DATA PROTECTION STRATEGY



MAINFRAME DISASTER RECOVERY PLANNING

BY

JON WILLIAM TOIGO

CEO TOIGO PARTNERS INTERNATIONAL

CHAIRMAN DATA MANAGEMENT INSTITUTE



SUMMARY

In organizations that have them, mainframes typically host the many of the applications regarded as most critical to their firms from a business standpoint. Ensuring the availability of these applications and their data against a host of interruption threats – as well as business and technological change -- is a non-trivial technical and procedural challenge.

In the past, recovering mainframe-based applications and data has been facilitated by the resiliency, manageability and standardization of the mainframe itself. However, as new (and in many cases, non-traditional) workloads find their way into the mainframe environment, the recovery of mainframes is becoming more technically complex.

This paper surveys the techniques we have used in the past and the new challenges we will all shortly confront in mainframe disaster recovery planning.

INTRODUCTION

For all of the changes in technology and business over the past 25 years, little has changed in either the methodology or the scope of mainframe disaster recovery planning. For example, a methodology, first advanced by this author in the early 1980s, is basically as applicable today as it was when first written. Its foundations are in systems development lifecycle management methodology (SDLC) from IBM circa 1980, and it still identifies the essential activities comprising any disaster recovery planning project.



As shown in the diagram, the DR planner begins by identifying the assets that need to be protected. Correctly performed, this activity begins at the business process level. Planners need to assess the criticality and priority of each business process. Then, they need to identify the data and technology assets associated with the support of that process.

Neither data nor technology is important in and of itself. Their criticality is inherited, like so much DNA, from the business processes they serve. Thus, only in the context of the business process can effective continuity planning objectives be set and appropriate recovery plans developed.

Once business process criticality (and by extension, application and data criticality) are defined, and objectives are set for their recovery, the next phase of DR planning involves the design of strategies and services that can be applied to recover access to data for the applications and end users who need them and within the timeframes determined to be appropriate and necessary for the business process.

"Time to data" (sometimes called Recovery Time Objective) is the ultimate metric for determining the fit of a strategy to a requirement. It is also the best standard against which to judge the efficacy of any recovery strategy.

Time to data refers to an aggregation of the time required to accomplish three basic things:

- 1. Re-host the application at an alternate location
- 2. Connect the application to a valid copy of its data
- Re-establish user connectivity to the application and data so that meaningful work can be performed.



In the past, a key strength of the mainframe was the comparative ease with which these three goals could be accomplished. Given the rigorous standards articulated by IBM around peripheral device integration, operating systems level management controls, and application hosting, recovering a mainframe was considerably less complex than, say, recovering a distributed server and distributed storage infrastructure where standards are less rigid, equipment interoperability is not guaranteed, and coherent management is generally lacking.

Following the design phase, planners need to develop a sustainable approach for testing recovery strategies and for managing the impact of business and technology change on the recovery strategies themselves. Testing is the long tail cost of planning, so it is important to planners to think about how they will test a recovery strategy as they select a strategy for use in the recovery plan.

This three-phase methodology remains germane to effective DR planning for IT. However, with the migration of non-traditional workload into the mainframe environment and the evolution of new hardware peripherals that resist centralized management via established mainframe utilities, the range of efficacious strategies for mainframe recovery is shrinking.

MAINFRAME RECOVERY STRATEGIES

Over time, a range of strategies have been articulated to recover mainframe operations – beginning with a *laissez faire* approach and culminating in strategy of full redundancy.

Laissez Faire

A laissez faire (or do nothing) approach is the "straw man" in most discussions of the spectrum of mainframe recovery options. Yet, to many organizations, the strategy continues to be viewed as having merit.

Essentially, organizations adopting the risk posture embodied in laissez fare make no advance provisions for a recovery facility, recovery hardware or network redirection (required for reconnecting users to apps over distance). In some cases, the selection of this strategy is based on empirical data about disasters – in particular, the fact that only 5% of disasters are of a type that compromise existing data center facilities or computing equipment. Given the low likelihood of a "smoke and rubble" disaster, the expenditure of OPEX and CAPEX to build a recovery capability that in all probability would never need to be used favors a "do nothing" approach.



No advanced network redirection (location unknown)



No replacement equipment: "next box off the line"





"Laissez faire"

Most laissez faire strategies are not technically "do nothing" strategies. In many cases, while provisions are not made for a recovery site, processor and peripherals, planners work out deals with equipment manufacturers to drop ship "the next box off the line" to whatever destination the client designates following a disaster.

Moreover, a laissez faire strategy does not mean that measures are not being taken to safeguard the data asset – next to personnel, the most irreplaceable asset of any firm – through a program of backup or replication. Moreover, it may reflect increased attention to careful management of infrastructure operations and security and increased spending on hazard mitigation to prevent avoidable disasters from materializing into the disruptive variety.

Since there is no pre-designated recovery site, however, a key impediment of this strategy for successful business recovery within defined time to data requirements is the lack of pre-defined redirection of voice and data networks. Facilities may be found on the fly, equipment may be shipped and installed in an acceptable timeframe, and data may be reloaded from tape – but

none of these activities means much to recovery if the applications cannot be accessed and used.

Service Bureau and Mutual Backup Options

Two additional mainframe disaster recovery strategies that have appeared in the past are the service bureau option and mutual backup option. These are similar to each other in terms of what they accomplish and their potential shortcomings.

Service bureau-based disaster recovery entails the shifting of workload associated with specific applications between the site affected by a disaster and the vendor of the application itself. It reflects a mostly by-gone era in which enterprise application software vendors maintained infrastructure of their own, intended to provide software as a service to customers who preferred not to field their own IT infrastructure.

While software-as-a-service is garnering renewed attention as a part of the current "cloud computing" discussion, the benefits and drawbacks of relying on a software service bureau to "take up the slack" if a customer's local processing capabilities are compromised remain basically the same today as in the past.

On the upside, a service bureau provides a predefined facility for recovery. However, it may not

provide a one-for-one replacement of the user's own infrastructure. Running applications and databases in their own logical partitions (LPARs) does make them nominally transferable to LPARs in similarly configured mainframes. There are, however, many nuances to consider, from specialty processors that are leveraged by the LPAR workload, to customized DASD configurations, to access security mechanisms and customized backup processes, and even to version levels of OS and application software. These differences can impact the portability of application workload and the performance of applications once re-hosted.

A big issue with service bureau recovery is the difficulty in testing the strategy. Questions abound: how much "test time" can the service bureau provider offer the planner to validate the strategy? Another practical issue has to do with the number and availability of multiple service bureaus to re-host all of the critical applications of a given firm: will expeditious

"Service Bureau"



Network redirection possible with careful planning.



Partial workload recovery on like or unlike equipment.



Service bureau facility pre-defined.

recovery require the coordination of many service bureaus processing different portions of the production workload? If so, how will user networks be configured to provide the right access to the right applications?

Many of these same concerns apply to mutual assistance agreements. This strategy essentially makes two (or more) companies with similar workload and similar infrastructure mutually responsible for all or a portion of each other's workload in the event of an interruption.

While this approach seems "neighborly" and could reduce the costs to each company that would naturally accrue to building their own redundant facility, the challenges of host equipment configurations, network re-direction and testing persist. While two or more firms could work together to build a common redundant facility (a strategy variant) that any company could use for testing and actual recovery, sharing resources effectively usually involves an intimate coordination of configuration details and considerable on-going coordination of configuration changes and maintenance.



Network redirect may be an issue.



Host equipment configuration may be an issue.





"Mutual Backup"

In the experience of many firms that have tried to develop mutual recovery arrangements with a peer, the answer has been gleaned over time to the eternal question of why tall fences make the best neighbors.

Cold Sites

A cold site is a subscription-based or privately-owned facility providing necessary environmentals, power, networking and physical security to receive a mainframe and its peripherals following a disaster. This strategy addresses two shortcomings of a laissez faire approach:

- 1. It provides a known recovery facility, and
- 2. It provides an end point for network re-direction.

However, a cold site strategy's success is still based on the ability of a vendor to provide equipment in a timely way following an interruption event – similar to laissez fair. Moreover, testing a cold site strategy with anything more than a paper-based procedural walkthrough is

impossible. Depending on the location, and the type of disaster event, the cold site strategy may be useless.

One variant of the cold site strategy provides the site with data storage peripherals, enabling recovery DASD to be mirrored (kept synchronized with) production DASD via WAN-based data replication or via on-going tape shipments. However, what a cold site is primarily selected to avoid is the labor cost associated with disaster recovery. Any capability that requires personnel to be permanently or occasionally positioned at the cold site facility would mitigate this value.

Commercial cold site facilities are typically offered by commercial hot site facility providers (see below) as an additional service that will be used to gracefully re-host applications some weeks following a disaster event that has caused a firm to leverage its hot site. In those intervening weeks, the hot site service provider works

Network redirect can be prepared in advance.



Replacement host must be obtained on the fly



Recovery facility pre-determined.

"Cold Site"

with vendors to populate the cold site with the equipment needed to host the displaced customer's workload for a more protracted period of time.

Commercial Hot Sites and Private Redundant Facilities

Hot sites are a fixture in the lore of contemporary mainframe DR. They are subscription facilities with pre-installed hardware and software that are ready to host the workload of the distressed subscriber within hours of a disaster declaration. In essence, hot sites provide the same value as privately-owned redundant data centers, but they are shared among multiple companies -- with all of the benefits and drawbacks of such a scenario.

Hot sites appeal to planners who confront applications with short time-to-data recovery requirements, but who lack the interest, resources or budget to field a corporate-owned redundant facility. Hot site subscriptions are designed to be affordable, with many providers offering an *a la carte* menu of additional services intended to deliver measurable recovery speed improvements. Some providers offer optional WAN-based data replication or electronic tape vaulting so that data restore timeframes can be abbreviated in an actual disaster. Others offer high availability WAN-based clustering and failover options.

Basic to any hot site service is the provisioning of a basic hosting platform, basic peripherals and voice and data networking equipment. Subscribers with specialized equipment needs must

either provide the additional gear themselves or lease it from the hot site provider at an additional cost.

Some vendors offer remote command and control facilities to mitigate the need to send many personnel to the hot site facility in the wake of an outage -- a problematic requirement when the disaster event has a broad geographical footprint (such as a hurricane), or an impact on public transportation (as in the case of airline travel suspension following the 9/11 attacks).

One downside of most hot site contracts is a disaster declaration clause. Before taking possession of contracted space, a subscriber typically must "declare" a disaster through a



Network redirect Prepared (assuming site available)



Recovery host ready (assuming site available)



Recovery facility known (first come, first served)







"Redundant Facility"

formal process that includes the payment of a "declaration fee." Even then, access to the contracted facility is not guaranteed.

The reason is simple. Most hot site contracts include a "first come, first served" caveat in their declaration procedure: if multiple subscribers are impacted by the same disaster event, recovery facilities are provisioned in the order in which disaster declarations are received. If the provider has many customers within the geographical area impacted by a disaster event, they typically promise only to find hosting *somewhere* for those clients who didn't declare first. While this strategy is viable if the hot site vendor has multiple recovery facilities, it raises many potential issues. For one, the cadre of hot site personnel who have tested with the client over the years preceding the disaster, and who are intimately familiar with the client's procedures, provide an affinity that can be useful in expediting disaster response. This affinity is lost if the client is forced to recover elsewhere. So too is the immediate access to special peripherals and components that the client has pre-positioned at the primary recovery center. Network redirection plans may also need to be re-worked on the fly – not something a company wants to have to do in the hours following a disaster.

Hot sites also continue to work under the cloud of the misbehavior of some operators in the early pioneering days of the business model. In the late 1970s, one firm sold contracts for a hot

site that didn't exist, then skipped the country with the money! While most hot site vendors today bear reputable brands, and consumers are more savvy about contracting for services site unseen, there are reasons to continue to be vigilant.

Most of the above potential shortcomings of a commercial hot site are addressed, of course, by a private hot site. A private hot site is a redundant data center, usually (but not always) positioned at least 50 Kilometers from the primary production site to avoid being impacted by the same disaster event.

Of course, building a redundant data center fully equipped with processor, DASD, networking and other peripherals, and establishing on-going data replication between sites to enable failover, represents a significant capital expense. Staffing that site with the necessary personnel to keep it minimally functioning until a disaster occurs adds more cost. Companies that have them have worked out ingenious ways to justify the expense, from using the second site as a testing and development location when not serving as a hot site, to using the second site as a production facility when maintenance and upgrades are performed at the primary data center – thereby avoiding "planned downtime" altogether. Still, for firms that lack the deep pockets required for full redundancy strategies, one of the other options for mainframe recovery may be the second best choice.

The selection of the appropriate recovery option is complex, but generally speaking the value proposition of each strategy is clear – bounded by cost and time-to-data (speed of recovery).



GUIDELINES FOR PLANNING SUCCESS

Successful mainframe disaster recovery planning, regardless of the strategy that is formulated, is guided by the following:

1. Eschew "scenario-driven" plans: there is no actuarial table for disaster events. Despite over a hundred years of data on severe weather events such as hurricanes, there is simply no way to tell whether, when or where a hurricane will make landfall. Building an effective plan requires preparation for a worst case disaster, disabling the physical plant in its entirety. Structurally, the plan should be developed in a modular way to appropriate portions of the plan can be activated in response to "lesser" disasters.

This guidance is offered based on over 25 years of planning experience. Disasters do not unfold according to planned scenarios. The plan is at best a guide to recovery, never a script. This underscores the need to test plans as often as possible, not only to ensure that changes in the business or technology infrastructure have been accommodated, but also to familiarize recovery personnel with their roles and the interdependencies between their recovery tasks and those of their peers. Truth be told, when a disaster happens, no one reads the plan.

 Focus on prevention: Whether you have a large budget for planning or you are operating on a shoestring, it is important to recognize that most outages – 95% by some analyst estimates – are not caused by smoke and rubble events. They are the result of avoidable problems.

In 2004, Gartner issued a report on the leading causes of downtime that has been validated in survey after survey of companies impacted by outages in the first decade of the new Millennium. As shown in the chart below, the biggest percentage of outages accrued to software, hardware and people errors.

The trends today suggest that the



situation is getting worse, not better. Software complexity is growing and staff time for maintenance is steadily declining for a number of reasons. The trend is for software related faults to grow as a percentage of outage causes. Hardware is also getting more unwieldy with more and more products being attached to mainframes that are not directly managed by the systems level management technology that has been the source of much of the mainframe's vaunted stability and reliability. Human error is increasing, both as a function of lean staff and the requirement to shoulder the workload once spread over several staff, and as the increasing access provided to end users via networks to mainframe resources. The only outage component that is declining is "planned" downtime, outages required to perform maintenance and upgrade work, which is being deferred in the current economic reality. Mostly unchanged is the percentage of downtime accrued to milieu-level disaster events.



Based on these observations, a central focus of disaster recovery planning must be disaster avoidance planning, which attacks the root causes of the preponderance of downtime events.

Disaster avoidance strategies focus on issues that, at first glance, may not seem to fall within the domain of disaster recovery planning. These include facility security, application and network security, infrastructure management, hazard detection, annunciation and suppression, power protection, and of course data protection.









Fire Detection and Suppression Systems



Power Protection Systems Surge/UPS/Generator

Active Security Practice And Software

Infrastructure Monitoring And Management Systems

Leak Detection Systems

3. Avoid "holy wars" over data protection strategies: Central to disaster prevention and disaster recovery is data protection. Data is an irreplaceable asset and the only way to protect it is to make a copy. While there is general agreement on this point, disagreements can become noisy when selecting the best method for making the safety copy. The methodology you choose should be based on the time-to-data requirements of the application itself, first, and on practical issues such as the solvency of the method, its cost and its ability to be tested.

Data protection methods have undergone some refinements over the past decade. For example, we have seen an evolution of tape backup from a disk to tape meme to a disk to disk to tape configuration, with the second stand of DASD used as a virtual tape subsystem (VTS). The VTS has evolved as well, from a location for aggregating backup datasets so that operators can fully fill a tape cartridge, to emulating multiple tape drives in order to expedite backups, to providing a location where services such as encryption and de-duplication can be applied to data.

Similarly, disk to disk mirroring has been extended with asynchronous software tools and onarray replication technology to make data copy over distance a reality. But, despite the hype, neither disk-based, nor tape-based data protection methods are inherently superior. For some applications, whose time-to-data recovery requirements are short, WAN-based disk to disk replication may be the best choice, but for most other applications, tape-based backup and restore is more than adequate – and considerably more cost effective! Most data centers today use a combination of both technologies, demonstrating the fact that there is no one-sizefits-all approach.



4. Business-savvy required: The strategies settled upon for both disaster prevention and disaster recovery must co-exist with the business realities in which they operate. While most competent engineers can create strategies that are reasonably comprehensive and effective in accomplishing their goals, business realities need to be considered.

From the beginning, the case for disaster recovery needs to move beyond its narrow focus on risk reduction value solely. Planners need to provide a fuller business value case that includes cost-savings and improved productivity (top line growth in business management speak) if funding is to be granted.



For every strategy considered, cost factors need to be carefully considered and articulated to business management. If the same result or outcome can be accomplished by several techniques, it is incumbent upon planners to identify the options considered and why the selected option provided a better ROI or was more affordable than other options. Key to this evaluation is the question of how the strategy can be tested, given the fact that testing and change management costs usually represent a much bigger expense than do the initial acquisition of DR products and services.

Improved productivity is a natural by-product of disaster prevention strategies, but the relationship will not be explicitly understood unless it is contextualized that way. By investing in outage prevention measures, such as better infrastructure management and security, the amount of idle labor cost associated with outage events – estimated to average about \$1 million per hour across vertical industries – can be reduced. Many savvy planners contextualize

their facility and equipment redundancies for disaster recovery as a technique for reducing planned downtime for equipment and software maintenance and upgrade.

Finally, the risk reduction argument for DR planning needs to go beyond a statement about the continuity that will be provided for business processes following a cataclysmic event. Smoke and rubble disasters are actually rather rare. Risk reduction should be interpreted to include compliance and legal risks, which are addressed by data classification, policy-driven data protection, and programs of data encryption and secure offsite storage – all of which are commonly part of effective disaster recovery planning.

We need to make a full business value case, using the parameters that management has defined for value, if plans are to be funded.

5. Prepare for new challenges: Testing and change management are critical to mainframe DR planning, not only to spot coverage gaps that result from the gradual evolution of the business and technology infrastructure over time, but also to try new techniques and technologies for disaster prevention and recovery as they enter the market. Truth be told, with every new technology innovation, there are new risks and new ideas for how to cope with those risks. While the distributed computing world tends to see waves of new technology more frequently than do mainframe environments, times are changing.

For example, mainframe workloads are changing. We are seeing an increase in the hosting of non-mainframe applications as virtualized guests inside LPARs as part of a shift in many companies away from expensive or less resilient distributed computing paradigms. Many of these guest applications and their data are not adequately managed or protected by tried and true systems level utility services. In fact, new hardware is being attached to mainframes that were never a part of the traditional mainframe peripheral set and that do not conform to traditional integration standards in order to support the data associated with guest applications.



These new guest workloads and infrastructure will require different, and often stand-alone, management and protection services delivered by software that may seem quite foreign to the traditional mainframe DR planner. Indeed, coordinating multiple backup methodologies, multiple security and encryption processes, and multiple management systems is likely to become a huge challenge in mainframe DR in the near term.

Mainframe DR planners will need to become more directly involved in workload migration planning, in application design and integration projects, and in infrastructure vetting processes to identify potential risks and recovery requirements and to ask the right questions about the manageability and recoverability of the new workloads and technology that are being introduced to the environment. This is a significant change from the traditional role of the DR planner.

CONCLUSION

Mainframe disaster recovery has benefited over the years from hardware standardization, systems-level management, and well-defined and disciplined architectural and operational procedures that contributed to the reputation of the mainframe as "disaster proof." This is quickly changing.

It is no longer sufficient, if indeed it ever was, to consider the replacement of the mainframe box with another mainframe box in a timely way the sole definition of good DR planning. Going forward, planners will need to become more directly engaged in the mainframe planning process and more business-savvy in their presentation of DR requirements and options to succeed in building the right set of strategies for business continuity.

TO THE READER

Parts of this paper are excerpted from a forthcoming book, Disaster Recovery Planning 4th Edition, by Jon Toigo. The book is being published throughout 2010 as an on-line "blook" (blog+book) at http://book.drplanning.org, so that readers can contribute their observations and experience. Please visit the site. There is no cost to register or to contribute your insights!

