

The Future Direction of High Availability Support Technologies

Suzanne Pherigo
Hewlett-Packard
3400 E Harmony Rd, MS E-8
Fort Collins, CO 80528
Phone: (970) 898-4943
Fax: (970) 898-7126
Email: suzanne_pherigo@hp.com

Bill Dieterich
Hewlett-Packard
3400 E Harmony Rd, MS E-8
Fort Collins, CO 80528
Phone: (970) 898-2640
Fax: (970) 898-7126
Email: bill_dieterich@hp.com

Introduction

In the not so distant past, there was a clear distinction between Mission Critical environments and non-Mission Critical ones. Financial institutions, airline reservation systems, and telecommunication companies were obvious examples of companies with mission critical environments. Servers in these environments were usually dedicated to specific tasks and the environments were kept as stable as possible. Any loss in system availability meant a break in business. A downed airline reservation system could cause chaos around the world. While these companies still represent important mission critical environments, the list of mission critical enterprises continues to grow. In addition, many “old world” mission critical installations are changing their environments to respond to new business models such as internet B2B and B2C paradigms. These new models cause changes their mission critical needs as well.

In today’s internet-based world, most businesses now have some mission-critical needs. Establishing and/or expanding e-business initiatives and exposing more business processes to their customers via Web-enabled applications make high availability and the reduction of unplanned downtime more important. System and application downtime can cost incredible amounts of money and the loss of a company’s reputation. The quality of their Information Technology environment directly impacts their quality of business and success for many companies depends on their ability to acquire, organize, and exploit information.

“New” or “old” business models demand high availability for the hardware and software stack. HP has long recommended that the best way to achieve high availability is through the careful design and balance of 1) the hardware and software infrastructure; 2) the IT processes to manage the environment; and 3) the support partnerships. Environments must be architected and built with the right level of hardware and software reliability, redundancy, performance, scalability, manageability and security. IT processes, such as configuration management, change management, performance and capacity management, security management, backup and disaster recovery must exist and work effectively. Partnerships must be defined and built between the IT organization and the support

provider so that the appropriate reactive and proactive service elements can be applied to the specific environment. Support technologies play a key role in achieving high availability in all mission critical environments and as environments continue to become more flexible, dynamic and “always on”, the importance of support technology is likely to increase.

This paper will examine the role of support technologies in achieving high availability today and how HP’s Support Technologies will continue to address the growing needs of both traditional mission critical customers as well as the newer internet-based ones.

Reducing Downtime

High availability is all about reducing downtime, especially unplanned downtime. While planned downtime is inconvenient and can cause disruption in services, most enterprises have policies to help manage planned downtime and the consequences that it brings. Unplanned downtime, however, can cause severe disruptions in an organization as well as lost employee productivity and revenue.

Infrastructure, processes and support technology are all key elements in reducing downtime. Infrastructure components, such as servers, software and network components are being architected for redundancy (e.g. by providing mirrored disks, clustered systems, etc) as a method to minimize downtime effects. Processes, such as change management and backup and restore processes, are becoming better defined and controlled to minimize the impact of human error and other problems on environment availability. Support technologies augment the built-in system and environmental capabilities. Support technologies are aimed at both reducing the number of downtime events as well as the duration of each downtime event. This requires proactive monitoring and analysis as well as reactive diagnosis capabilities to exist. Other important elements of support technologies include measurement of service level agreements (since you can’t change what you don’t measure) and root cause analysis to help recognize and prevent class problems.

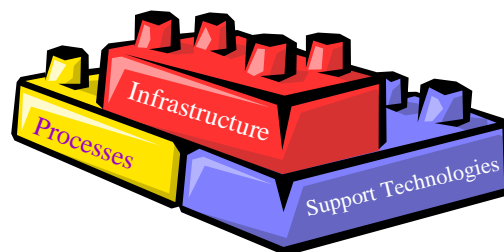


Figure 1: Elements in Achieving High Availability

Hewlett-Packard's High Availability Observatory (HAO) is a leading example of today's support technology. The HAO is aimed at both preventing downtime events and reducing the elapsed time of downtime events that do occur.

Worldwide Telecom, Inc. Today

Worldwide Telecom, Inc., a fictitious company, will serve as a typical example of today's Mission Critical environment. Worldwide will be used to illustrate today's mission critical environments and needs and how those environments and needs are changing over time.

Worldwide Telecom has been a Hewlett-Packard Mission Critical customer for several years. Their mission critical systems/applications include electronic mail, financial systems, desktop applications and customer relationship management systems. They are a distributed company with offices in five cities in the United States, Germany, and Japan. Because of their distributed nature, they have employees working around the clock, who need to access the mission critical applications and systems. Downtime at any time of the day (or night) will cause a business disruption for them. In partnership with Hewlett-Packard, Worldwide recently installed the High Availability Observatory to help monitor their environment in order to keep it running and fix it fast when any problems are experienced.

The High Availability Observatory (HAO) is a suite of support technologies, tools and processes that Hewlett-Packard provides to its mission critical customers today to aid in supporting their mission critical computing environments. The HAO consists of a Support Node workstation and network router that resides on the customer site, a secure connection back to HP's Mission Critical Support Center, and equipment and software within the Mission Critical Support Center (MCSC) to maintain and analyze information about the customer systems. Configuration and status data from systems, software, and network interconnect devices is collected via the support node and securely transmitted to the MCSC. This data can then be viewed by qualified HP support engineers or securely shared with HP experts to help solve customer problems quickly. The data is also analyzed at the MCSC to alert HP and the customer of potential issues in their environment, such as bad patches or out-of-date firmware revisions. Hardware failure events within the customer mission critical environment are detected by the HAO and alerts are immediately sent back to HP so that appropriate action can be taken. System availability is tracked to monitor downtime events, in order to help understand causes and take preventative measures. Using the HAO, HP can also remotely access the customer environment (with appropriate customer authorization) to analyze and resolve problems quickly, without having to travel to the customer site.

What is HAO ? A High level architecture overview

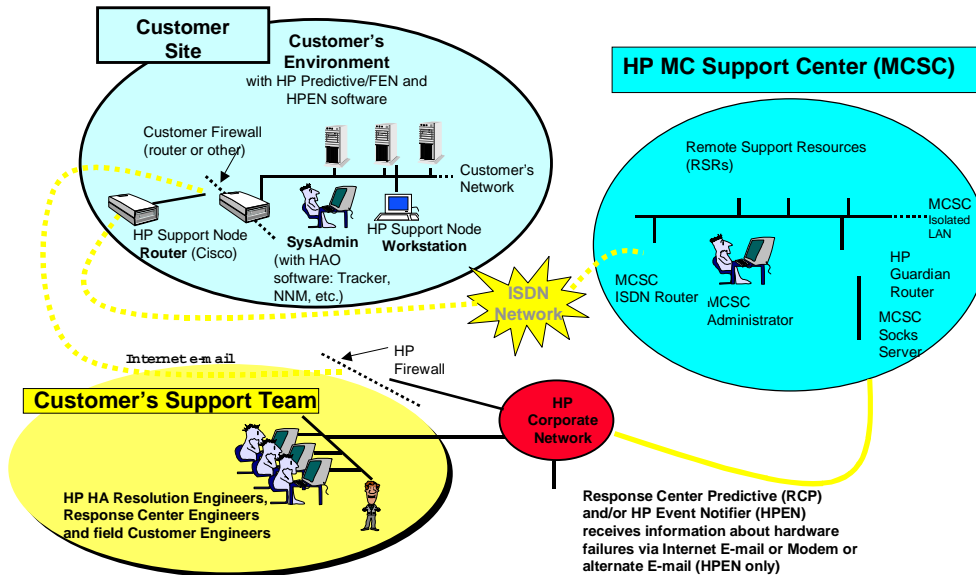


Figure 2: The HAO High Level Architecture Overview

Once the HAO had been installed in Worldwide Telecom's environment, they began seeing the benefits immediately. The first automated analysis of their server and storage configuration data was done within 24 hours of the HAO installation. This analysis discovered 124 disk mechanisms in Worldwide's mission critical environment with out of date firmware that needed replacing. Another automated analysis tool determined that six critical patches needed in the environment had not been installed. The firmware was replaced and the patches installed, thus eliminating potential downtime and data corruption.

Over time, Worldwide continued to experience the benefits of the HAO. Within the first three months of operation, proactive and reactive features of the HAO resulted in downtime prevention and a reduction in downtime duration.

- At one point, the HP representative for Worldwide discovered a growing PDT table while doing routine proactive monitoring from her HP office. The entries in the PDT table pointed to an impending memory problem. The memory was fixed before the system failed, thus preventing unplanned downtime.
- On another occasion, a system panicked. The HP Event Notifier received the message from the panicked system and generated an alarm that caused the HP engineer on call to be paged. The engineer gained authorization from Worldwide and then established the remote connectivity link between HP and Worldwide, in

order to debug the problem. He could not determine the cause of the problem, so called in an HP software specialist in Brussels. The software specialist used the remote link, along with the stored configuration data, to diagnose and fix the problem. The system was back up and running in less time than it would have taken the HP engineer on call to travel to the customer site.

- Over the course of one day, Worldwide began to experience very bad performance with one of their critical applications. The HAO Configuration Tracker was unable to collect information for the system on which this application runs and generated a warning. Upon investigation, it was discovered that a runaway database process had completely consumed the system processor. After correcting the situation, the environment returned to normal operation.
- The HAO data is used in monthly planning and status meetings held by HP and Worldwide's IT leads. Monthly uptime statistics are reviewed. Downtime events, hardware and software calls are correlated and proactive recommendations are made to help avoid downtime in the future.

These examples demonstrate the range of benefits that can be achieved using today's tools and technology¹. Other examples of the uses for the HAO and the data it collects include the following:

- Generating alerts when system components fail, reducing repair time.
- Archiving configuration information, which can aid in restoring systems to their previous state after failures, inadvertent configuration changes, etc.
- Tracking when configuration changes were made
- Performing patch analysis to ensure the right set of patches are installed
- Tracking cloned systems and flagging deviations from "golden" configurations
- Monitoring deviations from known good models
- Providing remote troubleshooting capabilities.

In addition, the data collected from Worldwide Telecom is analyzed along with data collected from other customers to identify class problems (e.g. bad firmware revisions, bad configurations), in order to fix these problems before they impact other customers and to help determine future capabilities that will most positively impact the customer experience with HP's products.

Future Trends

Several trends are developing today that impact mission critical computing environments and the methods that are required to achieve and maintain high availability within these environments.

¹ In fact, each of these examples comes from a real customer situation.

The definition and need for “mission critical” for most companies continues to expand and today encompasses any system or process that is customer facing or will negatively impact employee productivity if unavailable. Unplanned downtime within these environments can have a large negative impact. Outages not only cost lost employee productivity but also can cause negative press, and loss of sales and customers that is not easy to recover.

In addition, the definition of “downtime” is expanding, as well. Today, system up time is often used as the measurement of availability. If a server cannot be utilized effectively, (e.g. if applications are too slow or if the system is unreachable due to network difficulties) the fact that it is up and running is unimportant. The implication is that the measurement tools will need to change along with the monitoring and diagnosis tools.

Business Models for many companies are changing. In addition to embracing the Web and e-business, companies are using partnerships and outsourcing to augment their core competencies. The xSP model, in which technology and business services are accessed “over a wire” is much more prevalent. In this model, the xSPs plan, build, operate and maintain the data centers needed by their end-user clients. The end-user organizations don’t have to purchase, own or manage the technology, but they still care about availability and need to have downtime issues resolved as quickly as possible. And, of course, the xSP providers care deeply about high availability, performance and security for all of their customers.

Technology is becoming much more complex. It is rare to find single-vendor, homogeneous environments. Open systems and mixed environments are the norm. Many companies are global, and distributed applications and environments are required for effective business operation. Dynamic allocation of resources to more effectively manage performance and capacity is becoming more common, as is redundancy of all technology components.

Labor is becoming more expensive and harder to find. As technology becomes more complex, it is impossible for one person to build up all of the expertise required to support a complete environment. Thus, the experts are not always available at the right time and place to solve critical problems. In addition, the experts cannot be used to solve known problems over and over again, because their time is so valuable.

All of these trends impact the requirements of support technologies to support the mission critical environments of the future. The following figure depicts some of the changes that will need to be made in support technologies.

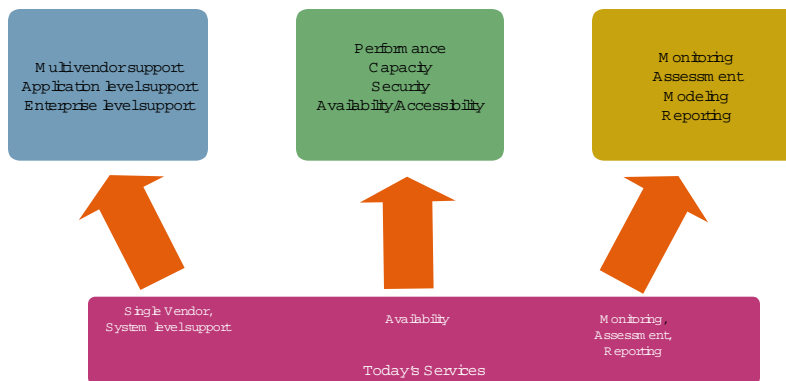


Figure 3: The Evolution of Support Services

Today’s services are primarily focused on system level data analysis and support. As environments become more distributed and more redundancy is built in, it becomes more important to provide application support and enterprise level support, rather than single system support. Environmental health, rather than system availability, becomes the measure of an environment’s ability to perform effectively. Focusing on the environment also means providing multi-vendor support, since environments are not usually homogeneous. In order to do this effectively, it is likely that partnerships and outsourcing will become increasingly common in the support world, as well. Because of this complexity, customers will put a high value on support providers that can deliver seamless support across multiple vendors.

Availability is still a key measure of environmental health. Other measures, such as performance and capacity, are important as well. In fact, in a “redundancy everywhere” model, system availability is secondary. If a system crashes, another one can transparently take its place without an interruption in service, so an immediate fix of the crashed system is not required². System performance and capacity problems can be even worse than system crashes. If a system remains running, but does not have the capacity or performance requirements to effectively support its applications and users, an immediate fix to the situation is required. Thus, performance and capacity monitoring, and being able to take action when thresholds are met, is critical in the future mission critical environment.

Security is becoming increasingly important as well. In the xSP model, for example, the data center is no longer isolated from the rest of the world. Multiple customers are supported out of one data center and their data must be isolated from each other as well

² It is still important to capture the system crash event, in order to ensure that the system does get fixed before it is needed again and to perform root cause analysis on the system crash.

as protected from outside attacks. Even in non-xSP business models, the internet is being used for business transactions, intra-company email and data exchange, and customer service. Sabotage to any of these areas can cause irreparable harm, and thus security becomes a key element of the environment to monitor.

Several capabilities are required to manage these attributes of the environment: real time event monitoring and notification, pro-active assessments, including trending and correlation capabilities, change modeling, and automated problem assessment and resolution.

Real-time event monitoring and notification is important to track the current operations. With real-time monitoring, security violations can be acted on immediately. Performance and capacity issues can be addressed quickly, and failures in systems or components can be addressed as soon as necessary.

Pro-active assessments help tune the environment for peak performance, as well as prevent future problems. For example, trend analysis might determine that on the first and fifteenth of every month, the payroll application experiences heavy usage and the performance degrades significantly (just when it is needed the most). Dynamically allocating additional performance to the payroll system for that time-period will significantly increase the ability of payroll employees to do their jobs effectively. Without knowing the performance trends, however, it is difficult to tune effectively.

Change modeling allows the impact of changes to be modeled before the changes are applied to the environment. This can help prevent inadvertent side effects from adversely affecting the live environment. In addition, it can help determine whether additional resources (e.g. hardware, network bandwidth) are required to keep the environment running at peak performance after the change is applied.

Finally, automated problem assessment and resolution is key when expert labor is hard to find. By capturing expert knowledge and applying that knowledge to known problems, situations can be resolved more quickly without requiring expensive and hard to find labor. Although there are situations where the personal touch is still important, automated help desks and expert systems are effective in many non-critical situations.

Reporting capabilities are important for all of these capabilities. The ability to present information succinctly, at the level and format required, and delivered to the right people at the right time is essential. Information overload remains a large problem for most IT personnel and if the appropriate information is not received and acted upon, it is essentially useless.

The next section will examine how Worldwide Telecom is changing and how new support technologies will continue to help Worldwide remain a competitive company with a highly available and useful environment.

Worldwide Telecom, Inc. Tomorrow³

Worldwide Telecom continues to operate electronic mail, financial systems, desktop applications, and customer relationship management systems as mission critical systems/applications. Over the past year, they have created a new customer portal that allows customers to view and manage account information, review service plans and make service changes. In addition, they have created a data center to host this customer portal, along with several smaller web-based businesses, for which they provide customer management and billing services. Both of these new endeavors are critical to the future success of Worldwide Telecom. They require a very high level of availability in these applications, along with their existing internally focused ones.

Worldwide has re-architected the data center for redundancy. All components have a fail-over mechanism, to minimize the effects of any downtime events experienced. In addition, they utilize capacity on demand and performance on demand capabilities to quickly modify the environment to support peak demands from various customers at various times.

They have expanded their support partnership with HP who continues to provide support solutions for their environment. The High Availability Observatory continues to operate within their IT infrastructure. In addition, the HAO and new support technologies have been embedded in the data center to provide support solutions for that environment.

The basic architecture for the support technologies has not changed. Information about the customer site is collected and securely transmitted back to HP, where it can be analyzed and used to help troubleshoot problems. The breadth and depth of the data collected and the analysis that is done has expanded, however. Over the first few months of operation, Worldwide Telecom experiences the following new benefits from the support technologies installed in their data center:

- Whenever a system or component fails, a failure event is generated and the information about the failure is sent back to HP. In most cases, no immediate action is required because of the redundancy built into the environment. If immediate action is required, the HP support engineer can log into Worldwide's environment to troubleshoot the problem and hopefully resolve it without having to travel to the Worldwide site. In all cases, however, the events are captured. Once a day a "Parts List" report is generated, which indicates which components have failed that need to be replaced (most days the report is empty). This ensures that the redundancy is not compromised within the environment.
- Performance and capacity thresholds are monitored and if exceeded, an alarm is generated. In some cases, the environment can automatically adjust compute

³ This section examines a potential future and is not a reflection of committed support technologies from Hewlett-Packard.

- power and capacity to take care of the problem, but in other cases, manual intervention is required.
- Periodic assessments are done to monitor the health of the environment over time.
 - Security analysis determines if there are potential security issues present in the environment. Information about the violations and how to rectify them are automatically sent to the appropriate personnel on a periodic basis. Issues that are not resolved are flagged as a continuing problem.
 - Performance and capacity analysis monitor the performance and capacity needs in the environment over time, as well as trends in increasing (or decreasing) needs and spiking trends over time. Suggestions are generated for better ways to configure the environment to lessen the impact of spikes or to suggest the best methods to meet the increasing needs.
 - Configuration analysis compares the environment to “known good models” (and “known bad configurations”) in order to suggest ways to tune the environment for peak performance.
 - Availability or accessibility of transactions and applications is monitored. Problems are reported immediately via alarms. The environment is assessed periodically for potential availability problems (e.g. non-redundant components, bad configurations, etc.). Periodic reports are generated which summarize availability over time and point out potential issues that could impact the environment in the future.
 - When a change to the environment is proposed, change modeling is done before the change is actually implemented. Example situations include:
 - How will installing a new application impact the environment? What additional components will be required to support it adequately?
 - How will tuning certain parameters impact the performance of the data center as a whole?
 - Non-critical issues are often addressed on-line. Expert knowledge and problem resolution techniques are available any time of night or day. Of course, live experts are still available when the problems are critical or cannot be resolved using the on-line technology.
 - Back at HP, consolidated customer information is analyzed to help find and resolve class problems. Recalled components can be quickly found at all monitored customer sites and replaced before they fail. Bad configurations can be flagged and fixed before causing problems in multiple customer environments. Root cause analysis on failures is done in order to fix software and components before they impact other customers.

Conclusions

The role of support technologies for mission critical environments is becoming increasingly important. A recent Gartner Group study concluded that in mission critical environments, hardware failures account for a quarter of all unplanned downtime; software accounts for about another quarter; and the network for just more than a fifth.

Increasing environmental availability involves reducing both the number and duration of downtime events experienced within the environment. In order to do this effectively, all of these components (hardware, software, and network) must be made redundant, proactively monitored and fixed quickly when problems do occur. Both preventing problems from occurring and fixing problems that do occur are extremely important. Thus, reduction in unplanned downtime is still a major goal in supporting mission critical environments. It is not the only one, however. Keeping performance at an acceptable level, having adequate capacity to support growing needs, and ensuring security within the environment are all requirements for healthy mission critical environments. Support technologies must encompass all of these aspects of environmental health to continue ensuring the success of mission critical environments in the future.