operators to better manage SP systems and their computing environment.

One of the key parts of PSSP is Reliable Scalable Cluster Technology (RSCT) which provides a consistent, flexible framework for managing, controlling and monitoring cluster resources. This framework addresses the requirements for cluster availability, scalability, and manageability.
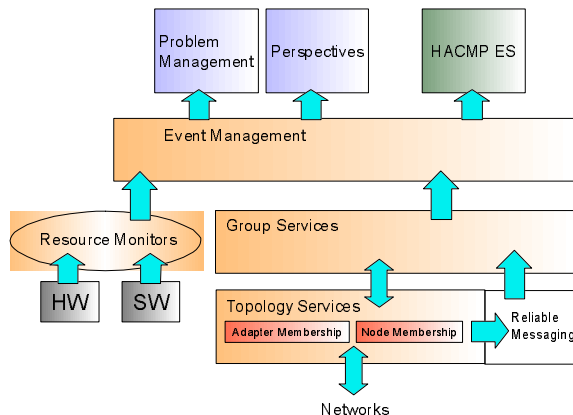


Figure 13: Reliable Scalable Cluster Technology components

The first component of RSCT is Topology Services. This provides a scalable heartbeat algorithm and is used to determine the available communications paths between cluster nodes and to determine which nodes are not accessible at all. From this, adapter membership and node membership of the cluster may be determined. Topology Services provides this information to other cluster subsystem components such as Group Services and it's own Reliable Messaging subsystem.

Group Services provides coordination services for use by subsystems that must precisely control actions across multiple processes executing on more than one node of the cluster. In addition to coordination services, Group Services also provides distributed messaging and synchronization services. These services are typically used to coordinate, not provide the recovery services of those subsystems that use it such as HACMP ES or Event Management. Other applications can use Group Services facilities through it's Application Programming Interfaces (APIs).

Event Management is the third RSCT subsystem which implements services for monitoring resources. Event Management is used for notification and reaction to important events such as failures and resource deprivation. This provides a set of advanced error detection and recovery features that allow proactive rather than reactive management of events to still further reduce the impact and occurrence of unplanned outages.

The SP cluster may be subdivided into different partitions if required. Partitioning the system in this way allows different levels hardware and software to be concurrently executing within the system. This allows different production environments to be concurrently executed for workload isolation or for test environments to coexist, yet not interfere with production workloads. Coexistence of up to three releases of PSSP software within an SP partition allowing for easier software migration.

Sets of software tools and related utilities have been grouped together to offer easier administration. These tools include the consolidation of error and status logs, accounting and performance information for expedited problem determination. When problems do occur, or simply to speed up maintenance activities, parallel system management tools and commands for enabling concurrent parallel performance of system management functions across multiple SP nodes. Keeping files synchronized between cluster nodes is always an issue. PSSP provides file collections for managing duplicated files and directories on multiple nodes. Login control and authentication based upon Kerberos for blocking unauthorized user or group access to a specific SP node or a set of nodes is, of course, an integrated service.

The full breadth of system management features offered by PSSP is unrivaled in the industry in terms of it's capabilities, yet these are not the only offerings from IBM.

HACMP also provides a Cluster Single Point Of Control (CSPOC) component. Just as the Control Workstation provides the ability to manage an SP system from a single point, CSPOC does the same for HACMP, the main difference being the ability to run CSPOC from any cluster node rather than requiring a dedicated workstation. CSPOC is also less extensive in it's capabilities than PSSP, it's main function being the management of shared disk devices, users and cluster services. However, for the smaller cluster, CSPOC provides they key features needed to simplify cluster management.

Whether you are managing SP systems using PSSP or HACMP clusters using CSPOC, IBM provides a full suite of cluster management software to get the best from your clusters.


**Making availability simpler: packaged solutions and ClusterProven™**

When purchasing a clustered solution, you need to be sure that all of the components will work together. Determining this for the hardware, the operating system and the cluster software is relatively straightforward. Selecting a set of components from one vendor or from multiple vendors who have tested their products together is a good start, but you still run the risk, unless you are getting someone with more experience to do this for you, of forgetting one or more components when designing the cluster. A better solution for you in such as case, is a packaged cluster solution that ensures you have all the necessary parts 'in the box'.

IBM offers packaged cluster solutions based upon most of it's RS/6000 servers. These contain all the necessary components to required to build a fully functional, and very highly available cluster, the IBM HA-H70 cluster solution, for example, has an estimated availability for the hardware, AIX and HACMP of 99.999%! Unlike some other packaged offerings which force you to take a particular configuration that is unlikely to meet your needs, an IBM packaged cluster solution can be modified to exactly meet your requirements. Adding processors, disks and memory is only the start, the cluster also comes with your choice of networking hardware to ensure that it may be seamlessly integrated into your existing production environments with the minimum of effort.

Getting the hardware right is just the start though. You also need to select software which will work in this environment. Given that software failures are a greater cause of outage than hardware, it is key that the software you choose not only works, but has been proven to work. This is where the IBM ClusterProven programme comes in.

The ClusterProven programme provide customers with clear criteria for selecting highly available solutions on IBM server platforms. IBM has defined standards that are applicable to all IBM servers and are formulated in terms of their value to end-users. An application solution receives ClusterProven status if it has been tested to perform in a clustered server environment to meet predefined criteria, delivering availability and scalability characteristics beyond that achieved on a single system, and maintaining application availability in the event of failure. There are two levels of ClusterProven applications: ClusterProven, which applies to the basic exploitation of the high availability features of IBM servers, and Advanced ClusterProven, which applies to application implementations that move customers even closer to true continuous operations, leveraging clustering technology to further improve availability, performance, scalability and manageability.

To assure quality and integrity in the use of the ClusterProven trademarks, IBM has set up validation and registration procedures. Customers purchasing ClusterProven applications do so with the knowledge that the application they have chosen has been proven to work correctly in a clustered environment to the high standards you would expect from IBM.

In addition to the benefits that the programme offers customers, it also provides support to solution developers including seminars, training, documentation and high availability centers of competence, designed to provide help in testing and exploitation of IBM clustering technologies.

Choosing a packaged cluster solution ensures that you get the hardware you need to both meet your requirements and function correctly as a highly available cluster. Adding a ClusterProven application, one which has been proven to work in this environment, completes the system. What could be simpler?

**Unrivaled support**

All of IBM's cluster offerings, as with any IBM product, are backed by IBM worldwide service and support. IBM's commitment is behind every product we sell, helping ensure the highest possible customer satisfaction. Unlike many vendors who claim to offer worldwide support, only IBM is truly able to deliver this. With over nine million systems under contract with IBM, nearly one million of these being non-IBM, we have an unsurpassed breadth of skills in professionals who undergo continual training to keep their skills at the leading edge. With more than 116,000 services professionals located in 164 countries, IBM is truly positioned to deliver the support that customers need.

In addition to the wide range of services suitable for all environments, IBM recognizes the special requirements of highly available systems and so, also offers over thirty standard service offerings specifically designed to support high availability systems for customers whose business objectives critically depend on system availability. IBM's High Availability Services are uniquely positioned in the marketplace in this respect.

Support is delivered through a wide range of services to support both hardware and software products at every stage of the product life cycle.

Before the cluster goes live, services for planning, designing and installing clusters are available. These services might also entail the migration of existing environments that are being replaced by the new cluster environment or the integration of the cluster into an existing IT infrastructure. The SmoothStart offerings allow customers to get standard cluster configurations installed and operational in a timely and efficient fashion. Customers requiring design and implementation of more advanced or custom solutions containing both IBM and other leading vendors' hardware and software. These are designed to meet your target availability level and are developed for your unique business requirements and systems environment. All of these services are built from various combinations of nine core service elements, all designed to enable customers to tailor solutions that specifically meet their business requirements.

Once the cluster is installed and running, IBM offers a wide range of hardware and software maintenance and support services to ensure that your high availability systems are always available and delivering real business benefit. These services include in-depth availability assessments of customer environments from end-to-end, not just assessing the individual IT components. Such an assessment produces recommendations for change that will improve the reliability and availability of the environment still further.

For those customers who do not wish to run and manage their own clusters, IBM offers a range of operational support and facilities management services to run clusters, and indeed, entire IT environments.

Finally, IBM's services include not just the availability of your IT equipment, but also your buildings and infrastructure. IBM, through it's Business Recovery Services has an unrivaled wealth of experience in helping customers build, manage and run their disaster recovery plans. Whether this involves providing emergency equipment should disaster strike through to hosting and / or managing your disaster recovery systems themselves, IBM provides outstanding emergency recovery facilities throughout the world plus world-class hardware and software testing laboratories.

Employee skill, or lack of it, is often a major cause of computer outage. Serious mistakes made by system administrators can often result in systems being unavailable for hours or even days. To ensure that your employees have the skills that they need to manage and run your high availability systems, IBM provides a wide range of education courses specifically pertaining to High Availability systems. These include courses for the designers, installers, and systems administrators of HACMP clusters as well as specific courses for RS/6000 SP systems and HAGEO.
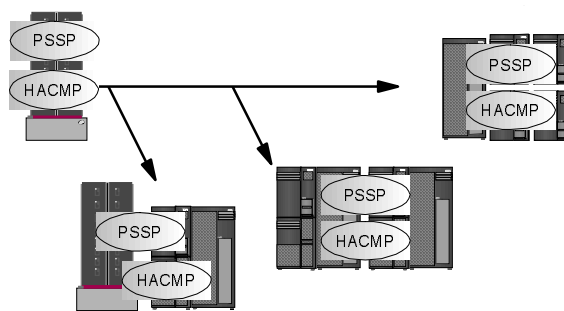
## Part III.  Future directions

Having industry-leading products and offerings within the UNIX clustering arena, puts IBM in an incredibly strong position when it comes to defining the future of it's AIX clustering offerings.  In effect, we can pick and choose between the best technologies available to deliver the best for our customers.  In 1998, IBM stated its intention to combine PSSP and HACMP into a single cluster software offering.  This brings the proven and mature availability benefits of HACMP and the scalability and manageability benefits of the RS/6000 SP and PSSP to all RS/6000 servers, dramatically reducing overall system complexity.  In short, IBM intends to provide a single comprehensive cluster environment for all AIX systems providing improved system manageability, availability, scalability and ease of use.

In order to preserve customer investment in their existing platforms and to utilize new, but proven, technologies as they become available, this merging of the two software offerings is being slowly phased in. Where the key is providing stable and reliable systems, frequent or dramatic change is most unwelcome. The intent of the remainder of this document is to provide a view of what has been achieved to date and to explain how the remainder of IBM's overall AIX cluster strategy will be delivered.

### Extending the scope of PSSP and the RS/6000 SP

PSSP originally ran only on the RS/6000 SP.  As customers demanded increasing processing power within a single SP node, the capability to use RS/6000 S80 processors as SP nodes was added, thus extending PSSP outside of the SP frame for the first time.  This capability to attach standard RS/6000 servers to an SP, and indeed to be able to build SP-like systems (clusters of RS/6000 servers running PSSP) without the requirement for SP frames will be further enhanced in the future with the addition of further servers capable of being either attached to the SP system or of running PSSP outside of the SP.  This brings the scalable systems management facilities of the SP to clusters of RS/6000 servers.

Longer term, clustered systems composed of standard RS/6000 building blocks will be built. These will have the similar characteristics as the RS/6000 SP system today.  In effect, what you w have here is a unified system consisting of multiple hardware nodes running multiple operating systems. These nodes will have common physical attributes to allow complete mixing and matching to fully meet customer requirements and protect investment.  Based upon a common rack-mounted format, these modular 'building blocks' will utilize technology common to users of RS/6000 servers today such as the highly dynamic and flexible configurability and concurrent maintenance.



Figure 14:  Extending PSSP outside of the SP frame

Just as an SP system today is highly configurable, this next generation clustered system will be likewise. Individual nodes within the system may run as standalone entities, multiple server workloads being consolidated into a single or multiple rack package.  Collections of nodes may also be logically separate from others via hardware and software partitions.  These maybe being used, as SP partitions are today, to run different software or application environments.  Individual nodes may be loosely coupled together in an HACMP-like fashion or more tightly coupled to create NUMA systems.

The system will also provide unified system network(s) capable of supporting NUMA, more loosely coupled clusters and I/O.  These networks will be based upon a future version of the SP switch technology available today.  Given than most clusters require some sort of shared storage, and that the future of storage is Storage Area Network (SAN) based, SAN management capabilities will also be fully integrated with

systems configurability in the same modular fashion.  Systems with higher I/O than compute requirements simply replacing compute building blocks with storage building blocks.

Regardless of the roles in which the individual system nodes find themselves, the system will provide a single unified management architecture. Managed from a single point of control, derived from the SP Control Workstation today, a full suite of remote hardware and software operations will be provided.  This will allow hardware control and service alongside the consolidated and distributed administration and operations facilities provided by a PSSP-like operating environment, a dedicated service network providing connections to all cluster nodes.

## Hardware management

Today, the SP Control Workstation (CWS) offers the ability to remotely manage and monitor SP hardware.  This capability includes powering on and off a node, changing settings in the ROS of a node, monitoring the state of the LEDs or LCDs of the nodes, and providing remote access to the service processor.  This capability will be extended and made available across all RS/6000 servers from a single Hardware Management Console (HMC). There is also evolving a need for a hardware management station for control of NUMA and logically partitioned (LPAR) SMP systems. From an HMC, customers can control any RS/6000 server; including SP nodes, frames and switches, LPAR SMP servers; and NUMA systems. This control point may also eventually control other IBM servers. This single point of management is key to enhancing the manageability of clusters to make them appear, and be managed, like single systems.

## High performance cluster interconnects

Through a staged implementation, the high performance interconnect fabric will be decoupled from SP system boundaries so that we are able to provide connectivity to any RS/6000 system.  Logical extension of this will allow NUMA systems to be built from RS/6000 components, the differentiation between the two being solely whether the NUMA capabilities are present and have been enabled, the underlying technology being identical.  The customer will have the choice of whether to utilize NUMA or not, based upon the applications they intend to run and their ability to work in a NUMA environment.  Centralized switch hardware management will be provided from the HMC, again working to make management simpler.  Nodes should be able to be connected to the switch, discovered, and configured without requiring explicit administrator action or requiring a secondary LAN.

This system will bring together and builds upon the best features of the current hardware and software offerings.  From the SP side of the house come the packaging density, the ability to get a lot of computing power into a small footprint and the scalability that this brings.  Also, the single point of management, the high performance interconnect technology  and the flexibility of the system to act equally well in both commercial or technical / parallel environments.  RS/6000 brings it's superior price performance and an extremely stable platform architecture.  This stability, coupled with the scalability it offers, makes it an ideal choice for the basic building blocks, as it does today.  The NUMA technology from NUMA-Q® and their experience in delivering enterprise-class systems based on commodity technologies is their contribution.
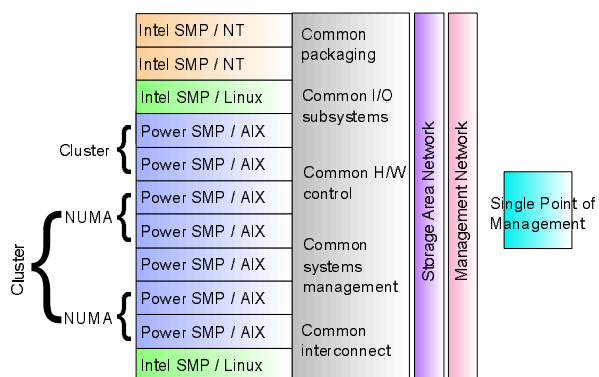


*Figure 15:  A future clustered server system*

22

What hasn't been mentioned here so far, more because it is really outside the scope of the document, is the fact that in addition to RS/6000 building blocks, systems based upon Intel processors will also be able to participate in this server complex. This provides any customer with a choice of hardware (Intel® and PowerPC) and Operating System (AIX/Monterey, Linux, NT). This allows a customer to choose those components that exactly meet their application requirements whilst allowing the management and running of mixed environments in single system. It also permits system evolution over time protecting investment in infrastructure components such as disks and networks and allowing flexible allocation of computing resources as workloads change.
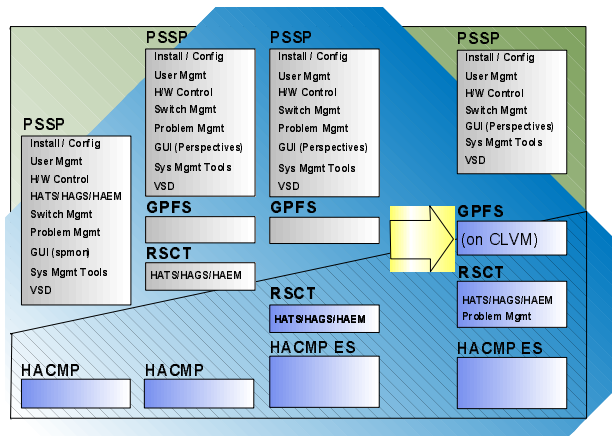
From the software side, single systems and highly scalable system management and the availability features of PSSP come from the SP. This coupled with the proven strengths of AIX, the application portfolio that it supports and the integrated SAN management strengths of NUMA-Q are also key to the system.



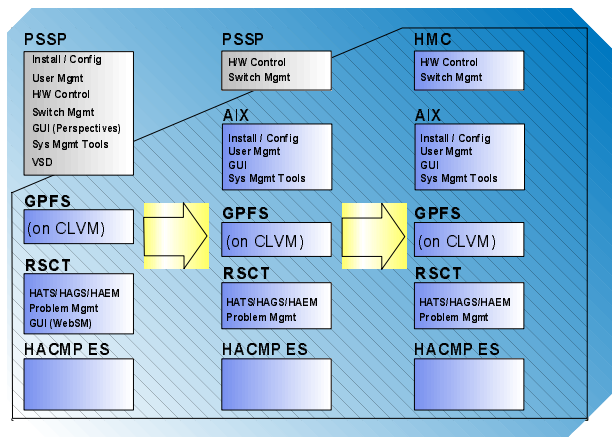*Figure 16: Software evolution - Phase 1*



*Figure 17: Software evolution - Phase 2*

## Evolving HACMP

Just as PSSP is being enhanced to allow SP-like systems to be built from standard RS/6000 components, so other products are being enhanced with parts of PSSP. The original HACMP product, so called HACMP 'Classic', was never intended to support clusters of more than eight RS/6000 systems. As customer requirements in this respect grew, HACMP needed a more scalable underlying infrastructure to address these customer requirements. This infrastructure is available in the components which make up PSSP's Reliable Scalable Cluster Technology (RSCT) subsystem. Replacing the original HACMP infrastructure with that of RSCT, whilst still maintaining a consistent user interface to preserve customer investment in systems and skills, gave rise to HACMP Enhanced Scalability (HACMP ES). This, more robust and scalable version of HACMP, allows much more advanced monitoring and recovery actions over and above

those of it's earlier 'Classic' versions.  As well as offering these additional advanced features, HACMP ES also provides the functionality of the 'Classic' product[7] for compatibility reasons.
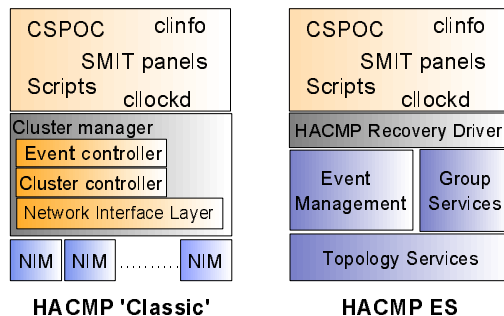


*Figure 18:  HACMP 'Classic' vs. HACMP ES*

For many years, it has been possible to migrate a cluster from one version of HACMP 'Classic' to another or from one version of HACMP ES to another, a node at a time.  This allows a single cluster node to be taken out of the running environment, upgraded to the latest level of AIX and HACMP etc. before rejoining the cluster.  The other nodes are then upgraded in turn.  This is only possible because HACMP, unlike many competitive offerings, is capable of running in a 'mixed' configuration whilst this migration is occurring.  This node-by-node migration allows the services provided by the cluster to be available to the end users throughout the migration procedure, a key requirement when running a cluster in a true 7x24 environment.  From HACMP Version 4.3.1, it has also been possible to migrate HACMP 'Classic' to HACMP ES in a node-by-node fashion, even though they have a fundamentally different underlying infrastructure.

Many of these advanced capabilities are a direct result of the 'built-in' RSCT components and consequently have not been able to be added to the original HACMP 'Classic' product.  In addition, many of the future planned enhancements to HACMP will build directly upon this RSCT functionality.  These include the ability to fail over to a cluster node meeting a defined set of criteria, for example, the node that is least heavily loaded.  RSCT provides the resource monitoring to collect this information to allow HACMP to make a decision as to which node to failover to.

Over time, whilst IBM will continue to enhance both the 'Classic' and ES variants of HACMP, many functions will only be available in the ES version.  As HACMP 'classic' will, over time, lag further and further behind HACMP ES in terms of functionality, we foresee a day when customers will always choose HACMP ES.  Until that time comes, IBM will enhance the 'Classic' variant only when the functionality may be added with code common to both versions.  HACMP ES, is and will be, the base for all future HACMP releases.

In addition to HACMP, IBM used to offer a product called High Availability for NFS (HANFS).  This provided a set of facilities to build a highly available NFS server based upon a cluster of two RS/6000 systems.  Originally, it and HACMP were separate products.  Over time, these were merged, such that HANFS used HACMP 'Classic' as its underlying technology.  When this occurred, HANFS became a feature of HACMP.  As more and more customers used HACMP to provide highly available NFS services to their end-users, it became apparent that maintaining HANFS as a separate feature did not make sense.  Driven by our customer's requirements the HANFS functionality has been fully merged into HACMP as of HACMP Version 4.4.  The HANFS feature has consequently been withdrawn  whilst maintaining customer investments throughout.  As with other HACMP migrations, it is possible for an existing HANFS customer to migrate their clusters using the node-by-node migration feature, preserving service to the end users throughout this process.

Future HACMP directions

One of the key areas for future HACMP enhancement is in the area of usability and configurability. Disk and networking resource relationships will be discovered dynamically from the base AIX configuration without the need for explicitly being configured to the cluster.  In addition, HACMP ES will simplify the handling of IP addresses by taking advantage of IP aliasing. Each LAN adapter will have a fixed base or "base" address which will never be replaced.  This address is guaranteed to always map to the node where the adapter is installed. The IP address associated with an application will be established by an alias.  During recovery, the alias will be reassigned to a new LAN adapter, and a gratuitous ARP will be issued to
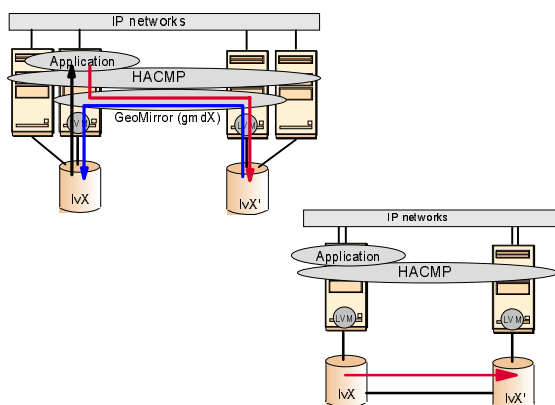
resolve the ARP mapping for that address. Multiple independent resource groups can take advantage of the same adapter by using different IP alias addresses. MAC address takeover will also be offered for compatibility with existing clusters.

Disaster recovery

To support disaster recovery within the AIX environment, IBM has two main product offerings today; the Geographic Remote Mirror (GeoRM) and the High Availability Geographic Cluster (HAGEO). GeoRM provides the capability to mirror disk data between two RS/6000 servers connected by an IP network. HAGEO is an add-on to HACMP which builds upon the GeoRM disk mirroring functionality and, in addition, provides the necessary 'hooks' into HACMP to provide automatic failure detection and recovery and to support the automatic reintegration of a failed site.

Despite working together, the integration between HAGEO and HACMP can be made much closer still. Answering our customer's concerns in these areas, it is our intent to more tightly integrate HACMP ES and HAGEO. This will ensure that these two products provide a single consistent interface for a clustered system, whether the cluster nodes are physically located together or separated by geographic distances. The rationale behind this is that alternative technologies such as SSA optical extenders or the Peer to Peer Remote Copy (PPRC) functionality to be offered by the IBM Enterprise Storage Server (ESS) also provide the ability to mirror disk data across longer distances between buildings or sites. Customers are already, and will increasingly wish to combine these hardware facilities within an HACMP cluster. Clusters supporting these hardware-mirroring environments have the same characteristics as systems running HAGEO. That is, they require the capability to determine the difference between failures of nodes, networks, network adapters and sites.



*Figure 19: Comparing HAGEO and HACMP with remote mirroring solutions*

HAGEO provides these capabilities today but HACMP also needs them in order to support split-site configurations. It thus makes sense to add these HAGEO functions (DBFS support, the concept of site (locality) and the site failover logic) to a future release of the base HACMP product. That way, HACMP can support either split-site environment. GeoRM will remain a standalone geographic disk mirroring product for those customers requiring it.

A customer might decide to install GeoRM today as a mechanism to 'bunker' their data to a remote location. If a site failure occurs, manual intervention will be required to make the bunkered data available for use. As their availability requirements increase in the future, they are likely to need faster access to the data on the remote site. In response to these requirements, they might decide to install HACMP to automate failure detection and recovery processing. The installation of this future HACMP release onto a system already running GeoRM will configure itself to provide similar functionality as is offered by HAGEO today. A similar process will occur when a customer installs this version of HACMP onto a system which is running with SSA extenders or ESS-PPRC.

**Usability and manageability enhancements**

The usability of clustered systems today is more complex than some customers would like. To that extent, major steps are being taken to provide a suite of tools and utilities that will dramatically improve the usability

of clustered systems, both in terms of configuration and ongoing management. The goal is to provide a system that requires little more effort to manage than that required for a standalone RS/6000 system. This will take time to achieve, but we're well on our way to getting there.

HACMP already offers a considerable suite of cluster management facilities via it's Cluster Single Point Of Control (C-SPOC) capabilities. These will be further enhanced and extended in the future by the provision of additional tools. Some of these will be derived from PSSP or RSCT functions, others will be new tools. All of them, however, will be integrated to provide consistent interfaces for clustered systems management. The end result is a new PSSP-like set of facilities applicable equally to SP-like as well as standalone RS/6000 systems.

Administrators will be able to manage all aspects of the entire cluster from a single point, the complete suite of AIX single system management applications also being available through this framework for use with individual nodes of the cluster. This single point, much as an CWS does today, providing a collection point for all monitoring data and control activities. For those users wanting graphical access to the system, a single consistent GUI will be provided. For those users who don't like GUIs, the cluster system management environment will also support command line access which is important for the automation of cluster management through scripting, and access for management from low bandwidth connections.

These cluster management facilities will encompass the creation of clusters from any combination of RS/6000 components, the identification of cluster resources and their actual configuration being dynamically discovered. Other configuration and management utilities to simplify the tasks of configuring and managing clusters to the level of a single standalone system will be provided along with tools to provide increased control and flexibility over the management of resource failures and the policies to react to those failures.

Just as PSSP does today, parallel installation and management will be possible for all managed servers. As these may be Intel or PowerPC based systems running a variety of operating systems, providing a single consistent interface will dramatically ease management of complex heterogeneous environments. Whereas today an administrator might be required to be familiar with many tools and interfaces, each with a different look-and-feel, in the future, tasks such as hardware control, resource monitoring, problem and topology management will be accessible via a single interface for all managed servers.

### Single systems image

There is much debate within the industry at present about clusters which implement a Single Systems Image (SSI). We have talked with many customers about their real requirements and consequently have determined that a Single System Image is unlikely to be the best solution to the vast majority of customer problems. A true single system image implies many things including:

- Consistent user access
- A consistent view of data
- A single device space
- A single memory space
- A single process space
- A single point of administration

Implementing a true SSI requires massive changes to be made to the underlying operating system. This is likely to make it almost unrecognizable to a user familiar with AIX today and results in an environment which is likely to be less reliable than a standalone system and will require significantly more skill to implement and maintain than comparable cluster offerings. As a result, IBM intends to offer a clustered environment that provides a single system image for those things that our customers tell us are most important to them, we originally called this Global System Image (GSI) an initiative based upon globalizing only those components that it is correct to do so.

IBM already offers it's Concurrent Logical Volume Manager (CLVM) and Virtual Shared Disk (VSD) products providing a single system image for raw logical volumes. The RS/6000 SP also offers a single system image for file systems through the use of the Generalized Parallel File System (GPFS). This is offered as a cluster file system for application and subsystem use and is designed for scalability and availability and provides concurrent access across the cluster with full coherence. Files can be striped across multiple disks for high performance, parallel access, and load balancing. Access to files is mediated through an efficient distributed Token Manager. GPFS provides many file system RAS features, including built in recovery disk replication, meta-data journaling, file replication for availability, and on-line reconfiguration for adding, removing or replacing disks, or restriping the file system. Although GPFS today runs only on the RS/6000 SP, we will provide a GPFS cluster file system that uses the Concurrent Logical Volume Manager for generic RS/6000 servers in the near future.

Further enhancement of this is likely to utilize enhancements to components in the base AIX operating system to provide a common look and feel and systems image for cluster management making management of a cluster, no more complex than the management of a standalone RS/6000 server.
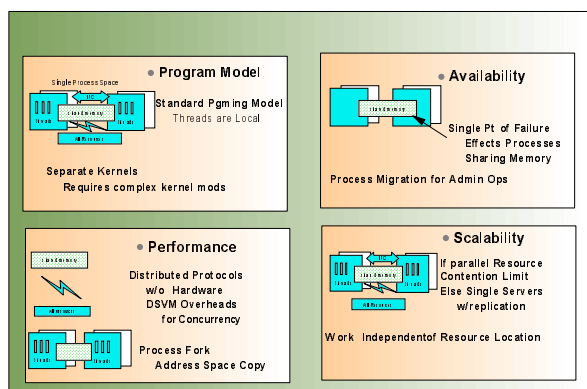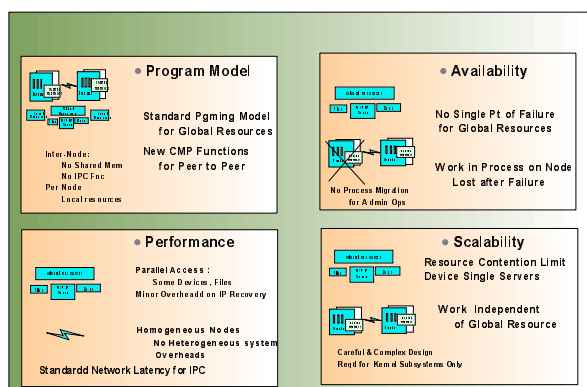


Figure 20: Single System Image environment



**Clustered solutions**

Figure 21: Global System Image environment

People buy clusters to perform useful work. IBM's existing packaged cluster solutions (HA50, HA70 etc) have been very well received by our customers. Consequently, IBM will offer packaged highly available cluster solutions based on all of it's future RS/6000 server systems. Hardware is just part of the story though. Applications are the real key. IBM has announced it's ClusterProven program to assist ISVs in enabling their applications to work with our cluster offerings with some success. One of the criticisms of this program however, is that the resultant scripts are often difficult to obtain so steps are being taken to address this by adding application support packages as part of the HACMP ES product. These packages will include the necessary scripts to start and stop the

application, scripts to automatically configure the cluster to support the application and code to perform application monitoring.  These packages will be built within a standard framework and will be supported by IBM.


## Project Monterey

Project Monterey is the name given to the version of AIX which is being ported to Intel's IA-64 architecture.  Components of SCO UNIX and NUMA-QPTX are also being added to make Monterey the leading IA-64 UNIX offering.  Applications written to the Monterey APIs will be source code compatible across both POWER and IA-64 platforms.  IBM will port HACMP ES to the IA-64 Monterey platform ensuring customers who need high availability functions for their IA-64 servers need look no further than the industry's leading high availability offering.  HACMP ES will be available on Monterey soon after general availability with HAGEO and GeoRM following after.


## NUMA-Q clusters

In 1999, IBM purchasedSequent, whose hardware platforms became the IBM NUMA-Q range of servers.  Today, these servers run thePTX operating system and provide their own clustering software.  In the future, NUMA-Q servers will run the Monterey operating system and consequently will have also run the Monterey version of HACMP ES.

As was stated earlier, many technologies and competencies from NUMA-Q are also being used in delivering the next generation clustered server systems, the Intel-based building blocks being similar in concept to the NUMA-Q 'quads' used today.  This future system can be further regarded as an expansion of NUMA-Q with RS/6000 and SP technology as well as an RS/6000 SP-like system with Intel nodes.

NUMA systems themselves can be thought of as being tightly coupled clusters of servers with an SSI management environment.  This allows them to appear easier to manage than a conventional cluster.  The application and availability issues however, prevent them from being used for the same sorts of tasks as a general purpose cluster.  Clustering NUMA systems together provides the benefits of both environments.


## Linux and the open source movement

Linux is an open-source operating system designed originally to provide PC users with a free or low cost UNIX operating system.  It has since been ported to most other hardware platforms including the RS/6000 and IBM S/390®.  Linux is open and extensible and provides a remarkably complete operating system offering many, but not all, of the features of commercial-strength UNIX operating systems such as AIX.  One of the key drivers encouraging use of Linux by software developers is it's conformance to open systems standards.  Consequently, software written on a cheap development platform with a good range of development tools is easily portable to other UNIX platforms.

IBM is fully behind customers and their choice of operating environment.  For those customers who are undecided upon a full Linux strategy, IBM will provide interoperability and portability to AIX systems, far in excess of other vendor offerings.  IBM is also porting many of it's keymiddleware offerings to Linux to allow customers to use proven applications in their Linux environments.  Clustering is part of this story.

Customers who wish to combine Linux systems together with their AIX servers will be able to do so.  The common distributed systems management functions discussed earlier that will be available to manage clustered AIX systems will be extended to also be able to manage clusters of Linux servers as well.  This will bring PSSP-like capabilities and tools to Linux systems.

## Software Fault Tolerance:  FT CORBA

High availability implementations are not sufficient to meet availability requirements for certain applications and environments such as air traffic control or emergency medical systems.  For these, fault tolerance is required.  Previous fault tolerant implementations have, to a large extent, been proprietary, have consequently been high-cost and hence only suitable for large systems.  To address fault tolerance requirements applicable to a larger range of systems environments, the Object Management Group (OMG) have published a specification [OMG document 00-01-19] for a software fault tolerant environment.
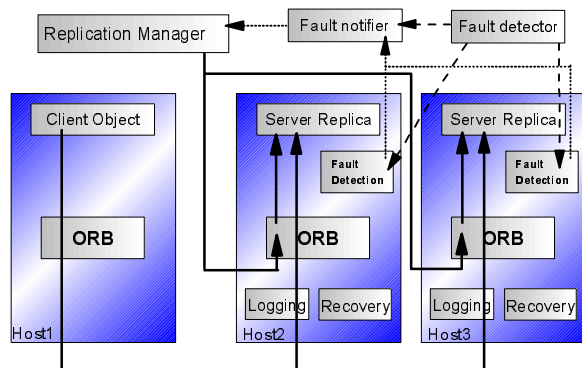


*Figure 22:  An FT CORBA implementation*

The FT-CORBA specification defines an environment based on a Fault-Tolerant Object Request Broker (ORB).  In an environment with no single point of failure, fault tolerance is provided via replicated server objects.  A client will send Requests to each of the active servers in a replicated cluster. Each server will process the Request as though they were the only participant, producing a corresponding Reply that is returned to the client ORB. This will process the first message received for any given Request. If the Reply contains a successful response, it is delivered back to the client program, subsequent Replies from the other servers being discarded.  The FT-CORBA specification supports a variety of fault tolerant strategies including request retry, redirection to an alternative server and both passive and active replication between the server components.

Using the active replication model defined by the specification, both server replicas are kept in passive synchronization by replicating the entire sequence of state transitions in both server processes.  A failure in one process can be recovered by simply switching to the replicated server process whose state is already in sync with the failing server.  Whilst both servers are actively replicated, they are passively synchronized. No attempt is made to ensure that both servers are absolutely consistent with each other. Both servers are presumed to be in sync by virtue of having both received and processed the same set of state transitions, and based on the assumption that the life span of any given element of state is ephemeral.  Any potential to get out of sync, if it does occur, will not persist for longer than this life span. Failover recovery is such an environment should be sub-second.

IBM has an extensive involvement at OMG and in defining CORBA-related standards. IBM has been involved with OMG since it's inception, and has contributed significantly to the CORBA ORB specification, language bindings, Object Services and Component Model, to name a few.  IBM is working with customers to define their requirements for FT-CORBA and is investigating providing a C++ ORB conforming to the CORBA Fault-Tolerant ORB specification using WebSphere™ as a delivery vehicle.  WebSphere already provides exceptional security, reliability, availability and data integrity features.  The addition of FT-CORBA would further enhance this in the future.

## Conclusions

IBM has a leadership set of cluster offerings today running on it's AIX systems, be they RS/6000 SPs or HACMP clusters.  A large number of customers are already finding out just what it is that makes these cluster offering special and how they really help deliver the benefits in availability, scalability, ease of management and investment protection that clustering offers.  Many more customers are joining this select group every day as they realize that only IBM is truly capable of delivering the systems they need to help them succeed.

In common with all of IBM's offerings, our cluster offerings are backed by the most experienced organization in the world when it comes to delivering true business computing.  Whether a software or hardware developer, a system designer building and installing systems for customers or a support professional helping customers with the ongoing management and maintenance of their cluster, the experience and dedication of these professionals in meeting our customer's needs is unmatched within the industry.

This combination of products and experience has made IBM the leading supplier of clustered solutions worldwide.  Furthermore, working with our customers has helped us make these products better still. Listening to our customer's needs, using our wealth of experience and building upon these leadership products, IBM intends to create the most comprehensive set of clustered solutions to date.  The best are just about to get better!

But this is not all about the future.  With the compatibility and seamless migration between versions of software offered by these leadership products, you don't have to wait to take advantage of IBM's cluster offerings.  Just as you can build a small cluster today and grow it as your business needs change, so you will be able to take full advantage of these future capabilities as they become available in both the short and long term.  If you suddenly find that the business needs faster processors, more storage, higher levels of availability or disaster recovery, these capabilities can be added to your existing clusters as required whilst still maintaining the investment in skills and equipment already in operation.  There is no need to wait.  The benefits of an IBM clustering solution that provides availability, scalability, manageability and investment protection are available to you today so you can make the choice, safe in the knowledge that the best is still to come.


## For more information

Product information for today's IBM clustered systems can be found on the World Wide Web at:

 **http://ibm.com/servers/unix**

## Sources and other notes

1. Standish Group study

2. GartnerGroup 10/13/99

3. Contingency Planning Research Inc.

4. GartnerGroup 10/13/99

5. DH. Brown 3/00

6. Meta Group 1999

7. HACMP ES does not provide an exact functional equivalent of HACMP 'Classic'. For example, certain functions required to support hardware devices that were withdrawn from marketing before HACMP ES became available are not present in HACMP ES. To all intents and purposes however, in just about any customer environment, HACMP ES may be regarded as offering all of the function provided by HACMP 'Classic'.

## Special Notices

This document was produced in the United Kingdom of Great Britain and Northern Ireland. IBM may not offer the products, programs, services or features discussed herein in other countries, and the information may be subject to change without notice. Consult your local IBM business contact for information on the products, programs, services, and features available in your area. Any reference to an IBM product, program, service or feature is not intended to state or imply that only IBM's product, program, service or feature may be used. Any functionally equivalent product, program, service or feature that does not infringe on any of IBM's intellectual property rights may be used instead of the IBM product, program, service or feature.

Information in this document concerning non-IBM products was obtained from the suppliers of these products, published announcement material or other publicly available sources. Sources for non-IBM list prices and performance numbers are taken from publicly available information including D.H. Brown, vendor announcements, vendor WWW Home Pages, SPEC Home Page, GPC (Graphics Processing Council) Home Page and TPC (Transaction Processing Performance Council) Home Page. IBM has not tested these products and cannot confirm the accuracy of performance, compatibility or any other claims related to non-IBM products. Questions on the capabilities of non-IBM products should be addressed to the suppliers of those products.

IBM may have patents or pending patent applications covering subject matter in this document. The furnishing of this document does not give you any license to these patents. Send license inquires, in writing, to IBM Director of Licensing, IBM Corporation, New Castle Drive, Armonk, NY 10504-1785 USA.

All statements regarding IBM's future direction and intent are subject to change or withdrawal without notice, and represent goals and objectives only. Contact your local IBM office or IBM authorized reseller for the full text of a specific Statement of General Direction.

The information contained in this document has not been submitted to any formal IBM test and is distributed "AS IS". While each item may have been reviewed by IBM for accuracy in a specific situation, there is no guarantee that the same or similar results will be obtained elsewhere. The use of this information or the implementation of any techniques described herein is a customer responsibility and depends on the customer's ability to evaluate and integrate them into the customer's operational environment. Customers attempting to adapt these techniques to their own environments do so at their own risk.

IBM is not responsible for printing errors in this publication that result in pricing or information inaccuracies.

The information contained in this document represents the current views of IBM on the issues discussed as of the date of publication. IBM cannot guarantee the accuracy of any information presented after the date of publication.

All prices shown are IBM's suggested list prices; dealer prices may vary.

IBM products are manufactured from new parts, or new and serviceable used parts. Regardless, our warranty terms apply.

Information provided in this document and information contained on IBM's past and present Year 2000 Internet Web site pages regarding products and services offered by IBM and its subsidiaries are "Year 2000 Readiness Disclosures" under the Year 2000 Information and Readiness Disclosure Act of 1998, a U.S statute enacted on October 19, 1998. IBM's Year 2000 Internet Web site pages have been and will continue to be our primary mechanism for communicating year 2000 information. Please see the "legal" icon on IBM's Year 2000 Web site (www.ibm.com/year2000) for further information regarding this statute and its applicability to IBM.

Any performance data contained in this document was determined in a controlled environment. Therefore, the results obtained in other operating environments may vary significantly. Some measurements quoted in this document may have been made on development-level systems. There is no guarantee these measurements will be the same on generally-available systems. Some measurements quoted in this document may have been estimated through extrapolation. Actual results may vary. Users of this document should verify the applicable data for their specific environment.

The following terms are registered trademarks of International Business Machines Corporation in the United States and/or other countries: ADSTAR, AIX, AIX/6000, AS/400, CSet++, CICS, CICS/6000, DB2, ESCON, IBM, Information Warehouse, LANStreamer, LoadLeveler, Magstar, MediaStreamer, Micro Channel, MQSeries, Netfinity, Parallel Sysplex, POWERparallel, PowerPC (logo), RS/6000, S/390, Service Director, ThinkPad, TURBOWAYS, Videocharger, VisualAge. The following terms are trademarks of International Business Machines Corporation in the United States and/or other countries: AIX VMe, AS/400e, DB2 Universal Database, Deep Blue, e-business (logo), HACMP/6000, Intelligent Miner, Intellistation, Network Station, POWER2 Architecture, PowerPC 604, PowerPC Architecture, SmoothStart, SP. A full list of U.S. trademarks owned by IBM may be found at www.ibm.com/legal/copy/trade.html.

Microsoft, Windows, Windows NT and the Windows 95 logo are trademarks or registered trademarks of Microsoft Corporation in the United States, other countries or both. UNIX is a registered trademark of The Open Group. Java and all Java-based trademarks and logos are trademarks of Sun Microsystems, Inc. in the United States and other countries. Lotus, Lotus Domino and Lotus Notes are trademarks or registered trademarks of Lotus Development Corporation. Tivoli, TME, TME 10 and TME 10 Global Enterprise Manager are trademarks or registered trademarks of Tivoli Systems, Inc. Other company, product and service names, which may be denoted by a double asterisk (**), may be trademarks or service marks of others.