



AWS Architecture Guide

Commvault® Version 11 SP4

VERSION 2.1 / JUNE 2016



Table of Contents

Abstract	7
The Cloud Difference	7
Infrastructure as Programmable, Addressable Resources.....	7
Global, Flexible and Unlimited Resources	8
Transforming The Disaster Recovery Model For A More Agile, Cost-Conscious Solution	8
Design Principles.....	9
Native Cloud Connectivity	9
Scalability.....	9
De-duplication Building Blocks.....	9
Client-side De-Duplication.....	9
Design For Recovery	10
Crash Consistency vs. Application Consistency	10
Storage-level Replication vs. Discrete Copies	10
Deciding What to Protect.....	10
Automation	11
Programmatic Data Management	11
Workload Auto-Detection and Auto-Protection	11
Self-Service Access and Restore.....	11
Cloud Use Cases with Commvault® Software	12
Backup/Archive to the Cloud.....	12
Disaster Recovery to the Cloud	13
Protection in the Cloud	14
Architecture Sizing.....	15
Amazon	15
AWS CommServe® Specifications	15
AWS Media Agent Specifications	15
Architecture Considerations.....	16
Networking.....	16
Virtual Private Cloud.....	16

Bridging On-Premise Infrastructure – VPN & DirectConnect	16
Infrastructure Access.....	17
Hypervisor access in Public Cloud	17
Amazon VPC Endpoints	17
Data Security	17
In-flight	17
At-rest.....	18
HTTPS Proxies.....	18
Data Seeding	18
“Over-the-wire”	18
Drive Shipping	18
Consumption / Cost.....	19
Network Egress	19
Storage I/O	19
Data Recall.....	19
Performance / Storage	20
Multi-Streaming with Object Storage	20
Cloud Connector Best Practices	20
Compression vs. De-duplication	20
Block Storage vs. Object Storage	20
Micro Pruning	21
Selecting the right Storage Class.....	21
Performing Disaster Recovery to the Cloud.....	23
Amazon	23
Restoring Applications (Automated or On-Demand).....	23
Virtual Machine Recovery into AWS EC2 Instances	23
Replicating Active Workloads.....	24
Using Workflow to Automate DR	24
Protecting and Recovering Active Workloads in the Cloud	25
Amazon	25
Agent-less EC2 Instance Protection (Virtual Server Agent for AWS)	25
Agent-In-Guest.....	28

Continuous Data Replicator (CDR)	29
Deployment	30
Remote Access / Bring Your Own Software	30
Installation Basics	30
CommServe® Disaster Recovery Solution Comparison	30
Pre-packaging Commvault® Software within a VM Template	31
Automating Deployment with Continuous Delivery	31
Cloud Library Configuration	32
Additional Resources	33
Documentation	33
Books Online – Cloud Storage	33
Videos	33
2 Clicks to the Cloud with AWS and Commvault	33
Appendix A: AWS Concepts and Terminology	34
Amazon (AWS)	34
Region	34
Availability Zone	34
Elastic Compute Cloud (EC2)	34
Instance Reservation	34
EBS-Optimized Instances	34
Elastic Block Storage (EBS)	35
EBS Snapshots	35
Simple Storage Service (S3)	35
Glacier	35

Notices

This document is provided for informational purposes only. It represents Commvault's current product offerings and practices as of the date of issue of this document, of which are subject to change without notice. The responsibilities and liabilities of Commvault to its customers are controlled by Commvault agreements, and this document is not part of, nor does it modify, any agreement between Commvault and its customers.

©1999-2016 Commvault Systems, Inc. All rights reserved. Commvault, Commvault and logo, the "C Hexagon" logo, Commvault Systems, Solving Forward, SIM, Singular Information Management, Simpana, Simpana OnePass, Commvault Galaxy, Unified Data Management, QiNetix, Quick Recovery, QR, CommNet, GridStor, Vault Tracker, InnerVault, QuickSnap, QSnap, Recovery Director, CommServe, CommCell, IntelliSnap, ROMS, Commvault Edge, and CommValue, are trademarks or registered trademarks of Commvault Systems, Inc. All other third party brands, products, service names, trademarks, or registered service marks are the property of and used to identify the products or services of their respective owners. All specifications are subject to change without notice.

Revision History

Version	Date	Changes
1.0	March 2015	<ul style="list-style-type: none"> Initial Version
1.1	May 2015	<ul style="list-style-type: none"> Updated AWS Architecture recommendations
1.2	June 2015	<ul style="list-style-type: none"> Added new Architecture Consideration sections - Networking (AWS VPC), Infrastructure Access, Performance / Storage Added new Installation sections - Bring Your Own License / Software sections (Installation), Video Tutorial links Added new Additional Architecture Resources section Updated document & layout for new Commvault branding Updated core cloud concepts and technology, AWS Sizing Recommendations and Security (Architecture Considerations) section Modified section layout Removed Data Aging caveats with SP11 Micro-pruning for Cloud Storage release, replaced text to refer to this only for pre-SP11 sites
1.3	July 2015	<ul style="list-style-type: none"> Updated with new trademark guidelines
1.4	August 2015	<ul style="list-style-type: none"> Minor reformatting Added new links to video content
1.5	September 2015	<ul style="list-style-type: none"> Added Selecting the right Storage Class section Minor reformatting
1.6	November 2015	<ul style="list-style-type: none"> New logo style Updated requirements for Disaster Recovery to the Cloud Updated cloud use case diagrams Updated VM Recovery to AWS feature Added Unsupported Cloud Configurations section
2.0	March 2016	<ul style="list-style-type: none"> Updated to reflect new Virtual Server Agent for AWS methodologies, deployment and changes to use cases Updated Backup to the Cloud, DR to the Cloud and Protection in the Cloud use case scenarios and requirements Updated Micro Pruning section Updated Drive Shipping to add note about Snowball support arriving 2016 Updated all BOL links to use Commvault Version 11 documentation Added new Documentation section to Additional Resources Added Automating Deployment with Puppet/Chef section Added Pre-packaging Commvault within a VM Template section Minor reformatting changes
2.1	June 2016	<ul style="list-style-type: none"> Updated Backup to the Cloud use case for more clear language around DDB requirements Added IntelliSnap functionality into VSA for AWS

Abstract

This document serves as an architecture guide for solutions architects and Commvault customers who are building Data Management solutions utilizing and the Commvault®Commvault's Cloud Solution sets.

It includes Cloud concepts, architectural considerations, and sizing recommendations to support Commvault®Commvault's Cloud Solution Sets. The approach defined in this guide extends existing functionality into easily sellable, re-usable architecture patterns to cover disaster recovery to the cloud as well as protecting running workloads in the cloud use cases.

The Cloud Difference

The Cloud megatrend is one of the most disruptive and challenging forces impacting customers' applications and infrastructure, requiring new business models and new architecture decisions, which impact how Commvault solutions protect and manage their data.

In general, Commvault believes the cloud contains these attributes that we can focus upon.

Infrastructure as Programmable, Addressable Resources

In a non-cloud environment: (i) infrastructure assets require manually configured, (ii) capacity requires manual tracking, (iii) capacity predictions are based on the guess of a theoretical maximum peak, and (iv) deployment can take weeks.

Within the cloud, these building blocks that represent the Infrastructure are not only provisioned as required, following actual demand and allowing pay-as-you-go, but can also be programmed and addressed by code. This greatly enhances flexibility for both Production/Dev/Test environments as well as Disaster Recovery scenarios.

Resources can be provisioned as temporary, disposable units, freeing users from the inflexibility and constraints of a fixed and finite IT infrastructure. Infrastructure can be automated through code, allowing for greater self-service and more automated delivery of desired business and technical outcomes. Consumption is measured by what you consume, not what you could consume, drastically changing the DR cost modelling challenges experienced today.

This represents a major, disruptive reset for the way in which you approach Disaster Recovery, testing, reliability and capacity planning.

Global, Flexible and Unlimited Resources

AWS offers global infrastructure available to Customers on a pay-as-you-go model, allowing for more flexibility in meeting requirements for Data Protection & Disaster Recovery.

Resources, bandwidth and their availability can now be localized to your corporate assets and human resources, allowing for a more distributed footprint that reduces backup windows and simplifies data protection that otherwise would be cost prohibitive with a physical datacenter or co-located approach, all while maintaining a simplified, unified pay-as-you-go billing approach.

Commvault® software is designed as a software-driven, hardware and cloud agnostic, highly modular, distributed solution that conforms with this new architecture reality, allowing Data Management solutions to be built to support and remain flexible with a highly distributed infrastructure built on-top of Cloud.

Transforming The Disaster Recovery Model For A More Agile, Cost-Conscious Solution

The cost model implications of Pay-as-you-Go don't just extend to Production workloads, but also to the ever present challenge of providing a flexible, agile yet capable DR solution for your applications.

Today, many physical DR environments have less capacity than their Production, or Dev/Test counterparts, resulting in degraded service in the event of a failover. Even more so, hardware is often re-purposed to fulfill the DR environment's requirements, resulting in higher than expected maintenance costs.

With the Public Cloud model, this hardware availability and refresh aspect is disrupted by removing the need to maintain a hardware fleet that can meet both your DR requirements and sustain your service level agreements.

You can provision instances to meet your needs, when you need them, and for specific DR events – both real and test – and the underpinning hardware is maintained and upgraded by the Cloud provider without any need for technical input, and no upgrade costs are incurred by the organization.

This dynamic shift allows you to begin costing per DR event, instead of paying for availability, improving your level of Disaster Recovery Preparedness through the application of flexible, unlimited resources to stage both DR tests and execute actual DR events – all without requiring pre-purchased hardware or disrupting production operations.

Design Principles

In this section, we provide design principles and architecture options for organizations planning to leverage the Cloud as part of their Data Management strategy.

Native Cloud Connectivity

The Cloud Connector is the native integration within the Media Agent module that directly communicates with AWS's S3 / S3-IA, without requiring translation devices, gateways, hardware appliances or VTLs.

This Connector works by communicating directly with Object Storage's REST API interface over HTTPS, allowing for Media Agent deployments on both Virtual and Physical compute layers to perform read/write operations against Cloud Storage targets, reducing the Data Management solution's TCO.

For more information on supported vendors, please refer to this comprehensive list:

[Cloud Storage - Support](#)

Scalability

Applications grow over time, and a Data Management solution needs to adapt with the change rate to protect the dataset quickly and efficiently, while maintaining an economy of scale that continues to generate business value out of that system.

Commvault addresses scalability in Cloud architecture by providing these key constructs:

De-duplication Building Blocks

Commvault software maintains a Building Block approach for protecting datasets, regardless of the origin or type of data. These blocks are sized based on the Front-End TB (FET), or the size of data they will ingest, pre-compression/de-duplication. This provides clear scale out and up guidelines for the capabilities and requirements for each Media Agent.

De-duplication Building Blocks can also be grouped together in a grid, providing further de-duplication scale, load balancing and redundancy across all nodes within the grid.

Client-side De-Duplication

As is the nature of de-duplication operations, each block must be hashed to determine if it is a duplicate block, or unique and then must be captured. While this is seen as a way to improve the ingest performance of the data mover (Media Agent), it has the secondary effect of reducing the network traffic stemming from each Client communicating through to the data mover.

In public cloud environments where network performance can vary, the use of Client-side De-Duplication can reduce backup windows and drive higher scale, freeing up bandwidth for both Production and Backup network traffic.

Design For Recovery

Using native cloud provider tools, such as creating a snapshot of a Cloud-based instance is easy to orchestrate, but does not deliver the application-consistency possibly required by a SQL or Oracle Database residing within the instance, and may even require additional scripting or manual handling to deliver a successful application recovery.

As part of any Data Management solution, it is important to ensure that you design for Recovery in order to maintain and honor the RPO and RTO requirements identified for your individual applications.

Crash Consistency vs. Application Consistency

While Crash-consistency within a recovery point may be sufficient for a file-based dataset or EC2 instance, it may not be appropriate for an Application such as Microsoft SQL, where the database instance needs to be quiesced to ensure the database is valid at time of backup. Commvault software supports both Crash and Application consistent backups, providing flexibility in your design.

Storage-level Replication vs. Discrete Copies

AWS offers replication at the Object Storage layer from one region to another, however, in the circumstance that bad or corrupted blocks are replicated to the secondary region, your recovery points are invalid.

While Commvault software can support a Replicated Cloud Library model, we recommend you configure Commvault software to create an independent copy of your data, whether to another region or cloud provider to address that risk. De-duplication is also vital as part of this process, as it means that Commvault software can minimize the cross-region/cross-provider copy by ensuring only the unique changed blocks are transferred over the wire.

Deciding What to Protect

Not all workloads within the Cloud need protection – for example, with micro services architectures, or any architecture that involve worker nodes that write out the valued data to an alternate source, mean that there is no value in protecting the worker nodes. Instead, the protection of the gold images and the output of those nodes provides the best value for the business.

Automation

The cloud encourages automation, not just because the infrastructure is programmable, but the benefits in having repeatable actions reduces operational overheads, bolsters resilience through known good configurations and allows for greater levels of scale. Commvault software provides this capability through three key tenants:

Programmatic Data Management

Commvault software provides a robust Application Programming Interface that allows for automated control over Deployment, Configuration, Backup and Restore activities within the solution.

Whether you are designing a continuous delivery model that requires automated deployment of applications, data collection and protection, or automating the refresh of a data warehouse or Dev/Test application that leverages data protection, Commvault software can provide the controls to reduce administrative overhead and integrate with your toolset of choice.

Workload Auto-Detection and Auto-Protection

The Commvault Intelligent Data Agents (iDA), whether the Virtual Server Agent for AWS, or the SQL Server iDA provide auto-detection capabilities to reduce administrative load.

Fresh instances, new volumes recently attached to a VM, or databases imported and created in a SQL instance are just examples of how Commvault software can automatically detect new datasets for inclusion in the next Data Protection window, all without manual intervention. Even agent-in-guest deployments can be auto-detected by Commvault and included in the next Data Protection schedule through intelligent Client Computer Groups.

This Auto-Detection and Auto-Protection level removes the requirement for a backup or cloud administrator to manually update the solution to protect the newly created datasets, improving your operational excellence and improving resiliency within your cloud infrastructure, ensuring new data is protected and recovery points maintained.

Self-Service Access and Restore

A common task performed by system administrators is facilitating access to recovery points for end-users and application owners, shifting their attention away from other day-to-day operations and strategic projects.

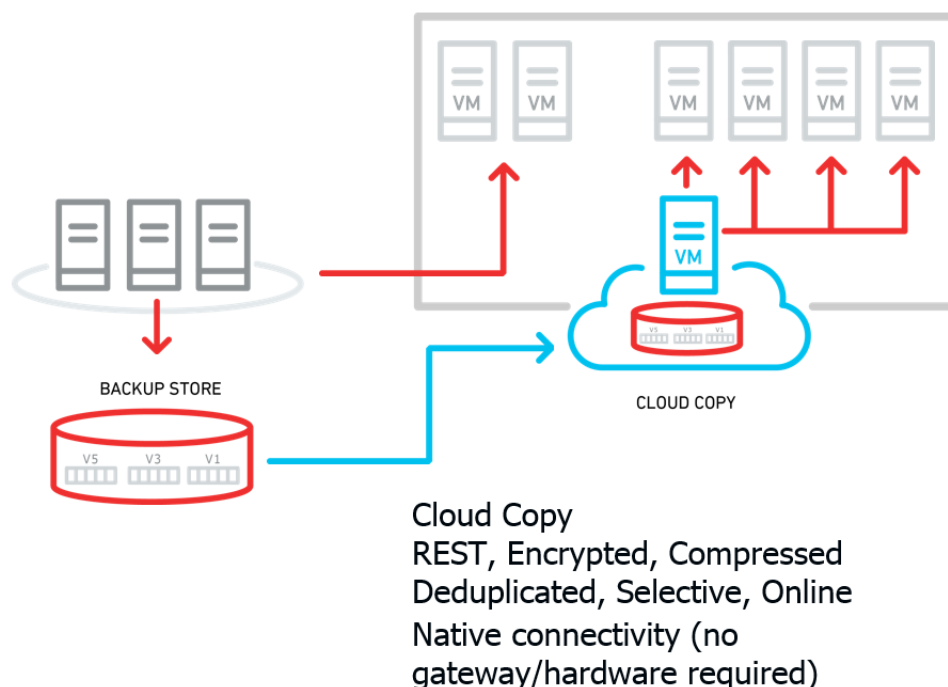
Commvault software's self-service interfaces empower users to access their datasets through a Web-based interface, allowing security mapped access to individual files & folders within the protected dataset, freeing up administrators to work on critical tasks.

Cloud Use Cases with Commvault® Software

There are three primary use cases when leveraging Commvault solutions with the cloud. These are backup to the cloud, DR in the cloud and protection for workloads running in the cloud.

Backup/Archive to the Cloud

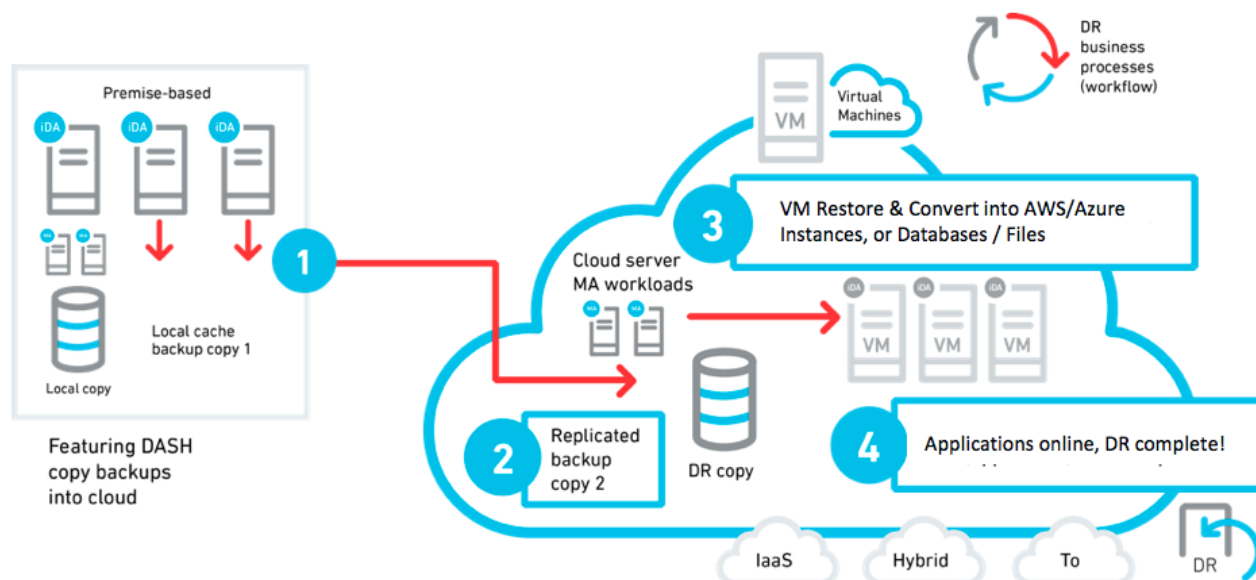
Protecting data at the primary on-premise location by writing directly to an external cloud provider's storage solution, or retaining a local copy and replicating the backup/archive data (either in full, or only selective portions of that data) into an external cloud provider's storage service.



Scenario / Suitability	Requirements
<ul style="list-style-type: none"> • Offsite Storage / “Tape Replacement” Scenario – no DR to the Cloud requirement, but can be extended if required • Native, Direct connectivity to 35+ Object Storage endpoints – no translation/gateway/hardware de-dupe devices required. • S3 or S3-IA Cloud/Object Storage target 	<ul style="list-style-type: none"> • Minimum 1x Media Agent On-Premise <i>No VM in Cloud required for B&R to the Cloud</i> • 1x DDB for the Cloud Library, hosted on On-Premise Media Agent. If a local copy is desirable, an additional DDB will be required. • Can use direct internet connection, or dedicated network to cloud provider for best performance (AWS Direct Connect)

Disaster Recovery to the Cloud

Providing operational recovery of primary site applications to a secondary site from an external cloud provider.

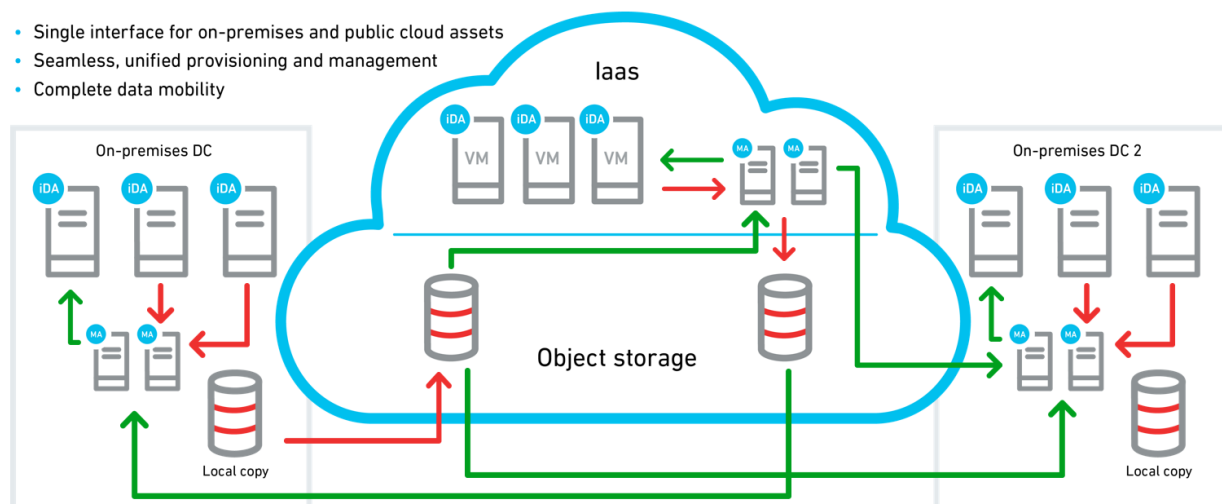


Scenario / Suitability	Requirements
<ul style="list-style-type: none"> • Off-site Storage Requirement & Cold DR Site in the Cloud – only use infrastructure when a DR event occurs, saving time & money (IaaS, DRaaS) • VM Restore & Convert – convert VMware and Hyper-V (Gen1) based Virtual Machines into AWS instances on-demand • Database/Files – restore out-of-place, whether on-demand or scheduled, to refresh DR targets • DR Runbook as Code – turn your DR runbook into a Workflow for easy simplified DR automation, whether for test or real DR scenarios 	<ul style="list-style-type: none"> • Minimum 1x Media Agent On-Premise, and minimum 1x Media Agent in Cloud • Media Agent in Cloud only needs to be powered on for Recovery operations • Highly Recommended to use dedicated network to cloud provider for best performance (AWS Direct Connect)

Protection in the Cloud

Providing operational recovery active workloads and data within an external provider's cloud.

- Single interface for on-premises and public cloud assets
- Seamless, unified provisioning and management
- Complete data mobility



Scenario / Suitability	Requirements
<ul style="list-style-type: none"> • Data Protection for Cloud-based Workload – protecting active workloads within an existing IaaS Cloud (Production, Dev/Test). • Agent-less Instance Protection – protect instances with an agent-less and automated/script-less protection mechanism through the Virtual Server Agent (AWS - requires v11 SP3 and newer) • DASH Copy to another Region, Cloud, or back to On-Premise – complete data mobility – replicate to another geographical region with IaaS provider, a different IaaS provider, or back to On-Premise sites 	<ul style="list-style-type: none"> • Virtual Server Agent and Media Agent deployed on a proxy within IaaS provider for agentless backup. Applications will require agent-in-guest deployed in each VM. (AWS – requires v11 SP3 and newer) • <i>(Applications requiring application-level consistency, and all other cloud providers)</i> Agents deployed in each VM within IaaS provider • Minimum 1x Media Agent in Cloud, and (optional) minimum 1x Media Agent at secondary site (whether cloud or On-Premise) • 1x DDB hosted on Media Agent • Recommended to use dedicated network from cloud provider to On-Premise for best performance when replicating back to On-Premise (AWS Direct Connect)

Architecture Sizing

Amazon

AWS CommServe® Specifications

Express / Workgroup	Data Center	Enterprise
<ul style="list-style-type: none"> • C4.2xlarge (8 vCPU, 15GB RAM) • 1x 100-150GB EBS “General Purpose” volume for CS Software & CSDB • Windows 2012 R2 	<ul style="list-style-type: none"> • C4.4xlarge VM instance (16 vCPU, 30GB RAM) • 1x 300GB EBS “General Purpose” volume for CS Software & CSDB • Windows 2012R2 	<ul style="list-style-type: none"> • C4.8xlarge instance • 1x 300GB EBS “General Purpose” volume for CS Software & CSDB • Windows 2012R2

AWS Media Agent Specifications

Express / 10TB FET	Data Center / 25-30TB FET	Enterprise / 60TB FET
<ul style="list-style-type: none"> • Up to 10TB estimated front end data • C4.2xlarge VM instance (EBS-optimized, 8x vCPU, 16GB RAM) • 1x 200GB EBS “General Purpose” volume for DDB • 1x 400GB EBS “General Purpose” volume for Index Cache • Linux or Windows 2012 R2 	<ul style="list-style-type: none"> • Up to 25-30 TB estimated front end data • C4.4xlarge VM instance (EBS-optimized, 16x vCPU, 30GB RAM) • 1x 600GB EBS “General Purpose” volume for DDB • 1x 1TB EBS “General Purpose” for Index Cache • Linux or Windows 2012R2 	<ul style="list-style-type: none"> • Up to 60TB estimated front end data • C4.4xlarge VM instance (EBS-optimized, 16x vCPU, 30GB RAM) • 1x 1TB EBS “Provisioned” volume for DDB @ 6500 IOPS • 1x 1TB EBS “General Purpose” for Index Cache • Linux or Windows 2012R2

Important: EBS-optimized instances are recommended as they provide dedicated network bandwidth for EBS volumes, improving De-duplication & Index Cache performance and freeing up bandwidth to send/receive from clients, other Media Agents & S3 endpoints

Bandwidth Considerations: Should additional network bandwidth be required on the Enterprise sizing, a C4.8xlarge instance can be used in-place of the C4.4xlarge

Architecture Considerations

Networking

Virtual Private Cloud

AWS has the capability to establish an isolated logical network, referred to as Virtual Private Cloud (VPC). Instances/Virtual Machines deployed within a VPC by default have no access to Public Internet, and utilize a subnet of the Customer's choice. Typically VPC's are used when creating a backbone between Virtual Machines, and also when establishing a dedicated network route from a Customer's existing on premise network via AWS Direct Connect.

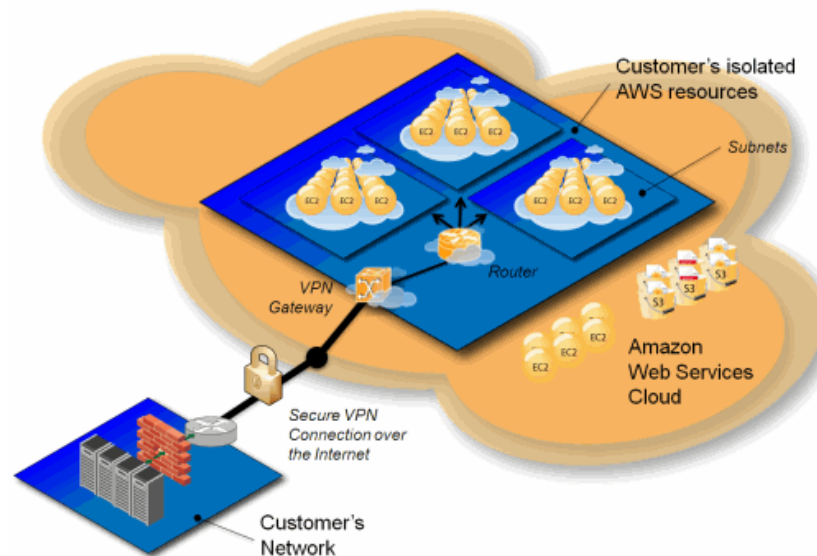


Figure 1 - AWS VPC Example Bridging On-Premise Infrastructure – VPN & DirectConnect

Customers may find a need to bridge their existing On-Premise infrastructure to their Public Cloud provider, or bridge systems and workloads running between AWS and on premise sites to ensure a common network layer between compute nodes and storage endpoints.

This is particularly relevant to solutions where you wish to Backup/Archive directly to the Cloud, or DASH Copy existing backup/archive data to Object Storage within a Cloud provider.

To provide this, there are two primary choices available:

- **VPN Connection** – network traffic is routed between network segments over Public Internet, encapsulated in a secure, encrypted tunnel over the Customer's existing Internet Connection. As the connection will be shared, bandwidth will be limited and regular data transfer fees apply as per the Customer's current contract with their ISP.

- **AWS Direct Connect** – a dedicated network link is provided at the Customer's edge network at an existing On-Premise location that provides secure routing into an AWS Virtual Private Cloud.

Typically, these links are cheaper when compared to a Customer's regular internet connection, as pricing is charged on a monthly dual-port fee, with all inbound and outbound data transfers included free of charge, with bandwidth from 10Mbit/s to 10Gbit/s.

Infrastructure Access

Hypervisor access in Public Cloud

AWS does not allow direct access to the underlying hypervisor, instead access to functionality such as VM power on/off, Console access are provided through an REST API.

Amazon VPC Endpoints

Amazon provides VPC Endpoints which enables you to create private connections between a given VPC and another AWS service without having to route via public Internet space. Support for S3 VPC Endpoints was announced May 2015, and while it is only supported within the same region as the VPC, it is highly recommended as it reduces availability risks and bandwidth constraints on the VPC's link through to public Internet.

An S3 VPC Endpoint must first be defined by creating an Endpoint Policy within the Amazon console, but there is no change to the FQDN hostname used to define the Cloud Library within Commvault. Instead, Amazon will ensure that DNS queries for the hostname will resolve against the S3 VPC Endpoint, instead of the public address, and apply appropriate routing (provided the Endpoint Policy was successfully created).

For more information on VPC Endpoints, please refer to AWS documentation:

[VPC Endpoints](#)

Data Security

In-flight

By default, all communication with Cloud Libraries utilize HTTPS which ensures that all traffic is encrypted while in-flight between the Media Agent and the Cloud Library end-point, but traffic between Commvault nodes is not encrypted by default. We recommend that any network communications between Commvault modules routing over public Internet space should be encrypted to ensure data security. This can be employed by using standard Commvault firewall configurations (Two-Way & One-Way).

At-rest

Data stored in a public Cloud is usually on shared infrastructure logically segmented to ensure security. Commvault recommends adding an extra layer of protection by encrypting all data at-rest. Most Cloud providers require that any seeded data be shipped in an encrypted format.

HTTPS Proxies

Please take note of any HTTP(S) proxies between Media Agents and endpoints, whether via public Internet or private space, as this may have a performance impact upon any backup/restore operations to/from an Object Storage endpoint. Where possible, Commvault software should be configured to have direct access to an Object Storage endpoint.

Data Seeding

Data Seeding is moving the initial set of data from its current location to a cloud provider in a method or process that is different from regular or normal operations. For seeding data to an external cloud provider, there are two primary methods:

“Over-the-wire”

Usually this is initially performed in small logical grouping of systems to maximize network utilization in order to more quickly complete the data movement per system. Some organizations will purchase “burst” bandwidth from their network providers for the seeding process to expedite the transfer process.

Major cloud providers offer a direct network connection service option for dedicated network bandwidth from your site to their cloud such as *AWS Direct Connect*

Please see the chart below for payload transfer time for various data sizes and speeds.

Link Size	Data Set Size							
	1 GB	10 GB	100 GB	1 TB	10 TB	100 TB	1 PB	10 PB
10Mbit	1m 40s	2.2 hrs	22.2 hrs	9.2 days	92.6 days	-	-	-
100Mbit	1m 20s	13m 20s	2.2 hrs	22.2 hrs	9.2 days	92.6 days	-	-
1Gbit	8s	1m 20s	13m 20s	2.2 hrs	22.2 hrs	9.2 days	92.6 days	-
10Gbit	0.8s	8s	1m 20s	13m 20s	2.2 hrs	22.2 hrs	9.2 days	92.6 days

Drive Shipping

If the data set is too large to copy over the network then drive seeding maybe required. Drive seeding is coping the initial data set to external physical media and then shipping it directly to the external cloud provider for local data ingestion.

Please refer to the Books Online *Seeding the Cloud Library* procedure for more information:

http://documentation.commvault.com/commvault/v11/article?p=features/cloud_storage/cloud_library_seeding.htm

Note: Amazon Snowball support for seeding data transfers into AWS S3 is available with both Version 10 and Version 11.

Consumption / Cost

Network Egress

Moving data into a cloud provider in most cases is no cost, however moving data outside the cloud provider, virtual machine instance, or cloud provider region usually has a cost associated with it. Restoring data from the cloud provider to an external site or replicating data between provider regions are examples of activities that would be classified as Network Egress and usually have additional charges.

Storage I/O

The input and output operations to storage attached to the virtual machine instance. Cloud storage is usually metered with a fixed allowance included per month and per unit “overage” charges beyond the allowance. Frequent restores, active data, and active databases may go beyond a cloud provider’s Storage I/O monthly allowance, which would result in additional charges.

Data Recall

Low-cost cloud storage solutions may have a cost associated with accessing data or deleting data before an agreed upon time period. Storing infrequently accessed data on a low-cost cloud storage solution may be attractive upfront, however Commvault recommends modeling realistic data recall scenarios. In some cases, the data recall charges maybe more than the potential cost savings vs. an active cloud storage offering.

As a best practice, Commvault recommends developing realistic use case scenarios and modeling cost against the identified scenarios to ensure the cloud solution will meet your organization’s SLAs as well as cost objectives.

Please see the links below for Amazon’s Cost Calculator here: [Amazon Cost Calculator](#)

Performance / Storage

Multi-Streaming with Object Storage

Object Storage performs best with concurrency, and as such with any Cloud Libraries configured within Commvault, performance will be best attained when configured for multiple readers / streams.

Cloud Connector Best Practices

There are additional Data Path settings and registry keys that can be modified to control the behavior of the Cloud Connector which will have an impact on the overall performance of the solution. For information on these settings/registry keys, please refer to *Cloud Connection Performance Tuning* within Books Online here:

[Cloud Connection Performance Tuning](#)

Compression vs. De-duplication

It is recommended that De-duplication should be used where possible, with the exception of environments where there are significant bandwidth concerns for re-baselining operations, or for Archive only use cases.

While additional compute resources are required to provide the necessary foundation for optimal De-duplication performance, using De-duplication even in a cloud context can still achieve greater than a 10:1 reduction.

Even with sealing of the DDB, reduction can be better than 7:1 reduction, providing significant network savings and reduced backup/replication windows (DASH Copy).

In comparison, Software Compression will achieve 2:1 reduction on average, and will constantly consume the same bandwidth when in-flight between endpoints (no DASH Copy).

Block Storage vs. Object Storage

While Public IaaS environments do allow for block-based storage to be provisioned and leveraged as Disk Libraries, the overall cost of those volumes can quickly exceed that of Object Storage. Based on AWS pricing June 2015, an internal case study showed that Object Storage could store 3x as much data as block-based storage (EBS “General Purpose”) for 33% less cost.

Additionally, with the inclusion of Micro Pruning in v10 SP11 for Object Storage, it is highly recommended that Object Storage be the primary choice for writing data to the Cloud, and other forms of storage by exception.

Micro Pruning

The Micro pruning support for Object Storage introduced in Version 10 SP11 is effective for new data written into the active store.

For customers who have upgraded from Version 10, but have not yet enabled micro pruning support, Macro pruning rules will still apply to existing data within the active store until the store has been sealed. But once the active store has been sealed, there will no longer be a need for continued periodic sealing against that store.

Micro pruning for Azure page blobs is supported, and does allow for more granular micro pruning within the stored objects, however it should be noted that Page Blobs incur a higher cost compared to standard Block Blobs and this cost should always be evaluated first.

Selecting the right Storage Class

Depending on the provider, there may be different tiers of Object Storage available which can significantly drive lower cost for your Cloud Architecture depending upon how you intend to access that data.

These tiers can be broken into three categories:

- **AWS S3 (Standard)** – this storage class represents the base offering of any Object Storage platform – inexpensive, instant access to storage on-demand at an avg. price of \$0.03/GB/month (as of October 2015 and depending on geographic region).

Typically it is expected that this tier would be used for Backup & Archive workloads in a short-term retention configuration.

- **AWS S3-IA (Infrequent Access)** – this is a relatively new offering that addresses what was a gap between Standard offering and Deep Archive storage tiers, in that it is offered at a lower price point than Standard storage (\$0.01 - \$0.012/GB/month) but is aimed at scenarios where data is infrequently accessed.

While the storage is always accessible, similar to the Standard offering, the cost model is structured to enforce an infrequent access use case by charging \$0.01/GB for any retrieval from this storage tier.

- **Glacier (Deep Archive)** – sometimes referred to as “cold storage”, this tier is intended for data that will probably not be accessed again, but must be retained in the event of compliance, legal action, or another business reason.

The cost of this storage class is the lowest compared to all three offerings – avg. \$0.007 to \$0.01/GB/month (as of October 2015, depending on geographic region) – but as with the

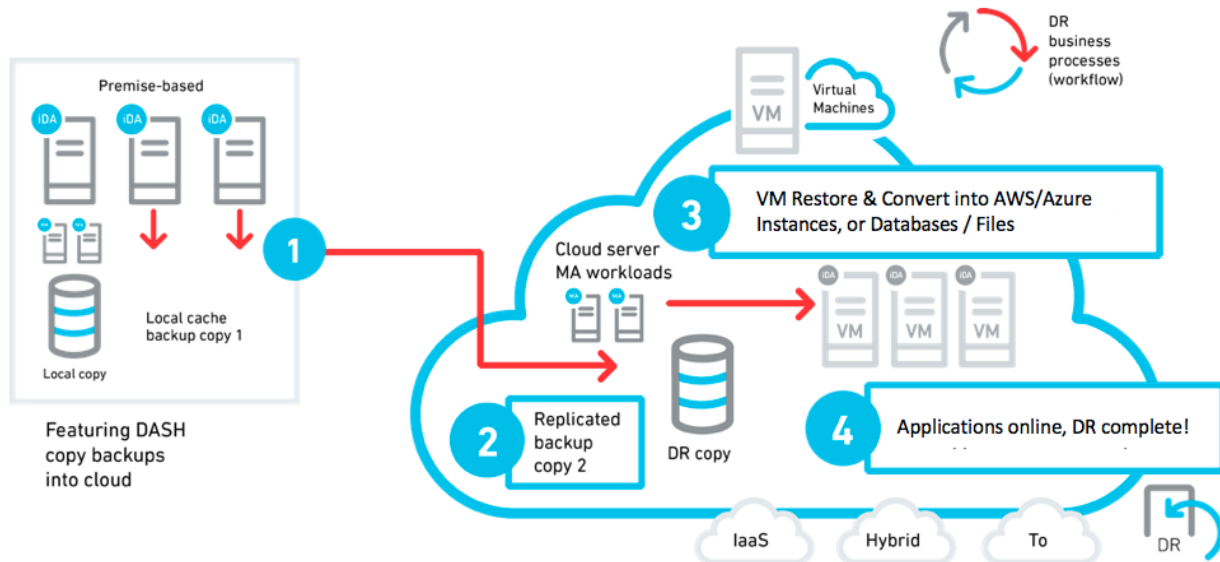
Infrequent Access class, the Deep Archive class's cost model is also structured with the expectation that retrievals are infrequent and unusual, and data will be stored for an extended period of time.

Additional charges apply if data is deleted prior to 30-90 days of an object's creation, and if more than 5% per month of your current managed data set in Glacier is retrieved, then additional costs may apply.

It is highly recommended that you review the cost options and considerations of each of these storage classes against the use case for your architecture in order to gain the best value for your cost model.

Performing Disaster Recovery to the Cloud

This section will cover the steps required to perform DR into the Amazon public cloud platforms. We will look at the recovery methods available for both image and agent based protection. This will also cover different recovery scenarios that may be needed to meet short recovery time objectives.



Amazon

Restoring Applications (Automated or On-Demand)

An agent in guest approach can be used to recover a wide variety of operating systems and applications. These can be captured at the primary site and replicated to the cloud based Media Agent in a de-duplication efficient manner. Once replicated the data can be held and restored in the event of a DR scenario or automatically recovered to existing instances for the more critical workloads.

Virtual Machine Recovery into AWS EC2 Instances

With the release of Version 11, the Commvault Virtual Server Agent allows for the ability to easily perform direct conversion of protected VMware or Hyper-V (Generation 1) virtual machines into AWS EC2 instances, from backups stored either within S3/S3-IA storage, another Cloud Library or from an on premise Disk Library.

This process could be used as part of a Disaster Recovery strategy using Amazon as a Cold DR site, or as a migration strategy (Lift-and-Shift).

Additional information on the Conversion feature can be located using the link below.

- [Converting Virtual Machines to Amazon \(from VMware\)](#)
- [Converting Virtual Machines to Amazon \(from Microsoft Hyper-V\)](#)

Replicating Active Workloads

Continuous Data Replicator (CDR) allows near time continuous data replication for critical workloads. These VMs will need a similarly sized EC2 instance running in AWS to receive any replicated data. In order for CDR to operate, an EC2 instance must be running at all times to receive application changes. Additional information on CDR can be located using the link below.

- [ContinuousDataReplicator \(CDR\)](#)

Using Workflow to Automate DR

The Commvault Workflow engine provides a framework in which the DR runbook process, covering the deployment of new instances, recovery of data and applications, and validation aspects of a DR operation can be automated to deliver a simplified, end-to-end GUI-driven DR process. This can be developed and maintained by your administrators, or with the assistance of Commvault's Personalization Services team.

For more information on Commvault's Personalization Services team, please contact your Account team.

For more information on the Workflow engine, please refer to the [Workflow Overview link](#).

Protecting and Recovering Active Workloads in the Cloud

This section will cover the basics on protecting active workloads running on both the Amazon and Azure public cloud offerings. Included will be the various protection approaches as well as replication and recovery to different geographic regions. We will also touch on cross platform recovery as well as recovery to onsite locations.

Amazon

Agent-less EC2 Instance Protection (Virtual Server Agent for AWS)

Introduced in Version 11 Service Pack 3, the Virtual Server Agent for AWS (VSA for AWS) delivers an agent-less, block-level capture of EC2 instances and their attached EBS volumes. Restoration options include both Full Virtual Machine recovery and limited Granular-level file recovery.

When to use the VSA for AWS

- *Agent-less protection approach for EC2 instances & file-level data* – no agents are required in-guest to perform a block-level backup to provide Instance and File-level recovery

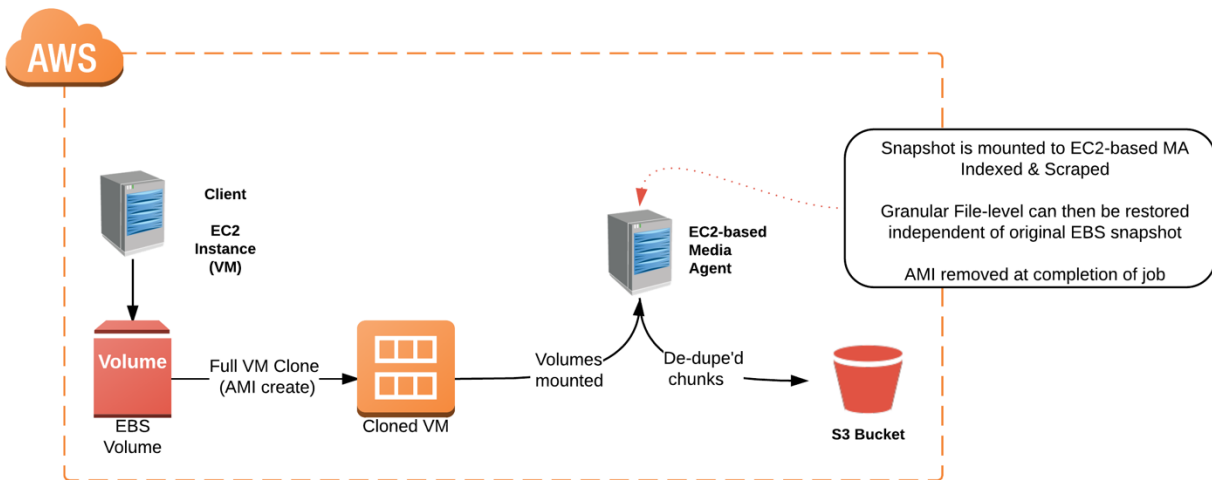
When not to use the VSA for AWS

- *When you require application-consistent backups* – the VSA for AWS approach creates a Crash-consistent image of the source EC2 instance and its EBS volumes. If you require application consistency, use an agent-in-guest either standalone or in conjunction with the VSA for AWS backup schedule.
- *Protecting worker/stateless instances* – Worker nodes may generate valued data that is moved to another centralized repository, and the nodes themselves do not require protection. It is recommended to instead target that centralized repository for Data Protection instead of the individual worker nodes, whether with VSA for AWS or agent-in-guest, depending on the required level of backup (crash vs. application consistent)

How Instances are qualified for protection

- Each VSA can be configured with 1 or more subclients. Each subclient defines a rule set on which to Auto-Detect and Protect EC2 instances, based on a user-defined criteria of Instance Name, Region, Availability Zone or Tags.
- During the Discovery phase of the backup job, the VSA will use the subclient's rule to qualify instances to add/remove for protection within that job.

The VSA for AWS is an Early Release feature, meaning that the feature is available for use in environments that meet all the necessary requirements validated by Commvault. For more information on Early Release products, requirements and status, please e-mail EarlyRelease@commvault.com



Commvault software does not require access to the AWS hypervisor-level, instead using the EC2/EBS REST APIs to create a Full VM Clone (AMI creation) of each EC2 instance, attaching the EBS volumes to a nominated proxy (EC2-based VSA / Media Agent) to read and de-duplicate the blocks before writing out to a S3 Bucket.

IntelliSnap® functionality introduced in Version 11 Service Pack 4 for the VSA for AWS modifies this behavior by simply creating an EBS snapshot and retaining the snapshot based on the Storage Policy's Snap Primary retention setting.

This has the effect of reducing backup windows considerably, providing fast, snapshot-based restoration capability and offloading the task of extracting blocks from the snapshot into S3 for longer term retention through the Backup Copy operation.

On snapshot cost: Use of the IntelliSnap method does mean that snapshots remain available for longer than the backup window, however we recommend that snapshots be retained only for as long as required. The Snap Primary can be configured to retain at least 1 snapshot, keeping snapshot costs at a minimum while providing fast backup and restoration capabilities.

Architecture Requirements for the VSA for AWS

- **Minimum 1x VSA/MA per Region,
Recommended 1x VSA/MA per Availability Zone**
 - Each 1x “VSA/MA” node represents a single Windows 2012 R2 EC2 instance with the Virtual Server Agent and Media Agent modules deployed. The EC2 instance specifications should match the Media Agent specifications within this Architecture Guide.
 - If the environment contains more than >25 VM’s within a single Availability Zone, it is recommended to scale out with additional VSA proxies (1 or more) placed inside of the source AZ to improve performance and reduce cross-zone traffic.

Architecture Recommendations

- Use of the IntelliSnap configuration is highly recommended (requires Version 11 Service Pack 4) to improve backup & restore times. Use of this method does mean that snapshots remain available for longer than the backup window, however we recommend that snapshots be retained only for as long as required. The Snap Primary can be configured to retain at least 1 snapshot, keeping snapshot costs at a minimum while providing fast backup and restoration capabilities.
- Use of the S3 VPC Endpoint is highly recommended to improve throughput to/from S3 buckets.
- Configuring more than 10 readers for a single VSA for AWS proxy may cause snapshot mount operations to fail. Consider scaling out with multiple VSA proxies if higher concurrency is required.
- Use of Provisioned IOPS for the De-Duplication Database used by the Media Agent module is highly recommended for optimal performance.
- Disable Granular Recovery of files if granular recovery is not required, or agents-in-guest are being used to collect large file system datasets. This will improve the backup window by removing the need to ‘walk’ the file system structure within the EBS volumes.

Notes on VSA for AWS Granular-level File Recovery: Granular-level file recovery is currently limited to NTFS (Windows), and ext2/3 file systems. Restoration of granular-files to an EC2 instances requires the deployment of a File System Agent (Recovery Only) within the EC2 instance. Agent-less file recovery in AWS is not supported due to API limitations.

Agent-In-Guest

An agent in guest approach can be used to protect a wide variety of operating systems and applications. These can be captured on the production workload and protected to the Media Agent residing in AWS, using Client-side De-Duplication to reduce the network consumption within the Cloud. These can also be replicated to a secondary Media Agent residing in a different geographic region. Once replicated the data can be held and restored in the event of a DR scenario or automatically recovered to existing instances for the more critical workloads.

When to use Agent-in-Guest approach:

- *When you require application-consistent backups* – use an agent-in-guest either standalone or in conjunction with the VSA for AWS backup schedule. Deployment of agents can either be pushed by Commvault software, baked into AMI templates using decoupled installation, or deployed as part of a continuous deployment method (ie. Puppet/Chef/Ansible).
- *When you require granular-level protection and restoration features for applications* – the Commvault iDataAgents can deliver granular-level protection for supported application workloads, such as SQL Server or Oracle Database, in comparison to a Full VM or File-level approach.

Architecture Requirements for Agent-in-Guest:

- Minimum 1x iDataAgent per Instance for the intended dataset (ie. SQL, File). Multiple iDataAgents can be deployed on the same instance.
- Minimum 1x Media Agent per region. Media Agents connect to the target Object Storage, and can either be deployed on the same instance, or on a dedicated host for a fan-in configuration. The EC2 instance specifications of the Media Agent should match the Media Agent specifications within this Architecture Guide.
- Check the *Systems Requirements* section in Books Online to determine if the iDataAgent supports your application (<http://documentation.commvault.com>)

Architecture Recommendations

- Use of multiple readers to increase concurrency to the Object Storage target is recommended
- Use of the S3 VPC Endpoint is highly recommended to improve throughput to/from S3 buckets

Continuous Data Replicator (CDR)

CDR allows near time continuous data replication for critical workloads. Replication can be configured as Direct Replication (1:1 source to destination host), or as a Fan-in or Fan-out based replication configuration.

Additional information on CDR can be located using the link below.

- [ContinuousDataReplicator \(CDR\)](#)

When to use CDR approach:

- *Fan-in replication of File data-sets from Remote regions to a Cloud-based host* – CDR supports a fan-in based replication approach (Many-to-1) in which data from remote sites can be replicated in to a target Cloud-based host
- *Guest-level Replication of File and SQL/Exchange data-sets* – CDR can be used to perform asynchronous block-based replication of File, SQL and Exchange datasets between hosts, irrespective of the source/destination hardware/hypervisor and storage. This allows data to be kept in-sync between an on-premise host and a cloud-based host.

When not to use CDR approach:

- *Replicating VM's/Hosts* – CDR is intended as a host-based data replication feature, and it should not be used to attempt to replicate Virtual Machines or System volumes as a replacement Bare Metal Recovery strategy.

Architecture Requirements for CDR:

- Storage for the replication journal is required for each source and destination host – ensure that this volume is available and is sized to support both the daily change rate, and sufficient storage in the event of loss of network connectivity between hosts
- For Direct Replication configurations, a similar sized EC2 instance should be configured as the destination target.
- Check that the data set and Operating System is supported by the CDR agent¹

¹ ContinuousDataReplicator System Requirements - http://documentation.commvault.com/commvault/v11/article?p=system_requirements/flr.htm

Deployment

Remote Access / Bring Your Own Software

As with all IaaS offerings, remote access to Virtual Machine instances can be achieved with your favorite protocol / software (RDP for Windows, SSH for Linux instances) and Commvault module deployment can be achieved with the current procedures listed in Books Online.

Installation Basics

The following links cover the steps when installing the CommServe in the cloud. This is only needed when the primary CommServe will be running on the hosted cloud VM or used for DR recovery. Multiple modules can be deployed in a single installation pass to streamline deployment.

- **Installation Overview -**
http://documentation.commvault.com/commvault/v11/article?p=deployment/common_install/c_install_overview.htm
- **Installing the CommServe -**
http://documentation.commvault.com/commvault/v11/article?p=deployment/common_install/p_commservice_install.htm
- **Installing the Media Agent -**
http://documentation.commvault.com/commvault/v11/article?p=deployment/common_install/p_mediaagent_install.htm
- **Installing the Virtual Server Agent (Amazon) -**
http://documentation.commvault.com/commvault/v11/article?p=products/vs_amazon/c_vamz_deployment.htm

CommServe® Disaster Recovery Solution Comparison

The following link covers CommServe® DR Solution comparisons for building a standby DR CommServe in the Cloud, or simply restoring on-demand (DR Backup restore):

[CommServe Disaster Recovery](#)

Pre-packaging Commvault® Software within a VM Template

For environments where deployment time is reduced by pre-preparing software and configuration within VM templates, such as Amazon Machine Images, the Commvault iDataAgents can also be deployed in *Decoupled mode*. This means that the iDataAgent is deployed within the instance, but will only be activated upon registration with the CommServe.

For more information, please refer to the *Installing the Custom Package* instructions within Books Online:

- **Installing the Custom Package on Windows -**
http://documentation.commvault.com/commvault/v11/article?p=deployment/install/custom_package/t_custom_package_install_win.htm
- **Installing the Custom Package on Linux -**
http://documentation.commvault.com/commvault/v11/article?p=deployment/install/custom_package/t_custom_package_install_unix.htm

Automating Deployment with Continuous Delivery

For environments using Continuous Delivery toolsets such as Puppet, Chef or Ansible, Commvault supports deployment methods that allow administrators to both control agent deployment and configuration to provide an automated deploy & protect outcome for applications and servers.

For more information on creating an unattended installation package for inclusion in a recipe, please refer to the Unattended Installation guide within Commvault Books Online:

- **Unattended Installation -**
http://documentation.commvault.com/commvault/v11/article?p=deployment/install/c_silent_install.htm

For more information on using Commvault software's XML / REST API interface to control configuration post-deployment, please refer to the Command Line – Overview section to review options available for each iDataAgent:

- **REST API – Overview -**
http://documentation.commvault.com/commvault/v11/article?p=features/rest_api/rest_api_overview.htm
- **Command Line – Overview -**
http://documentation.commvault.com/commvault/v11/article?p=features/cli/command_line_overview.htm

Cloud Library Configuration

This section covers the steps needed to configure cloud storage as a primary or secondary storage target. Please keep in mind that use cases outside of archive will require Commvault infrastructure in the cloud to recover any protected data.

For most backup use cases (except for very small environments limited to 100 GB in payload size), cloud as a direct storage target is not recommended. For performance and responsiveness, a primary copy should be stored on an on-site disk library and a secondary copy should be hosted on the cloud storage. The secondary copy should be setup as an encrypted network optimized DASH copy to the cloud.

The link below lists all of the supported direct cloud storage targets.

- [Supported Cloud Storage](#)

The link below covers cloud storage target setup and management.

- [Cloud Storage - Overview](#)

Details on performance tuning are covered below.

- [Cloud Connector Performance Tuning](#)

Additional Resources

Documentation

Books Online – Cloud Storage

http://documentation.commvault.com/commvault/v11/article?p=features/cloud_storage/cloud_storage_overview.htm

The *Cloud Storage* section from Commvault's Books Online documentation covers technical procedures and information on Supported Cloud Targets, Advanced procedures, Troubleshooting and FAQ sections for Commvault customers.

Videos

2 Clicks to the Cloud with AWS and Commvault

<https://www.youtube.com/watch?v=2IcygSNhzY0>

Focuses on creating an S3 container and configuring as a Cloud Library within Commvault v10. This video is applicable to both v10 and v11 environments. (An updated v11 video will be available soon.)

Appendix A: AWS Concepts and Terminology

The following section covers the basic concepts and technology used in the Amazon Web Services offerings.

Amazon (AWS)

Region

A Region is a separate geographic area where AWS services are offered. Services can be replicated between regions for geographic redundancy. Not all services are available in every region.

Availability Zone

An Availability Zone is an isolated service location within a region connected via low-latency links. Services can be replicated between Availability Zones for protection against single zone or datacenter failures.

[Region Product Page](#)

Elastic Compute Cloud (EC2)

EC2 is the product name of the virtual machine IaaS

Instance Reservation

By default, all instances requested are On-Demand, however if instances are reserved for a period of time then there can be significant discounts in consumption.

- **RESERVED** - Reserved instances is where your organization would commit to usage for a given time period (i.e. 1-3 years). Discounts can be significant depending on commitment type.
- **ON-DEMAND** - On-Demand instances is where your organization would have no commitment, can start or stop at any time, and pay a simple hourly rate.
- **SPOT** - Similar to On-demand, however your organization will have to “bid” for instances in the AWS market place. Your organization’s instances will keep on running until you stop the instance or the current spot price exceeds your organization’s bid price.

EBS-Optimized Instances

Specific EC2 instance types that have dedicated full-duplex network throughput for disk I/O.

[EC2 Product Page](#)

Elastic Block Storage (EBS)

EBS is the product name of the classic block storage service attached to EC2 instances, which traditional operating systems can lay a file system and use. EC2 instances utilize shared bandwidth network and storage operations. For dedicated bandwidth for storage I/O, select an EC2 that is EBS-Optimized.

- **GENERAL PURPOSE** - General purpose is a hybrid storage offering that utilizes traditional HDDs backed by SSD for performance. All new EBS storage provisioned with EC2 are general purpose by default.
- **PIOPS** - PIOPS is a pure SSD offering where your organization can pay for specific IOPS levels and latency.
- **MAGNETIC** - Magnetic is the legacy EBS offering that is comprised of traditional HDDs only, hence the name “magnetic”. Use of Magnetic EBS volumes should be discouraged.

EBS Snapshots

EBS snapshots are a point-in-time compact copy of data in an EBS volume, which can be access instantaneously, shared, or copied across regions.

[EBS Product Page](#)

Simple Storage Service (S3)

S3 is the product name for the object storage service. Object storage is fundamentally different than file or block storage. S3 storage can only be accessed through the REST/API from AWS.

- **BUCKET** - A bucket is a logical container of objects, where security permissions can be set and inherited by the child objects. You can't have objects in S3 without a Bucket created first.
- **OBJECT** - An Object is a reference to a single item of unstructured data such as file, backup, or anything else. The key term is single because you can't put two files in a single object. Each file would become an individual object that would be addressable separately including API access, permissions, and metadata.

[S3 Product Page](#)

Glacier

Glacier is the product name for AWS's cold storage service. In general, your organization's data will write to S3 and then will be “aged off” to Glacier by a data management policy set within S3.

- Recall is hours not minutes
- Recall can be expensive if your organization exceeds the allowed monthly allowance
- Deleting data before 3 months will incur additional charges

[Glacier Product Page](#)