

COHE^SITY

Buyer's Guide to Modern
Web-Scale Backup & Recovery

A Comprehensive Guide to Evaluating
Web-Scale Backup and Recovery



Table of Contents

Foreword1

Introduction2

Assessment: Backup and Recovery Challenges Facing Enterprise IT3

Evaluation Criteria for a Backup / Data Protection Solution4

 Hyperconverged Platform4

 Single, Simple Management Interface5

 Cloud Native6

 Limitless, Transparent Scale-out7

 Guaranteed Data Resiliency8

 Quick and Flexible Recovery for Any Size Environment10

 Maximize Storage Capacity and Efficiency11

 Online Upgrades and Expansion12

 Extensible Platform to Consolidate Other Secondary Workloads13

 Security and Compliance14

What Modern Backup and Recovery Looks Like15

Foreword

By: Edwin Yuen

*Sr. Analyst - Cloud, Data Protection, Systems Management, and DevOps
Enterprise Strategy Group, Inc.*

There can be few areas of mainstream IT technology that are more important and more in need of innovation than backup and recovery. While the need to keep data and applications safe as well as ensure their fast and reliable restoration remains of critical importance, the needs for data protection have changed with the modern digital business. Today's digitally transformed IT needs a backup and recovery solution that has kept pace with the demands of our digital world. Unfortunately, the truth is that backup is often treated as a burdensome overhead or a passive insurance policy rather than a source of strategic value.

The reasons are well documented. Exploding data volumes, 'keeping everything forever', heightened expectations of users, fragmentation across physical, virtualized and cloud environments, and a growing patchwork of different piece parts to control it, have led to massive complexity in managing service levels, unnecessary cost, and unpredictable recovery times. In other words, data protection is facing escalating business risk when the intended role for it is to reduce risk exposure.

It is not uncommon for IT to have an escalating number of different solutions from different vendors in place just for backup, accompanied by a battalion of dedicated servers, gateways, storage, networking components, and tape devices that all need to be maintained and coordinated. Then there's archiving, offsite disaster recovery, copy management, access for test & development and so on. The task of protecting and provisioning corporate data and apps has become both complicated and expensive and the changing digital landscape will only continue to make it worse.

But what about the cloud, I hear you ask? Surely the public cloud providers can offload the need for all this infrastructure and just take care of my data for me? Well, discounting the issues around data sovereignty and regulation that are growing in importance, the public cloud providers are not responsible for your data's integrity. Simply put, the public cloud providers are only responsible for providing the infrastructure and access to the storage for your data. At the end of the day, you are still responsible for your data, both backup and recovery - even in the cloud. At best, cloud can often just be yet another target for your existing set of data protection solutions. At worst, cloud adds more complexity, more expense, and more risk exposure to your environment.

All of which brings me back to my opening point: how can we modernize data protection to meet the needs of today's digital businesses? Cohesity has taken a fresh look at the old problem of backup and recovery to try to bring new innovation to the secondary data space. Cohesity took a blank-sheet approach rather than an incremental one, applying Google-class design principles such as hyperconvergence, software-defined infrastructure, and web-scale to design a platform specifically optimized for secondary data operations. The result is a modern, web-scale backup and recovery solution designed for the digital challenges of today and tomorrow.

This guide is intended to provide guidance and assistance for IT professionals looking to re-evaluate their current data protection practices. In addition to describing the new hyperconverged approach, it offers education and best practices for what to look for in a modern data protection solution, a self-assessment, evaluation criteria, and questions to ask when looking at alternative approaches. At the very least I hope you will be more informed as a result.

Introduction

As organizational needs change and workloads become increasingly distributed, a key realization is emerging: traditional approaches to backup and recovery may no longer be sufficient for many organizations. These companies may have discovered that their existing tools are not keeping pace with other advancements in their computing environments, such as scale-out storage systems and hyperconverged systems, which seek to reduce data center complexity and help manage surging storage costs.

Today's backup and recovery offerings are fragmented with point solutions for backups, target storage, and long-term data retention. It's a complex scenario that's less than optimal since each of these silos are designed on proprietary hardware and/or software that has its own maintenance and support contracts. In many cases, customers are only able to perform backup jobs once a day, with such jobs sometimes bleeding into their production windows, resulting in performance degradation. And should recovery be necessary, it can take multiple hours, which is inefficient and can have a serious revenue impact. This is a recipe for tragedy should disaster strike.

Even as the cracks appear in the backup and recovery foundation, organizations are creating and consuming more data than ever before. Data is the new lifeline for enterprises, and it's exploding all around us. Organizations across the globe are experiencing a data deluge, with the quantity of data increasing at an accelerating rate as new types of data are assimilated into existing data systems.

No longer are companies storing only traditional routine business data. Today, businesses are absorbing data from a myriad of sensors and machines; these Internet of Things (IoT) devices are distributed across the enterprise. For some companies, the need to store human-generated data is paramount. Regardless of the data source, the need for change is real.

A significant part of data growth challenges revolves around supporting a far broader set of applications than in the past. As cloud becomes an integral part of the overall IT, modern applications are residing both on-premises and in the cloud. This situation presents opportunities to leverage, as well as challenges to overcome.

The point is that data growth is real and will continue, as will the adoption of an ever-increasing number of applications shows. But, to paraphrase the old cliché, with great data comes great responsibility. Regardless of the source, enterprises need to keep pace with this data growth as they consider backup and recovery capabilities.

Historically, organizations have invested significant portions of their IT budgets in data protection tools. Unfortunately, many of these organizations are working with a severe disadvantage, and as a result, at significant risk. They continue to protect their data with legacy backup and recovery solutions that were designed for a different era of computing.

With some products having their origins tracing back decades, many such tools are simply unable to keep pace with modern business requirements and constant advancements. As business needs evolve, organizations must strive to stay ahead of growing IT complexities. Enterprises need IT infrastructure, including data protection and recovery tools, that's simple to manage, agile, and easy to scale.

If you're considering an upgrade to your existing legacy backup and recovery capabilities, you're far from alone. Gartner predicts that *"... By 2021, 50% of organizations will augment or replace their current backup application with another solution, compared to what they deployed at the beginning of 2017."*

Choosing the right enterprise-grade, modern, web-scale backup and recovery solution involves understanding the problem first, then comparing solutions to those problems. This guide is your handbook for selecting the product that best meets the needs of your organization, for today as well as tomorrow. As your business evolves, you need a backup and recovery platform that can evolve with you.

The primary objective of this Buyer's Guide is to help you determine criteria as you consider your next backup and recovery solution. Inside this guide, you will discover:

- The fundamental questions to ask about your current environment
- Common IT data protection and recovery issues currently faced by enterprises
- What to look for in your next backup and recovery solution, to future-proof your investments
- RFP-ready questions you can use to ensure that your selection process is complete

As is the case with so many technologies, the community around the data protection space is chock-full of often conflicting opinions; and some of the information out there is calculated to create fear, uncertainty, and doubt in the minds of buyers. This guide will help you cut through the noise, and provide direction and confidence as you journey through your data protection and recovery options.

As you review the evaluation criteria for a modern backup and recovery solution, take particular note of the questions you should ask vendors. If you're talking with them in person, make sure you get complete, clear, concise answers to these questions. If you're preparing an RFP for a new backup and recovery solution, these questions will help you get the answers you need.

Assessment: Backup and Recovery Challenges Facing Enterprise IT

Before you get started with a selection process, it's important to understand any potential shortcomings in your existing backup and recovery environment. Table 1 outlines the common backup and recovery challenges faced by IT. Take a few minutes to determine your starting point to help best discover where to kick off your search.

Issue	Description	Present in Your Environment?
Silo'd infrastructure	Does your current solution require the use of backup software along with separate media servers and dedicated storage targets?	<input type="checkbox"/>
Multiple, fragmented UIs to configure backup workflows	Does your current solution require you to use multiple user interfaces to create backup workflows? For example, do you need to create a backup job for a physical server in one solution's console; a separate product's console to configure a different backup job for a virtual server, and yet another product and console to help protect SQL Server?	<input type="checkbox"/>
Bolt-on cloud gateways	If your current system provides any kind of support for the public cloud, does it require you to deploy separate bolt-on cloud gateways that act as intermediaries between your on-premises and public cloud-based backup and recovery environments?	<input type="checkbox"/>
Forklift upgrades for scale-up	When the time comes to grow the backup and recovery environment, do you need to schedule downtime for the backup and recovery environment to add more nodes?	<input type="checkbox"/>
Slow restores, last point in time only	Does your current solution suffer from performance problems at restore time, resulting in the potential for RTO misses? Or does it just allow you to recover from the latest backup made?	<input type="checkbox"/>
Variable and/or fixed block deduplication with compression	Does your current solution lack data reduction features? Or, if it includes them, does it have small deduplication domains or fixed block deduplication capability only?	<input type="checkbox"/>
Disruptive, pre-planned upgrades	When a new software version becomes available for your current solution, do you need to go through a potentially disruptive process to adopt the new version? Or does the upgrade require a high level of planning?	<input type="checkbox"/>
Backup data cannot be reused; infrastructure remains silo'd	Does your current solution provide <i>only</i> backup and recovery services? Does it allow you to reuse the data for other use cases such as test/dev? Does your current product allow you to increase efficiency by consolidating other workload storage environments?	<input type="checkbox"/>

Table 1. A checklist for profiling your current backup and recovery environment.

If you checked one or more of the boxes in Table 1, consider reviewing options for replacing your existing backup and recovery software. There are options on the market today that have rethought backup and recovery with a modern, web-scale point of view to address all of these deficiencies in elegant, simple, and affordable ways.

Evaluation Criteria for a Backup / Data Protection Solution

The evaluation criteria you used for yesterday's backup and recovery solution no longer applies. A new way of thinking is necessary. It's time to stop thinking about data protection and recovery as a silo of services and responsibilities. Today, these capabilities need to be a core part of your infrastructure, not a bolted-on afterthought that adds more complexity to the environment.

Further, today's modern application-driven workloads require smooth data mobility between on-premises environments and the cloud. Organizations need backup and recovery solutions that are agile, offer infinite levels of flexibility, and can respond quickly to changing business requirements.

The sections that follow detail the critical attributes of a modern backup and recovery solution.

Hyperconverged Platform

With IT under more pressure than ever before, and as organizations seek to complete critical digital transformation initiatives, there is decreasing tolerance for complex data center environments to support what is considered routine operations. As crucial as backup and recovery solutions are to the ongoing health of the business, you can't afford solutions that introduce needless complications.

From a pure hardware perspective, what does a modern, forward-thinking backup and recovery architecture look like? As you consider how your backup and recovery environment correlates with your production environment, you'll discover that you need a backup and recovery solution that isn't static. It needs to be able to grow with you as your production environment grows, and it needs to be as flexible as you've come to expect from your production environment. For example, as you add new applications, your backup and recovery environment shouldn't take weeks of planning and hours of downtime to expand to accommodate them.

As you review your backup and recovery options, consider what you need to buy. With some solutions, you need to procure:

- Backup software
- Compatible tape devices for long-term archival
- Media servers
- Storage devices

A key goal in your data protection journey should be simplicity. Modern data center architecture options—such as hyperconverged infrastructure—can help you to bring simplicity to your backup and recovery environment. Hyperconverged infrastructure has emerged as a leading contender for primary production workloads; and thanks to the nature of hyperconvergence, it's an even better fit for secondary environments, including backup and recovery needs. Even if you're not using hyperconvergence in your tier 1 (or primary workload) setting, it's likely to be an excellent option for your data protection needs.

Why? What makes hyperconverged so well-suited to support comprehensive data protection needs? First, it eliminates silos. You already have enough silos in IT. You shouldn't need to break down your data protection environment into compute, storage, and software silos to be managed separately from everything else.

The opportunity to consolidate these disparate silos into a single, unified solution that integrates the backup software and target storage in a single solution makes it possible for you to quickly and easily deploy – and scale – your backup and recovery environment as needed. There's no messy process of choosing only the hardware that aligns with a set of hardware compatibility list requirements.

In the context of backup and recovery, hyperconverged can also apply to the kinds of services you can protect. A modern backup solution needs to ensure that all the workloads you're operating can be protected. As such, it needs to be able to

protect your bare-metal servers as well as all your virtual machines (VMs), regardless of hypervisor; provide full support for your enterprise applications and databases; and protect your primary storage systems and NAS devices.

While evaluating the next modern backup solution, the storage architect should look at an offering that eliminates silos by consolidating backup software and target storage on a single solution. And that solution should seamlessly sync with long-term storage media, whether it's public cloud or tape-out.

Questions to ask vendors about their backup and recovery abilities:

- How do you help me simplify my backup and recovery environment?
- Does your solution help me eliminate hardware and software silos?
- Can the solution protect all my workloads wherever they're running, whether they're on bare metal, or in a virtual machine?
 - For virtual machine-based workloads, how deeply does the solution integrate with the hypervisors in use in your organization? For example, if you're running VMware vSphere, does the solution leverage VMware VADP for changed block tracking?
- Does the solution provide application-level support for common enterprise applications, such as Windows Server, SQL Server, SharePoint, Exchange, and Oracle?
- Does the solution provide native protection for primary storage arrays from companies such as Dell EMC, Pure Storage, and NetApp?
- Does your solution span beyond the core data, i.e., core, cloud, and edge?
- Is your solution locking me into your proprietary hardware?

Single, Simple Management Interface

Ease of deployment is important in any new system. But once the systems are up and running, management becomes the next potential pain point. And unlike deployment, management is something you do every day, making ease of use not just a "nice to have" feature, but an absolute necessity.

Modern backup and recovery solutions need to go beyond user-friendliness, however: they also have to provide comprehensive support for all data protection use cases from a single console. At first glance, it may seem okay to have different consoles to back up different kinds of data, but backup and recovery isn't about backup. It's about recovery.

When you're in recovery mode, minutes count. Confusion can be high and answers in frustratingly short supply. You don't need an ornery recovery administrative experience thwarting your recovery efforts.

Don't take this to mean that it's OK for deployment processes and backup consoles to be complex and sport multiple interfaces. Nothing could be further from the truth. From beginning to end, a backup and recovery tool needs to provide a simple, intuitive, single interface that can manage all aspects of the environment, including initial deployment of the backup and recovery hardware and software, ongoing management of backups, and eventual recovery.

Questions to ask vendors about their solution manageability:

- Describe the administrative experience that users encounter during the deployment phase of your solution.
- How many consoles are required to provide comprehensive backup support for bare metal servers, virtual machines, and container environments?
- Do you provide a single console for recovery efforts?
- With your solution, do you use the same console for the creation of backup jobs as well as recovery tasks?

Cloud Native

Use of the public cloud is exploding and is expected to increase even more dramatically in the coming years. As a result, today's IT infrastructure must go well beyond the four walls of the local data center. This expansion includes your data backup and recovery environment.

According to Gartner, *"... by 2020, the number of enterprises using the public cloud (IaaS) as their data center backup destination will double, increasing from 10% at the beginning of 2017."* In other words, more companies than ever will turn to the public cloud as their backup destination of choice, rather than using a secondary on-premises data center.

Regardless of your current adoption of public cloud technologies, your next backup and recovery product needs to include native integration with the public cloud. You may tell yourself that your current legacy backup and recovery solution does fine providing support for the public cloud. Perhaps you've purchased a license or appliance that bolts public cloud functionality onto the product or added a gateway appliance that adds cloud support.

Beware, though: these kinds of bolt-on accessories are often incomplete, expensive, and add significant complexity to the backup and recovery environment. Given the importance of the cloud infrastructure, can an afterthought solution really address your evolving business needs? Unlike these legacy backup solutions, a modern backup and recovery solution eliminates the need for bolt-on cloud gateways and natively supports cloud integration.

But it's about more than just the underlying hardware and software combination. To achieve an optimal balance between cost and recovery metrics such as recovery point objectives (RPO) and recovery time objectives (RTO), you need to be able to store some protected data on-premises and put the rest into an inexpensive cloud repository.

A modern backup solution allows enterprise IT owners to take advantage of policy-based automation to seamlessly move their data between on-premises and cloud infrastructure for tiered long-term archival.

Recovery is all about getting workloads back into an operational state following an incident of some kind. An incident can be as simple as an accidentally deleted file, or as severe as the complete destruction of your data center. In the case of the former, recovery isn't that painful. For the latter, however, recovery can be incredibly challenging without the right tools in place.

To that end, the backup and recovery product you choose should enable provisioning of workloads to either the original cluster or an entirely new location, be it in the cloud, in a secondary data center, or in a collocated facility across town. In this way, immediately following a disaster, you'll be able to spin up your workloads while you work on your longer-term recovery efforts.

Questions to ask vendors about their cloud abilities:

- Does your solution provide integration with public cloud providers?
- Is that integration provided via an additional product, or is it a native and integral part of the platform?
- What options for leveraging the cloud are available in the solution? For example:
 - Can you choose to use the cloud as just an archive?
 - Can you choose to use the cloud as an active storage tier with the backup and recovery tools managing where data resides?
 - Can you use the cloud as a replication target to make it easier to use the cloud as a disaster recovery site?
- Can your solution truly offer free and flexible data mobility across on-premises and cloud environments?

Limitless, Transparent Scale-out

Just like your production environment, your backup and recovery environment isn't static. It's constantly changing in lockstep with the changes you make in your production environment.

Or at least it should be. For many, the backup and recovery environment is all too often inflexible. This inflexibility leads to increased costs and increased delays in deploying new production systems, as IT operators are forced to contend with adjusting the backup and recovery environment to accommodate new workloads.

In an era where speed to market makes or breaks the bottom line, it's not acceptable for the backup and recovery environment to hold a business hostage. It's even more intolerable when you consider the fact that hyperconverged solutions exist on the market. These kinds of solutions are purpose-built to make scaling the environment easy.

This type of scale-out to address business needs should also be true of modern backup and recovery solutions. With legacy backup and recovery solutions, you're forced to make critical sizing decisions at the start of your deployment, meaning that you have to attempt to predict your backup and recovery capacity needs three to five years in advance. Who can do that with any real degree of confidence (and potentially block valuable CapEx in the process)?

Your backup and recovery environment should be designed on web-scale principles similar to Google and Facebook. Most importantly, the solution you choose should allow you to start at the beginning and scale as your needs dictate. That means that you should be able to start small, without the need to overprovision infrastructure from the beginning.

Then, as your business grows, your production environment will also grow, which will necessitate growth in your backup and recovery environment. Make sure you can scale that environment on demand to address growing business requirements. A modern backup solution limitlessly scales linearly, without any impact on workload performance.

Questions to ask vendors about their flexible scalability:

- How do you grow the environment as operational needs dictate?
- Does the product use a scale up or scale out expansion methodology?
- To what level can you scale the backup and recovery environment?
- Will we experience any performance impact at scale?

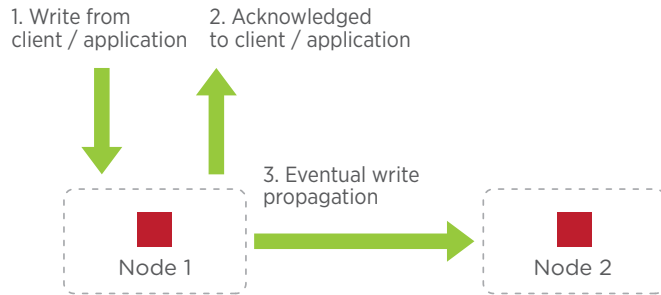
Guaranteed Data Resiliency

Organizations are not just backing up their data once a day anymore. Modern applications are accessing and modifying data continuously, so organizations back up their data more often.

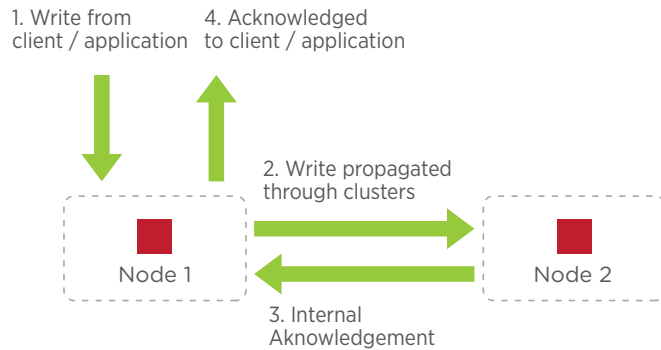
During a backup job, if the target device fails, most legacy backup solutions have the intelligence to continue the backup on another target device from the failed point; but that process isn't as smooth if it involves third-party backup software. Take Oracle RMAN, for instance; in that case, a failed node requires the administrator to restart the entire backup job.

This need results in extended backup windows, and can potentially bleed into production time, impacting RPOs due to the eventually consistent data state the legacy solution supports. The impact of eventual consistency can be far more significant if the node fails during a recovery process. While the data is being recovered, some applications/users might be writing to the same data volume, and a node failure at that point can result in inconsistent and corrupt data.

Eventual Consistency



- System eventually returns latest write
- Potential for data loss if node fails



- System always returns latest write
- Guaranteed data resiliency

To meet business SLAs, avoid “file not found” errors, and achieve guaranteed data resiliency, your next backup and recovery solution needs to be strictly consistent. In a strictly consistent model, the application or client only receives an acknowledgement of the write once the data propagates to multiple nodes. Thus, in the event the ingest node fails, the application or client won't receive a false acknowledgment that its data is protected or written; this results in keeping the data strictly consistent. Insist on a strictly consistent model, whether you're backing up directly or using a third-party application like Oracle RMAN or running a recovery process.

What does “eventually consistent” mean?

Bear in mind that a modern data backup and recovery solution implements a scale-out storage foundation, which enables the environment to grow as the needs of the business dictate. Any hardware environment is subject to hardware failure. Imagine a case in which you're actively writing data to a node in your backup cluster, and that node fails. How do you protect yourself? In general, scale-out storage systems write more than one copy of the data, with each copy going to different nodes. In this way, a hardware failure doesn't result in a loss of data, right?

It depends on the way that the vendor has chosen to implement their storage layer. Here's an example:

- An application writes data to node 1 (step 1)
- Node 1 acknowledges the write operation to the application (step 2)
- Node 1 does not have a chance to replicate that write operation to node 2 (step 3) since this solution uses an eventually consistent scheme.
- Node 1 fails and the data is not recoverable

The result is that there is data loss.

Now, consider this example of a strictly consistent operation:

- An application writes data to node 1 (step 1)
- Node 1 replicates that data to node 2 and acknowledges to node 1 that the operation is complete post internal acknowledgment (step 2)
- Node 1 acknowledges the write operation to the application (step 3 and 4)
- Node 1 fails and the data lives safely on node 2

Questions to ask vendors about data resiliency:

- Does your solution offer strict consistency?
- Can your solution help me meet my SLAs and guarantee data resiliency at scale?
- Does your solution only propagate the data, or does it also include the metadata?

Quick and Flexible Recovery for Any Size Environment

When (not *if*) you're in a recovery scenario, you'll discover that the solution you choose is only as good as the time it takes to recover the files, VMs, and other necessary elements you need to get back.

When disaster strikes, the administrator's primary goal is to locate and restore the production environment as quickly as possible: preferably, this means instantly.

To measure the effectiveness of backup and recovery tools, there are two metrics you need to understand:

- Recovery Point Objective (RPO): the amount of acceptable loss of data. A shorter RPO means less data is lost. For example, if you have a 5-minute RPO, you're saying it's acceptable to lose up to five minutes of data. To accomplish this RPO, you need a solution that ensures you're taking backups every five minutes.
- Recovery Time Objective (RTO): the amount of time it takes to recover your data once a loss occurs.

RTO and RPO are just the beginning



Often, backup and recovery tools focus on the simple, easy-to-grasp metrics of Recovery Time Objectives (RTO) and Recovery Point Objectives (RPO). RTO describes the amount of time your company is willing to be down during an incident that requires recovery, while RPO describes the tolerance you have for data loss. Both are typically measured in minutes or hours. Traditional thinking says that the closer you get to zero for these metrics, the more complicated and expensive the solution needs to be.

Although there is some truth to that, a modern backup and recovery system carries with it capabilities that make legacy backup and recovery products pale in comparison and that bring the potential for zero or close-to-zero RTO and RTO without legacy constraints and costs.

Modern backup and recovery software can provide low RPO and RTO metrics, thanks to the way data is protected under the covers. Choose a solution that combines fully-hydrated snapshots with short recovery points to ensure that you get the lowest RTOs possible.

But what if you need to go further back in time? Perhaps your most recent recovery point includes corrupted data. This situation isn't uncommon for companies that have fallen victim to ransomware attacks, for instance. A modern backup and recovery product allows you to recover data from any recovery point you need, not just the most recent one.

Better yet, a modern backup and recovery tool can provide instantaneous recovery times by provisioning clones from backup snapshots. So, rather than waiting for a complete recovery to take place before getting back in business, you can get a temporary clone operational while you await a recovery to the production environment.

Often, though, you need to be able to restore multiple services or VMs at the same time. Modern backup and recovery products enable you to instantly mass restore an unlimited number of VMs to any point in time.

Perhaps most importantly, you need to be able to identify the items that you'd like to recover. Modern backup and recovery tools fully index data being protected, making finding the files you need to recover as easy as performing a Google search.

Questions to ask vendors about their RPO and RTO abilities:

- What is the lowest RTO value offered by the platform, and what kind of environment do I need to build to enable a low RTO?
- What is the lowest RPO value offered by the platform, and what kind of environment do I need to build to enable a low RPO?
- Does the platform allow me to provision a clone from any historical snapshot I want instantaneously?
- What capability does the platform offer concerning identifying and locating files or elements that need to be recovered?
- Does the platform's indexing engine encompass both on-premises and cloud-based data to provide a complete restore function without having to access multiple consoles?
- How many virtual machines can the solution help me recover at one time? Do I need to break down my recovery job?

Maximize Storage Capacity and Efficiency

Data copies are exactly what they sound like: copies of your data that may exist in different parts of the environment. The problem of multiple copies has such an impact on businesses that an entirely new class of solutions were introduced to address it: data copy management. When you consider your backup environment, you may realize that there are dozens of copies of your data strewn about, particularly if you're using a relatively unintelligent backup product.

No matter how many copies of your data are in use or in your backup environment, the total capacity consumed by that data should not grow appreciably. Data deduplication is used to massively reduce the number of copies of your data so that your storage capacity goes much further than it otherwise would.

Most organizations today rely on deduplication appliances that create another silo and don't scale. In fact, most of these solutions only dedupe at the node level, which isn't efficient, given that the same data set might be stored in another node within the same cluster.

Modern backup solutions help storage administrators reduce their data center footprint with efficient data reduction techniques like global variable length dedupe. In this model, the deduplication is performed using variable-length data deduplication technology that spans an entire cluster and dedupes across all workloads, including data stored on physical systems, inside VMs, databases, and more. It also uses efficient protocols that result in significant savings across the storage footprint.

Improving on legacy offerings, modern backup solutions create variable length chunks of data, which optimize the level of deduplication no matter the type of data. In addition to providing global data deduplication, these solutions provide the option to dedupe inline, post-process or not at all. This flexibility allows you to choose the appropriate option for your data.

Inline vs. post-process deduplication



- **Inline deduplication.** With inline deduplication, as data is written to the system—ingested—the deduplication algorithm is applied. This algorithm compares the incoming block to all of the other blocks in the deduplication table. If a match is found, the block is discarded and a pointer written in its place, thereby resulting in immediate space savings. Inline deduplication can sometimes impose a minor latency penalty, but with modern processors, this penalty is generally minimal and often goes completely unnoticed.
- **Post-process deduplication.** Under a post-process deduplication methodology, as data is ingested, it's immediately written in its raw, un-deduplicated (hydrated) format. On a schedule determined by the administrator, the deduplication process is initiated. This process goes through all of the data that has been written to storage since the previous deduplication process ran and removes any duplicate data blocks that have been written. Post-process deduplication tends to have a very minor performance benefit over inline deduplication, but this comes at the cost of storage capacity, at least on a temporary basis. While the data sits on storage in its hydrated form, it's consuming space that will be freed up after the deduplication process finishes up.

You should also look for more when it comes to data reduction. In addition to deduplication, data compression is an additional feature that, when combined with deduplication, creates an environment in which data capacity is maximized. Whereas deduplication removes duplicate blocks of data from your environment, compression typically operates on files and reduces the amount of space consumed by each one. There is often a lot of extraneous or duplicate information inside individual files that compression can squeeze out.

Questions to ask vendors about compression and deduplication:

- Does your solution provide truly global deduplication across all nodes, blocks, clusters, etc.?
- Does the solution allow to dedupe in the cloud as well?
- If deduplication is not global, to what degree is your solution able to dedupe at the node, block or cluster level?
- Can I dedupe across the entire cluster, including across multiple workloads, with your product?
- Can your solution dedupe data both at rest and data in flight?
- Is data compression also a part of the product's capabilities?

Online Upgrades and Expansion

There are a number of emerging trends that, taken together, make it more important than ever to keep your hardware and software current.

- The modern security environment is increasingly hostile, with malicious actors constantly probing defenses in search of weakness that can be exploited.
- The software-defined nature of today's infrastructure solutions means that vendors are continuously adding features and services to existing products.
- As more and more businesses turn to scale-out architectures, such as hyperconverged infrastructure, the ability to move to a just-in-time data center architecture mindset becomes more feasible.
- For backup and recovery environments — particularly including those that encompass other kinds of secondary data — maintaining 24/7 operations is no longer optional.

The old method of taking systems down to patch, upgrade, or augment them is no longer palatable to organizations. At the same time, it's important to keep up with the latest and greatest software. Modern backup and recovery solutions allow backup admins to upgrade their clusters without any downtime, through the use of rolling upgrades. A "rolling upgrade" means that each node in the cluster is upgraded individually, leaving all services operational on the remaining nodes for the duration of the maintenance.

This rolling upgrade paradigm also extends to eventual node replacement, as refresh cycles come and go. The concept of "forklift upgrades" doesn't exist in modern backup and recovery solutions. The administrators can introduce or retire nodes on the fly. Never again should you have to rip and replace a backup and recovery environment just because a depreciation schedule demands it.

Questions to ask vendors about upgrades:

- What process do you use to patch or upgrade the software components in the solution?
- Is there any downtime required to perform a software upgrade to any of the components that comprise the backup system?
 - If yes, how much downtime is necessary?
- Will your solution allow me to add or retire nodes without bringing the entire cluster offline?
 - If not, what process is necessary to add or retire nodes in the solution?

Extensible Platform to Consolidate Other Secondary Workloads

As you consider your next backup and recovery solution, think bigger than just data protection. Although it's critically important, it's just one aspect of your broader secondary data needs. Secondary data, as the name implies, isn't as mission critical as primary data, and doesn't have to support IOPS-hungry applications.

It's important to remember that not everything needs millions of IOPS and gigabytes of sustained throughput. All kinds of secondary applications are perfectly happy to run in a saner storage environment. For example, you may need simple file storage, or a place to stand up test/dev environments, or a place to deploy an analytics application that's capacity-hungry but doesn't tax storage performance.

Backup and recovery APIs in a software-defined world

Simplicity via a single console is a desirable element of any enterprise service today. However, for full integration into IT-wide workflows, all services in the data center should have the ability to be addressed via REST APIs. This type of integration makes it possible to automate vast swaths of your ongoing operations, increasing the overall efficiency of IT in the process.

For example, if you choose a backup and recovery product that includes a comprehensive API, you make it possible to include your backup and recovery environment as a full-fledged programmable entity that can be included in your software development processes. By doing so, you bring the power of data from your backup environment to accelerate your new application development. To bring your IT infrastructure under one roof, you need a backup and recovery solution that integrates with third party automation solutions to simplify management.



As you consider the broad scope of the applications you need to support, it becomes clear that many are better suited for a secondary storage environment, rather than trying to keep them on expensive primary storage.

One such application is backup and recovery. Viewing it as just another application to support requires a different mindset, but changing your thinking here has the potential to yield substantial results: both in terms of cost savings, and, beyond that, to operational improvements.

Modern backup and recovery solutions go well beyond backup functions, too. No longer are these environments just expensive insurance policies; they add tremendous value. Modern backup and recovery solutions consolidate other secondary workloads such as files and objects on the same platform that supports your data protection needs.

Once the modern platform has the data, interesting things happen. Consider the aforementioned test/dev environment: your modern backup and recovery solution can make your data productive by instantly provisioning clones for accelerating application development. The same improvements can be applied to analytics, enabling you to gain deep insight into the data without having to figure out where to store it all.

Questions to ask vendors about the platform extendability:

- What services beyond backup and recovery does your solution provide?
- How can I use your solution to accelerate my application development?
- Does your solution allow me to directly run virtual machines from the production environment in the event of a failure in the primary storage environment?
- For any capabilities beyond backup and recovery, is there additional hardware or licensing needed to unlock these features?

Security and Compliance

Previously, security was mentioned in the context of software updates, but that's just one very small portion of the security story in your backup and recovery environment. The product you choose needs to take an as-designed, built-in approach to security that ensures that your company stays out of the news, out of court, and pleases your customers.

Encryption is considered a gold standard for security. There are multiple points at which data can be encrypted on most platforms. Most commonly, data is encrypted at-rest. That is, while data's sitting in storage and not being used, it's being kept safe from prying eyes.

But that single approach ignores the fact that data is actually used from time to time. After all, at some point, your CEO is going to open that financial report that's been sitting on storage in the data center. In some systems, once it leaves the storage, it's no longer protected as it traverses the network. Your goal should be to secure a backup and recovery environment that provides AES-256 encryption, the strongest encryption commonly available. It should also protect that data while it's both at-rest and in-flight between systems.

Moreover, data should remain encrypted while it's being replicated to the cloud. There should be no point where a potential attacker can gain access to an unsecured communications point to get to your data.

And, while keeping data safe from attackers is important, there are people you actually want to be able to access it from time to time. Ensuring that the right people have the right level of access to the right data is just as important as locking people out; it may be even more important!

Today, role-based access control (RBAC) mechanisms that integrate with existing directories, such as Active Directory, are the leading ways to accomplish access goals. Your next modern backup and recovery product should provide RBAC and allow you to grant access to subsets of your secondary data through different parameters.

Of course, when it comes to security, keeping the bad guys out is just part of the equation. Along with security often comes compliance. Today, compliance in backup and recovery products is becoming progressively more important, particularly as we continue down the path of cloud in an increasingly regulated world.

Today, backup administrators have to worry about such topics as data sovereignty and the EU's General Data Protection Regulation (GDPR). Data sovereignty is a regulatory requirement for some types of data, stipulating that certain kinds of data cannot leave a country's borders. In a cloud-centric world, it's often difficult to figure out where data is actually residing. GDPR is an overarching set of data protection requirements that aims to bring control of personal data back to the people; it also carries with it stiff penalties for companies that fail to follow its framework.

Finally, no security discussion around data protection would be complete without recognizing the growing threat of ransomware. This internet scourge is proving highly dangerous for organizations around the globe, as it exploits weaknesses in both technological and human-based systems. Modern data protection services need to put ransomware protection – to include detection, deletion, and recovery – as a front-and-center capability.

Questions to ask vendors about security:

- Does the solution provide data encryption at-rest?
- Does the solution provide data encryption in-flight?
- Does data remain encrypted while it's being replicated to the cloud, and does it stay encrypted once there?
- Does the solution offer Active Directory-integrated and other types of role-based access control?
- Describe in detail the capabilities that your solution offers around protection against ransomware.
- Provide details around how your solution helps my company adhere to data sovereignty and GDPR requirements.

What Modern Backup and Recovery Looks Like

As you journey down the road toward implementation of a modern backup and recovery solution, it becomes apparent rather quickly that the solution landscape in this area has changed in a dramatic way in recent years.

No longer are you confined to thinking of backup and recovery as a silo, but rather as a full-fledged member of your business-critical application set; and one that's self-healing, too. No longer do you have to dread the day that a ransomware attack strikes; you have a defense in place. No longer do you need to worry about the day when you max out your solution's capabilities and need to rip and replace it; modern solutions scale out, making expansion, upgrades, and node replacements a breeze.

But that's just the beginning. A modern backup and recovery product provides capabilities that are either impossible or prohibitively expensive with legacy approaches. Such capabilities include native cloud features, mass instant recovery of VMs, and a comprehensive API to enable advanced workflows. And it should be wrapped together in a simple, intuitive, single interface that manages your company's end-to-end backup and recovery services.

To learn much more about a leader in modern backup and recovery, visit Cohesity at

<https://www.cohesity.com/solutions/data-protection/>