

Reference Architecture: Cloud Volumes Service for AWS for High Performance Compute Workloads

Prabu Arjunan, NetApp
January 2020

Abstract

This architecture blog provides a brief overview of HPC and a reference architecture for a high-performance file system optimized for fast processing of workloads, such as video processing, financial modeling, and electronic design automation (EDA).

Table of Contents

Reference Architecture: Cloud Volumes Service for AWS for High Performance Compute Workloads.....	1
Table of Contents.....	2
Introduction	3
Customer Problem	3
<i>WuXi NextCODE:</i>	3
Customer Value of CVS.....	3
What are Some Other Key Values of CVS?.....	4
High performance.....	4
Increased resilience with snapshot copies	4
Speed up time to market: spin up cloud volumes in seconds with Instant Copy	4
Data durability	4
High availability	5
Security and encryption.....	5
Average cost savings of around 70%.....	5
General Architecture.....	6
HPC on AWS using CVS:.....	6
Benchmark Test with Comparison to Native File Services:	7
Benchmark Results	7
Proven Performance	9
The Cost Benefits of Cloud Volumes Service.....	9
Code Snippet	10
To list all cloud volumes and provision volumes for HPC workload	10
Customer Support for Cloud Volumes Service for AWS.....	11

Introduction

High-Performance Computing (HPC) drives the frontiers of science, finance, engineering, education, chemistry and genomics, and a host of other compute-intensive industry verticals. HPC workloads are mission-critical for most enterprises and feature prominently in discussions over enterprise cloud migration. The flexibility and heterogeneity of HPC cloud services provide a welcome contrast to the constraints of on-premises HPC.

This architecture piece provides a brief overview of high-performance computing workloads and gives a reference architecture for HPC workload on Cloud Volumes Service (CVS) and AWS. It also explains the benefits of running HPC workload on CVS and AWS.

Customer Problem

WuXi NextCODE:

WuXi NextCODE is focused on using genomics to identify underlying biology of disease and to advance the understanding of diseases to propel the next generation of transformative therapies.

Although all humans share a similar DNA sequence, it is not 100% unique to the individual. If you and a friend were to compare your DNA, you would find that in the roughly 3 billion letters of the DNA, you differ in about 5 million locations. WuXi NextCODE, a genomics information company that uses genomic data to improve healthcare around the globe, is a groundbreaking user of Cloud Volumes Service for AWS. Their work is immediate and requires enormous processing power to complete. According to Dr. Gudbjartsson, PhD, WuXi NextCODE, “the challenge is to take a dataset of 5 million and figure out the differences or mutations that are important—which ones are the causes of rare diseases, which ones are the causes of cancer, and how to treat patients”. At the heart of the WuXi NextCODE platform is the genomic relational database, the only relational database designed and optimized to query and analyze massive genomic data. When they attempted to use cloud-based NFS servers and cloud-native file share services to run data sets, they experienced timeouts or file failures.

Customer Value of CVS

By taking advantage of NetApp® Cloud Volumes Service, the genomics platform makes it possible to integrate data dynamically to deliver unprecedented computational efficiency.

“A benchmark analysis for analyzing genomic data is generating the allele frequency of every mutation found in a population of 100,000 individuals,” Dr. Hákon Gudbjartsson, PhD, the CIO of WuXi NextCODE, said. “With earlier storage solutions (or self-managed storage), we always had timeouts or file failures. But when we tested this using the NetApp Cloud Volumes Service, it actually finished in less than an hour. That was a great breakthrough for us.”

WuXi NextCODE found that only Cloud Volumes Service delivered the performance they needed to optimize their use of massive genomic data. With Cloud Volumes Service, they can rapidly integrate data on-the-fly to deliver unprecedented computational efficiency to their customers. With Cloud Volumes Service for AWS, you can run a high-performance workload with maximum data protection. Underlying

that security is NetApp® Snapshot™ technology, which offers a crucial option for rapid, efficient backup and restore. By design, CVS for AWS provides nine 9s (9.999999999%) of data durability.

What are Some Other Key Values of CVS?

High performance

With consistently high performance of **over 450 MB/sec**, Cloud Volumes Service provides shared persistent storage with high throughput and low latency. It easily meets the demands of large HPC workloads, with [SLAs that guarantee performance](#).

Increased resilience with snapshot copies

You can easily create a snapshot of an HPC database using NetApp® Snapshot™ technology.

Snapshots act as logical backups. They're point-in-time representations of your data, with a rapid revert function that allows you to restore your database without downtime. You create snapshots manually or schedule the creation of snapshots using the Cloud Volumes Service API or graphical user interface (GUI); rapid revert is only available from the API.

Snapshots are fast, plentiful, and nondisruptive. A snapshot simply manipulates block pointers, creating a "frozen" read-only view of a volume that enables your applications to access older versions of files and directory hierarchies without special programming. They do not make full copies, but rather record new writes (incremental). Snapshot creation takes only a few seconds (typically less than 1 second) regardless of the size of the volume or the level of activity within the environment. Since they are read-only, incremental copies, you only pay for the space consumed by new data written.

Speed up time to market: spin up cloud volumes in seconds with Instant Copy

Most organizations need multiple copies of data for testing and development. HPC landscapes are littered with system copies for variety of uses; creating and refreshing those copies is cumbersome. Typically, creating copies of HPC landscapes is a time-consuming and tedious process. Cloud Volumes Service for AWS allows you to instantly copy volumes, drastically improving the process of copying, backing up, and reverting. The process takes seconds, which ultimately leads to quicker time to market.

Data durability

With Cloud Volumes Service, data is protected not just against multiple drive failures, but also against numerous storage media errors that can harm your data durability and your data integrity. And with 99.9999999% durability—based on the experience of over 300,000 customers—you don't have to worry that your data is going to disappear, which is underpinned by [the product's SLA](#).

High availability

Built on industry leading hardware and software, NetApp Cloud Volumes Service offers four 9s (99.99%) of availability enabled by architectural features, such as redundant network paths, failover, and advanced data protection.

Because NetApp Cloud Volumes Service for AWS sits centrally in relation to each of the Availability Zones within an Amazon Web Services (AWS) region, your service is unaffected by Availability Zone outages. You can access your data from any Availability Zone within the region without having to replicate content. This availability is covered by [CVS's SLA](#).

Security and encryption

NetApp Cloud Volumes Service uses at-rest encryption, relying on the XTS-AES 256-bit encryption algorithm. CVS encrypts your data without compromising your storage application performance. NetApp manages and rotates encryption keys for you, thus, this single-source solution can increase your organization's overall compliance with industry and government regulations without compromising your user experience.

Average cost savings of around 70%

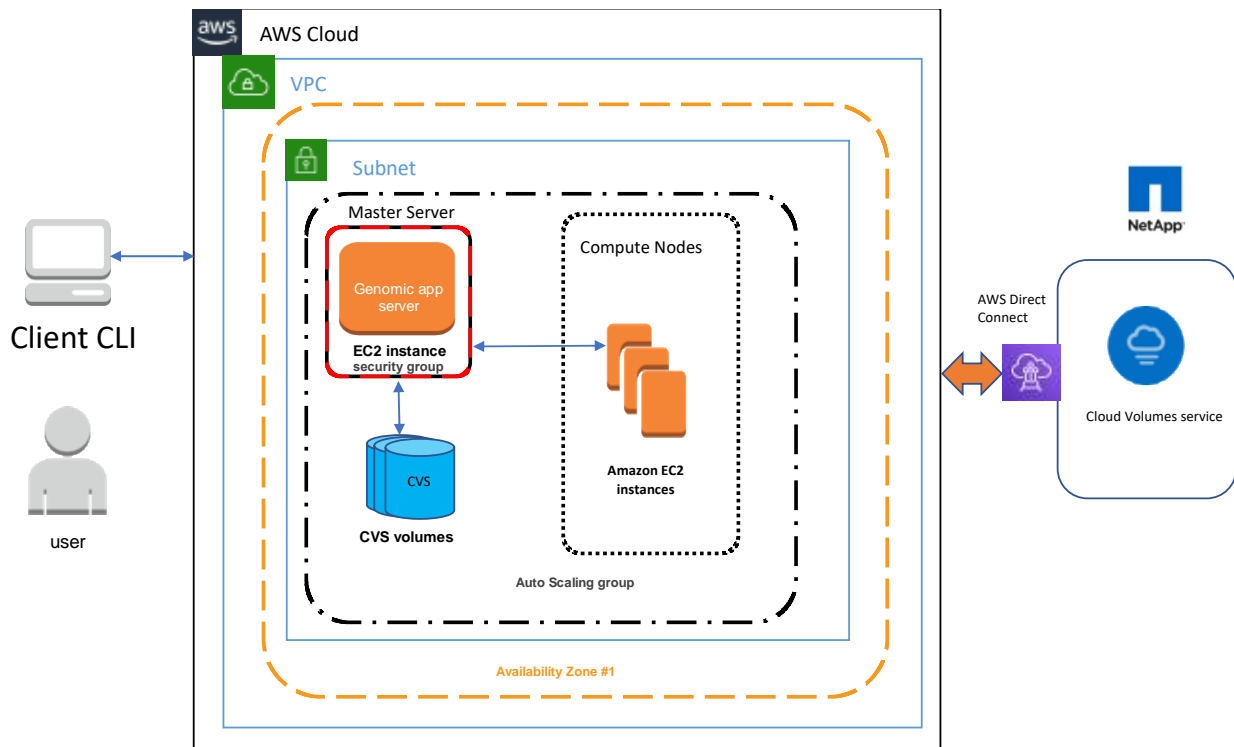
When you use CVS for AWS, you control your cloud performance by dynamically adjusting service levels. If you need to increase performance, you can increase the allocation (for example, 10TB provides 160MB/s) and/or choose a higher service level.

- The Standard service level offers very economical cloud storage, at just \$0.10 per gigabyte per month. It enables throughput up to 16MB/s for each terabyte allocated. This level is ideal as a low-cost solution for infrequently accessed data.
- The Premium service level delivers a good mix of cost and performance. At a cost of \$0.20 per gigabyte per month, it offers 4x the performance of the Standard level, with 64MB/s for each terabyte allocated. This is a good fit for many applications where data capacity and performance needs are balanced.
- The Extreme service level provides the highest performance. At a cost of \$0.30 per gigabytes per month, it enables up to 128MB/s for each terabyte allocated, and cloud volumes can scale to deliver several GB/s for reads and writes. Extreme is the best fit for high-performance workloads.

One of the unique features of NetApp Cloud Volumes Service for AWS is the capability to change performance on-the-fly. For example, if you need to have the Extreme performance tier for two hours a day and Standard performance for the rest, Cloud Volumes Service for [AWS can use API calls or a scheduler in Linux to facilitate that process](#).

General Architecture

HPC on AWS using CVS:



In the architecture diagram, you can see that the HPC application is configured with NetApp Cloud Volumes Service for AWS. With the combination of backups, snapshot copies, and right-sized throughput, you can easily host your high-performance workloads in the cloud with maximum data protection and nine 9s of data durability.

The HPC application is configured on an Amazon EC2 instance.

- Single or multiple cloud volumes are used as the dedicated storage for the datasets.
- The dataset volume(s) are provisioned using the Extreme service class because that class provides the highest throughput at a manageable cost.
- Dataset volume is backed up and restored from the backups rather than creating an individual volume.
- Cloud Backup Service backs data up to the S3 cloud.
- For more details on the configuration details, please refer the below link
 - [Cloud Backup Service](#)

The key components of the solution include:

- HPC application
- Amazon EC2 instances with autoscaling group.
- NetApp Cloud Volumes Service for AWS (storage)
- NetApp Snapshot Technology
- NetApp Cloud Backup Service

Benchmark Test with Comparison to Native File Services:

The benchmark use case calls for a massively scaled-out sequential read workload performed against thousands of genome files (GOR formatted). The files themselves are spread across one-to-many cloud volumes and many-to-many more Amazon Elastic Compute Cloud (EC2) instances. Because the content is genomically ordered, reads are efficient, resulting in the fastest possible analysis. A second use case calls for an equally scaled-out random read workload that's SQL-like and somewhat comparable to what one might expect from Apache's Hive project.

During testing, Cloud Volumes Service was compared against existing cloud-native storage present in Amazon Web Services (AWS), where the tests were run. The sequential read benchmark had a single hour to complete—with a stated goal of 10TiB per hour. The test itself comprised 2,500 GOR files representing 2000TiB of content, and was intended to evaluate several attributes, including:

- Four EC2 NFS servers atop Provisioned IOPS SSD (io1) Amazon Elastic Block Store (EBS) volumes
- A single EC2 NFS server configured similarly to the four servers.
- Four EFS volumes configured for a total of 4096MiBps of bandwidth using Provisioned Throughput mode.
- One EFS volume configured using Bursting Throughput mode.
- One Cloud Volumes Service volume provisioned with the maximum bandwidth supported by Cloud Volumes Service.

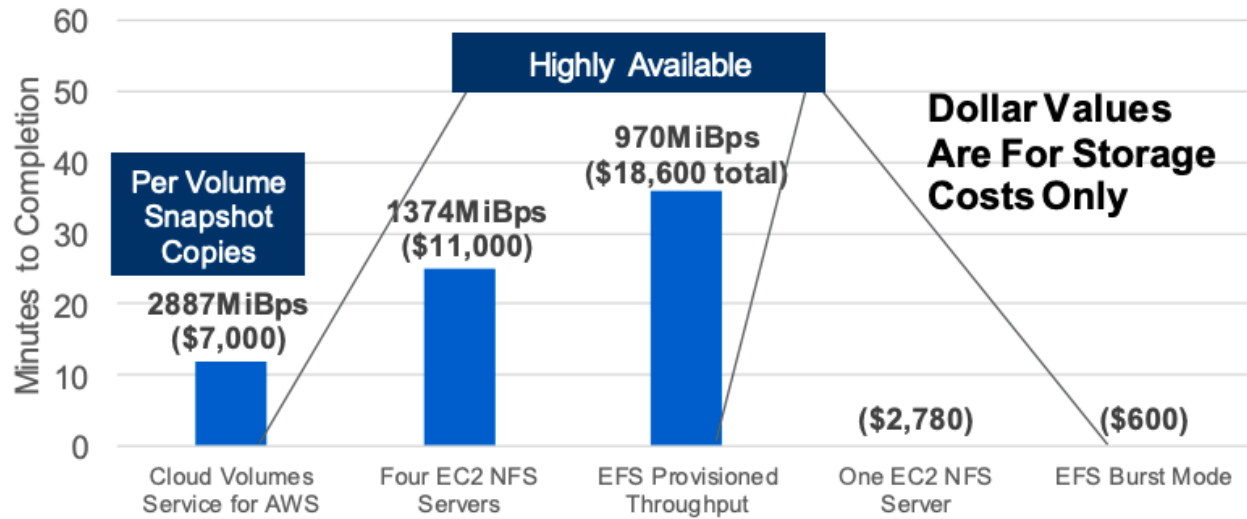
Benchmark Results

The results of the tests are shown in the following charts.

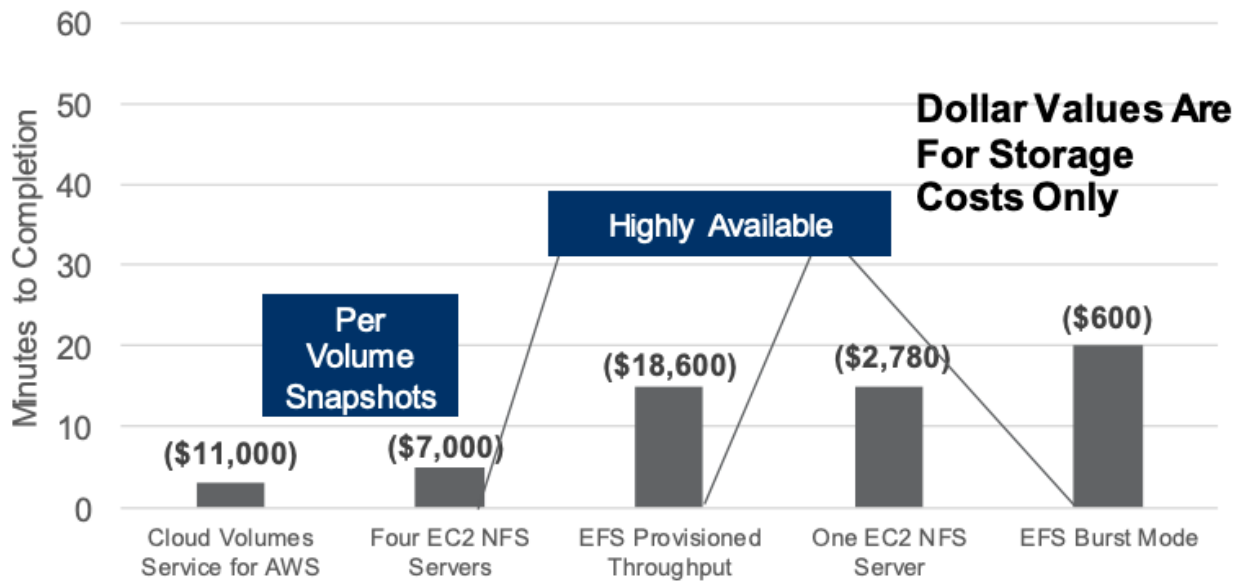
Time to completion is represented on the y-axis of the graph; the lower the value (expressed in MiBps), the better. With a little bit of math, you can see that the throughput achieved using the Cloud Volumes Service volume is equal to 9.91TiB per hour, which is 2.1 times the rate of the four self-managed NFS servers and 3.0 times that of the Provisioned Throughput configured EFS volumes.

The data from the tests themselves shows that a handful of workers took longer than the rest of the 2,500 workers, and that the throughput achieved throughout most of the Concurrent Versions System (CVS) test equaled roughly 3,500MiBps.

2 TB Sequential Read - 2,500 Files



Seek Query - 2,500 Files



Proven Performance

Check the [blog](#) illustrates the performance of an HPC applications on Cloud Volumes Service. We ran the benchmark with various workload mixtures and volume counts. The results were stunning.

The Cost Benefits of Cloud Volumes Service

Cloud Volumes Service is the lowest cost, highest quality solution for database hosting in the cloud. You can save a significant amount of time and money by changing performance levels on demand in CVS for AWS; other cloud storage solutions recommend that you configure performance and capacity to meet peak requirements, which means peak prices. Database performance requirements are rarely consistent—they require a system that’s adaptable, agile. But other cloud solutions offer up monolithic structures without swift-footed switching between performance levels, recommending static throughput of the highest level for the workload. You can use the NetApp API to change performance on the fly, or interface the service with a scheduler such as cron in Linux. That saves a lot of money.

For example, let’s say that you’re using Cloud Volumes Service and configured a volume at the Standard performance level (\$0.10/GB). If you realize that you need more performance, you can update the volume with an API call or scheduler and the change happens in seconds—it’s nondisruptive to clients. It’s just as easy to revert to the lower performance tier. So instead of continually paying for peak performance, you only incur added costs for the time you used the higher performance tier.

Think about it like this: If you have a performance intensive workload at certain times (such as online sales transactions on Black Friday, UBER, or Lyft during weekend peak times), you may need a volume to perform at the Extreme level for 30TB at \$.030/GB, but only during those peak periods. If you were to run at this level all the time it would cost \$9000/month. But with Cloud Volumes Service for AWS, when the intensive task finishes, you can quickly drop down to the Standard performance level for 160MB/s (16MB/s x 10TB) and meet the I/O needs for off-peak loads at a significantly lower cost. This performance level costs \$1,000 per month (10TB at \$0.10/GB). The cost savings vary, but if you run the processing intensive workload for 20% of the time and adjust the Service level, you can usually save about \$6,400 each month.

Note that the formula we used to calculate savings is: $\$9000 - ((\$9000 * 0.2) + (\$1000 * 0.8)) = \$6,400$, which **equals savings of more than 70%**.

Additionally, Cloud Volumes Service for AWS provides savings from:

- Space efficient snapshot copies, which only incur costs for unique data used in snapshots. This single 4KB copy is enough to protect all the data that is to be held in the snapshot the process very quick and extremely space efficient, regardless of whether your volume is a few megabytes in size or hundreds of terabytes.
- High performance storage that enables you to use fewer compute instances, which saves time and results in lower EC2 costs.
- Support for both NFS and SMB, which enables a dataset to be shared between Linux and Windows instances.
 - Alternative solutions require an expensive and slow data copy between multiple volumes.

Code Snippet

Cloud Volumes Service has rest APIs that can be called by various orchestration engines and scripting languages. Here are some example scripts that you can leverage to get started.

To list all cloud volumes and provision volumes for HPC workload

```
import requests
import json
import time
#Base URL
CVAPI_BASEURL="https://cv.us-west-1.netapp.com:8080/v1/"
CVAPI_APIKEY = "Supply your CVS API key here"
CVAPI_SECRETKEY = "Supply your CVS Secret key here "
#Headers
HEADERS = {
    'content-type': 'application/json',
    'api-key': CVAPI_APIKEY,
    'secret-key': CVAPI_SECRETKEY
}

getfilesystemDetailsHeaders = {
    'content-type': 'application/json',
    'api-key': CVAPI_APIKEY,
    'secret-key': CVAPI_SECRETKEY
}
filesystemURL = CVAPI_BASEURL + "/FileSystems"
filesystemCreateURL = CVAPI_BASEURL

class cvsAPI(object):
    # get FileSystems
    def get_fileSystems(self):
        getResult = requests.get(url=filesystemURL, headers=HEADERS)
        print("File system creation success, the response code : ", getResult.status_code)
        filesystemsData = getResult.json()
        for i in filesystemsData[:]:
            filesystemId = (i['filesystemId'])
            name = (i['name'])
            print("FileSystemId : ", filesystemId, " = VolumeName : ", name)

    # create Volume/filesystem
    def create_fileSystems(self):
        payload = {
            "name": "IAAS",
            "creationToken": "IAAS",
            "region": "us-west-1",
            "serviceLevel": "basic",
            "quotaInBytes": 1000000000000
        }
        postfileSystems = requests.post(filesystemURL, data=json.dumps(payload), headers=HEADERS)
        print("FileSystem Created", filesystemURL, postfileSystems.content)
        datafileSystems = postfileSystems.json()
        datafileSystems1 = datafileSystems['filesystemId']
        exportname = datafileSystems['creationToken']
        time.sleep(30)
        datafileSystems1 = datafileSystems['filesystemId']
        self.test = datafileSystems1
        self.export = exportname
        return datafileSystems1, exportname

volume = cvsAPI()
volume.get_fileSystems()
volume.create_fileSystems()
```

For more details, check out [our blog on Cloud Volumes Service APIs](#).

Customer Support for Cloud Volumes Service for AWS

Cloud Volumes Service for AWS is fully supported by NetApp, which has provided 24x7 enterprise class support for decades. If a customer has any questions or needs help with the service, they can open a ticket by sending an email to cvs-support@netapp.com.