# Understanding Recovery Point Objective (RPO) and Recovery Time Objective (RTO) in HA/DR Solutions

*What level of disaster recovery and high availability capabilities do you need? Understanding RPO and RTO can help you to find the answer that best balances costs, risks, and profitability.*

*By Craig A. Johnson*

For a few decades now, all prudent IT departments have had some form of disaster recovery (DR) plan in place. For many companies, that still solely entails sending nightly tape backups offsite and arranging for a site, possibly one owned by a recovery service provider, where data and applications can be loaded onto systems and run, should disaster strike the primary data center. Increasingly, because of more stringent regulations concerning data and system protection and because already high downtime costs are continuing to rise, organizations are augmenting DR with high availability (HA).

Before continuing, two definitions, DR and HA, are in order. In the context of this article, a "disaster" is any event that destroys at least one production server and/or renders all of one or more systems' online production data permanently unusable. This may be an incident that razes the whole data center, such as a natural disaster or terrorist attack, but it can also be something less catastrophic. For example, if a company maintains only one online copy of production data and applications, without using RAID or disk mirroring to protect it, a disk crash may also be classified as a disaster for the purpose of this discussion. DR is, then, the reloading of data and applications at a remote location in the event of a total data center loss or locally in the event of a disk failure.

HA maintains real-time or near real-time replicas of all data, applications, and other objects on a hot-standby backup server. This backup server can be used whenever a production server or its data fails or needs to be taken offline for maintenance.

HA and DR are not mutually exclusive. All HA solutions can provide rapid recovery from the lesser type of disaster (i.e., the loss of only the primary server and/or its data and applications). Furthermore, if the HA topology separates the primary and backup servers by a sufficient distance such that a disaster that strikes one will not affect the other, HA can also provide rapid recovery from the more catastrophic type of disaster.

It is thus appropriate to think of HA and DR not as unrelated technologies, but rather as points on an availability continuum. In fact, there is not just a single DR and a single HA point on that continuum. Rather, there are different levels of DR and HA across the spectrum.

Points on the continuum are defined by two variables, Recovery Point Objective (RPO) and Recovery Time Objective (RTO). RPO defines the organization's goal for the maximum amount of data
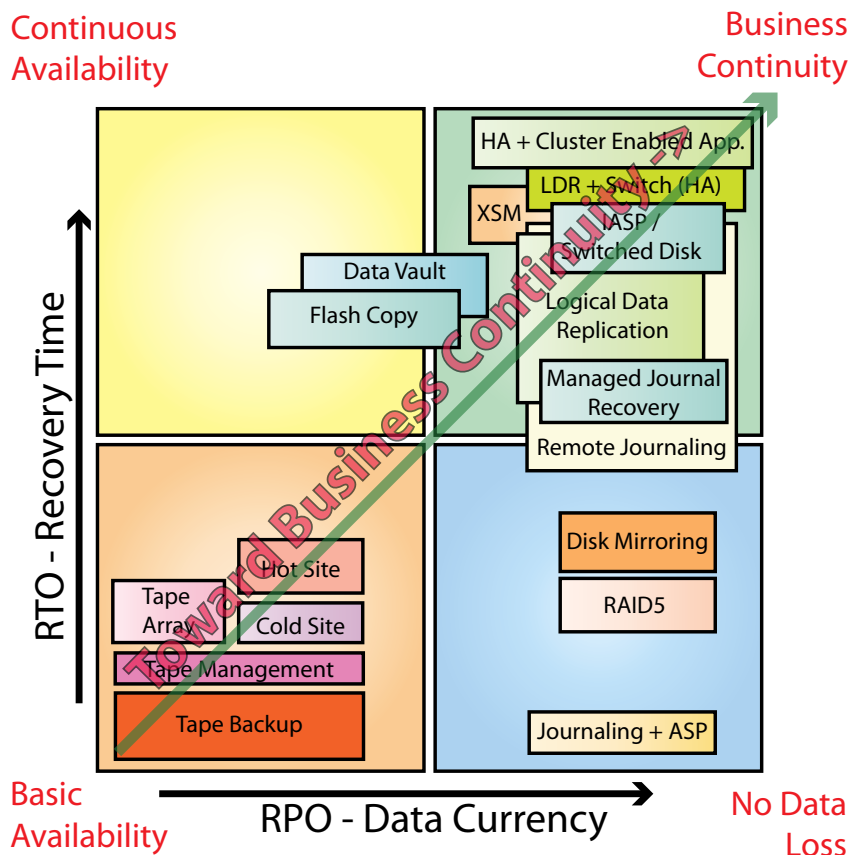
Figure 1: Availability Continuum

that will ever be lost as a result of a single disaster. This is referred to as a recovery point objective because it is the earliest point in the data stream that the organization is prepared to fall back to after a disaster. The loss of data updates applied after that point is tolerated if necessary.

RTO defines the organization's goal for the maximum downtime that the organization will have to endure during any one downtime event.

It should be noted that RPO and RTO are objectives, not certainties. For example, consider an organization that maintains replica servers in geographically remote locations. The organization may set an objective of being able to switch operations to the backup server in no more than 15 minutes. A powerful HA solution may reliably fulfill that objective in almost all circumstances. Nonetheless, it will be thwarted in the exceptionally unlikely event of simultaneous disasters at the primary and backup sites.

As illustrated in **Figure 1**, the availability continuum is two-dimensional, with RPO and RTO defining the axes. This space can be divided into quadrants: Basic Availability, No Data Loss, Continuous Availability, and Business Continuity. A variety of products fulfill the recovery point and time objectives that define each quadrant.

### Basic Availability

The Basic Availability quadrant contains traditional tape-based DR solutions. These serve organizations that are willing to accept long recovery times and the possibility of losing considerable data.



The 5th Wave By Rich Tennant

© The 5th Wave, www.the5thwave.com

"He saw your laptop and wants to know if he can check his Hotmail."

A solution in the Basic Availability quadrant is assumed to not include journaling. (The use of journaling moves it into the No Data Loss quadrant.) Without journaling, if a disaster occurs, data can be recovered only up to the last backup, which typically means the previous night. Thus, 24 hours worth of data may be lost if the production data is destroyed just before the nightly backup.

In reality, the potential loss is greater than that. If tapes aren't shipped offsite immediately, the lag adds to the data loss exposure because, depending on its nature, a disaster may destroy any tapes that are still onsite. In addition, tape is not a perfect medium. If the most recent tape is corrupted, as much as 48 hours of data may be lost.

Recovery times are long in this quadrant because tapes typically have to be retrieved from an offsite location. If the primary data center is intact and only the data is destroyed, the tapes can be returned to the data center. In the event of a disaster that destroys the data center, the tapes must be sent to the recovery location if they are not already stored there.

Data and applications must then be loaded from tape onto disk. Today's high-speed tape drives makes this a faster process than in the past, but companies with particularly large databases may still require several hours, and possibly a couple of days, to load the data.

Tape management tools that reduce the chance of errors and allow operators to find tapes faster, coupled with high-speed tape arrays that load multiple tapes in parallel, reduce recovery times further, but it may still take hours to bring the business back online after a disaster.

The use of a hot backup site (which contains all of the necessary hardware and operating software and needs only to load the data and business applications in case of a primary site disaster) instead of a cold site (an empty room into which all of the necessary hardware, software, and data must be shipped when a disaster is declared) can also shrink recovery times. Nonetheless, the time required to return to business in a hot site will still result in unacceptable business losses for many large enterprises.
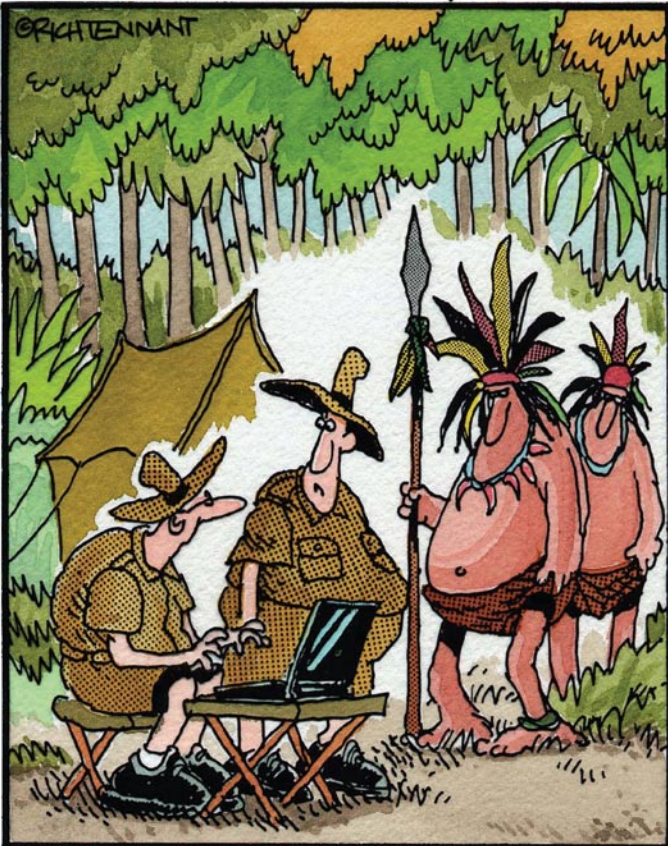
### No Data Loss

The use of journaling and modern application practices can eliminate virtually all data loss in many situations. However, this does nothing to improve recovery times. In fact, it might worsen RTO somewhat as the recovery process still requires the loading of data from backup tapes and then taking the extra step of loading subsequent data from the journal before operations can resume. If journaling isn't used, depending on the nature of the business processes, some operators may be able to begin processing new business immediately after the backup tapes are loaded, while other operators simultaneously reenter transactions lost since the last backup.

Even with journaling, there may still be a risk of as much data loss as for a Basic Availability solution. If journaling is done locally, a catastrophe that destroys the entire data center will destroy the journals along with the production databases. In

that case, the most recent possible recovery point will be the last tape backup that was sent offsite.

**Figure 1** also places RAID and disk mirroring in this quadrant. These technologies protect against data loss due to single point failures, but they can't safeguard data from simultaneous failures that defeat the RAID or mirroring protection. In addition, except when using cross-site mirroring to replicate data to a remote location, mirroring and RAID cannot protect against data loss due to a disaster.

With RAID and disk mirroring, data recovery after a single point of failure is instantaneous, but another issue in addition to the disaster and simultaneous failure exposures keeps these solutions in the No Data Loss quadrant rather than moving them up to the Business Continuity quadrant. A Business Continuity solution guarantees the availability of not only the business' data but also its processing capabilities. Neither RAID nor disk mirroring does that. They will not keep the business running if a server crashes or needs to be taken offline to upgrade or maintain the hardware, operating system, databases, or business applications.

### Continuous Availability

Solutions in the Continuous Availability quadrant focus on recovering business operations quickly but possibly at the cost of some lost data. For the most part, flash copy and data vaulting offerings fall into this space, but, depending on the specific product and how it is implemented, a particular solution may edge into the Business Continuity quadrant.

Flash copy technologies create periodic snapshots of specified objects and transmit them to a second system. These snapshots can be retrieved and loaded onto the backup system quickly to significantly reduce recovery times compared to tape recovery options. By increasing snapshot frequency, the potential data loss can also be kept to a minimum.

Data vaulting captures changes made to a production system since the last tape backup and transmits them to a backup system. Typically, changes are batched before being transmitted. The RPO that can be achieved is, therefore, dependent on the batch frequency. Some vaulting products also offer the option of continuous vaulting, which can achieve an RPO of close to zero data loss.

Recovery times for vaulting solutions depend on the nature of the product used and how it is set up. If the previous night's backup tape must be loaded and then the vault contents applied to it, the recovery time will be about the same as for a tape backup with journaling. However, some products allow for a completely online and fully automated recovery process that can significantly reduce recovery times.

### Business Continuity

The ultimate in availability, which is zero downtime and zero lost data, is achieved in the very top right corner of the Business Continuity

quadrant. The products in this quadrant are all referred to as HA solutions. No existing product can offer a 100-percent guarantee of zero downtime and zero data loss, but some HA solutions come very close.

A complete discussion of these HA technologies is beyond the scope of this article. They include products that maintain real-time or near real-time replicas of all production data and objects on a backup server along with the ability to swap the roles of the primary and backup servers when the primary server fails or it must be taken offline for maintenance. The speed with which this role swap can be performed determines the RTO that can be achieved. Various options within this product category, such as switched disk technologies and clustering, can help to reduce role swap times.

Along with very fast recovery times, HA solutions can fulfill near perfect RPOs. By replicating data and object changes to the backup system in near real-time, the backup system is almost completely current. Depending on bandwidth and processor loads, data and object changes made on the primary system may be vulnerable for a fraction of a second, or typically seconds at most, until they are written to the backup system.

Many HA solutions offer the option of ensuring absolutely no data loss by using synchronous replication, which writes data updates to the backup system before the user transaction is considered to be ⇨

complete on the production system. Thus, the backup system can actually be slightly ahead of the production system. However, most organizations forego this option because users will be kept waiting while the data is replicated to the backup system. This can be an unacceptably long wait if the backup server is unavailable or bogged down or if the network connection is overloaded or down.

When the backup system is located sufficiently distant from the primary system such that a disaster that strikes one will not affect the other, the HA solution is inherently a DR solution as well.

## Continuum Positioning

What is the ideal position in the two-dimensional availability continuum? If you could exclude the price of the solution from the equation, clearly you'd want to move your systems and data to the top-right corner of the Business Continuity quadrant. But faster recovery times and more complete protection against data loss are not free. Trade-offs must often be taken.

Regulations are the first consideration when answering the continuum positioning question. Sarbanes-Oxley (SOX), the Basel II accord, the Basel Committee's Capital Adequacy Directive (CAD III), the Gramm-Leach-Bliley Financial Services Modernization Act, the Health Insurance Portability and Accountability Act (HIPAA) and The Patriot Act, among other laws, all require that applicable organizations protect the availability of certain data and/or be prepared to deliver requested data to the proper authorities within a specific time frame. At a minimum, your RPO and RTO levels must be set such that you can meet the demands of those regulations. There is a price to be paid to achieve this, but if you can't afford to meet the relevant laws, you can't afford to be in business.

Beyond regulatory compliance, there are business reasons why you might want to raise RPO and RTO above their current levels. Downtime costs money. Much has been written elsewhere on how to calculate downtime costs, and there isn't room to repeat that discussion here. Suffice it to say that if you haven't gone through the exercise of summing all of your company's downtime costs (including lost business, idled labor, customer dissatisfaction, and penalties for late shipments or late regulatory filings, among other costs), you'll probably be shocked by the enormity of the number. Thus, if your systems and data are currently in the Basic Availability quadrant, an investment in higher availability will likely deliver a significant ROI.

In the end, choosing the appropriate RTO and RPO is a business decision. The solution you put in place must deliver sufficient value to meet, and preferably exceed, the minimum expected ROI that your company insists on before embarking on any project. To estimate ROI for a proposed availability project, you must calculate hourly downtime costs and forecast the amount of downtime that will be avoided. The product of those two numbers (hourly downtime cost times the forecasted number of avoided downtime hours) determines the value that you will receive from an investment in higher availability.

Some companies are now installing solutions that meet both their HA and DR business requirements by deploying additional servers into their clustered environment. Companies that have very high HA and DR requirements are implementing redundant servers within the datacenter to meet their HA needs and coupling those with remote replica servers to meet their DR and resiliency requirements. Also somewhat new to the HA/DR landscape are organizations deploying solutions that combine different availability options that fulfill business requirements by leveraging the benefits that these integrated solutions can provide.

The appropriate RPO and RTO levels are not the same for all organizations. A small company that receives all of its orders on paper and then keys them into the system may decide that protecting all data entered since the last backup is not a high-value project. That data can, after all, be recreated from paper documents on the extremely rare occasions when the production databases and local journals, if any, are destroyed. The company also might be willing to accept considerable downtime because orders will continue to flow through the normal paper-based route. In contrast, a large online retailer will be willing to invest heavily in a solution that protects all of its transaction data from loss. And, because unavailable systems can lead to millions of dollars worth of customer dissatisfaction and lost sales every hour, that retailer will also be willing to make a large investment in a solution that minimizes downtime.

Even within a single company, different RPOs and RTOs may be established for different systems. To take an extreme and half frivolous example, large financial institutions that offer Internet, phone, and ATM banking will be willing, if necessary, to invest heavily in solutions that protect the availability of their online banking systems and data, but they would not be willing to spend as much to protect the systems and data that maintain their house league baseball scores.

To put it simply, for every organization and every system within those organizations, the bottom line on setting RPO and RTO levels is ROI. It takes a little work to forecast ROI for a proposed availability project, but it is worth the effort.    TUG

**Craig Johnson** *is the VP for Vision Solutions responsible for the research and development of high availability and disaster recovery products for i5/OS. He joined Vision Solutions (formerly Lakeview Technology) in 1995. Craig is recognized as an expert in both IBM and Vision Solutions high availability and disaster recovery solutions available for System i today.*

*Craig is recognized as an expert in the development and delivery of products for System i, leading a development team that has adopted leading-edge development practices and technologies to deliver enterprise-level availability solutions to many of the top companies utilizing availability solutions for their companies' mission-critical business applications.*