

in

How PowerFlex Transforms Big Data with VMware Greenplum

f

Tue, 01 Nov 2022 21:18:15 -0000 | Read Time: 0 minutes

Tony Foster Sue Mosovich

🐦

Quick! The word has just come down. There is a new initiative that requires a massively parallel processing (MPP) database, and you are in charge of implementing it. What are you going to do? Luckily, you know the answer. You also just discovered that the Dell PowerFlex Solutions team has you covered with a solutions guide (<https://infohub.delltechnologies.com/t/vmware-tanzu-greenplum-on-dell-emc-powerflex/>) for VMware Greenplum.

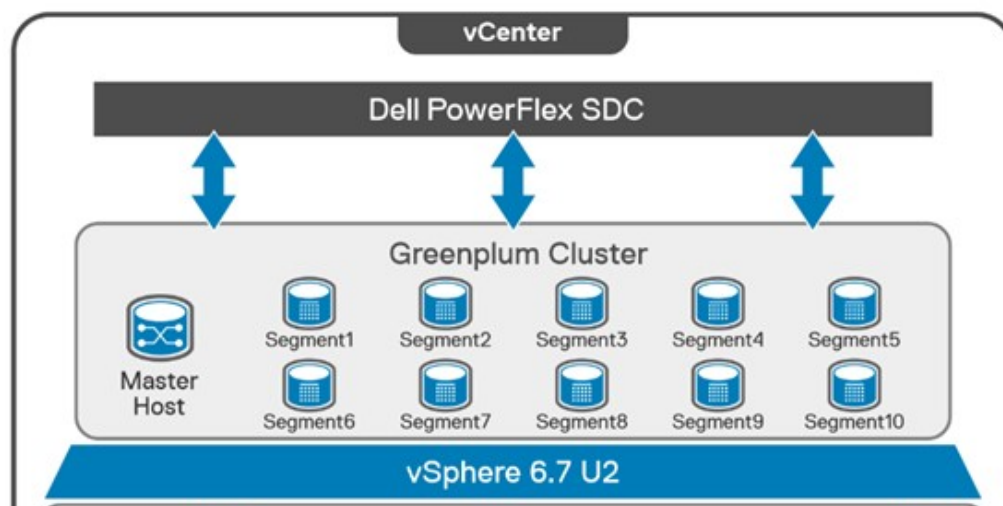
What is in the solutions guide and how will it help with an MPP database? This

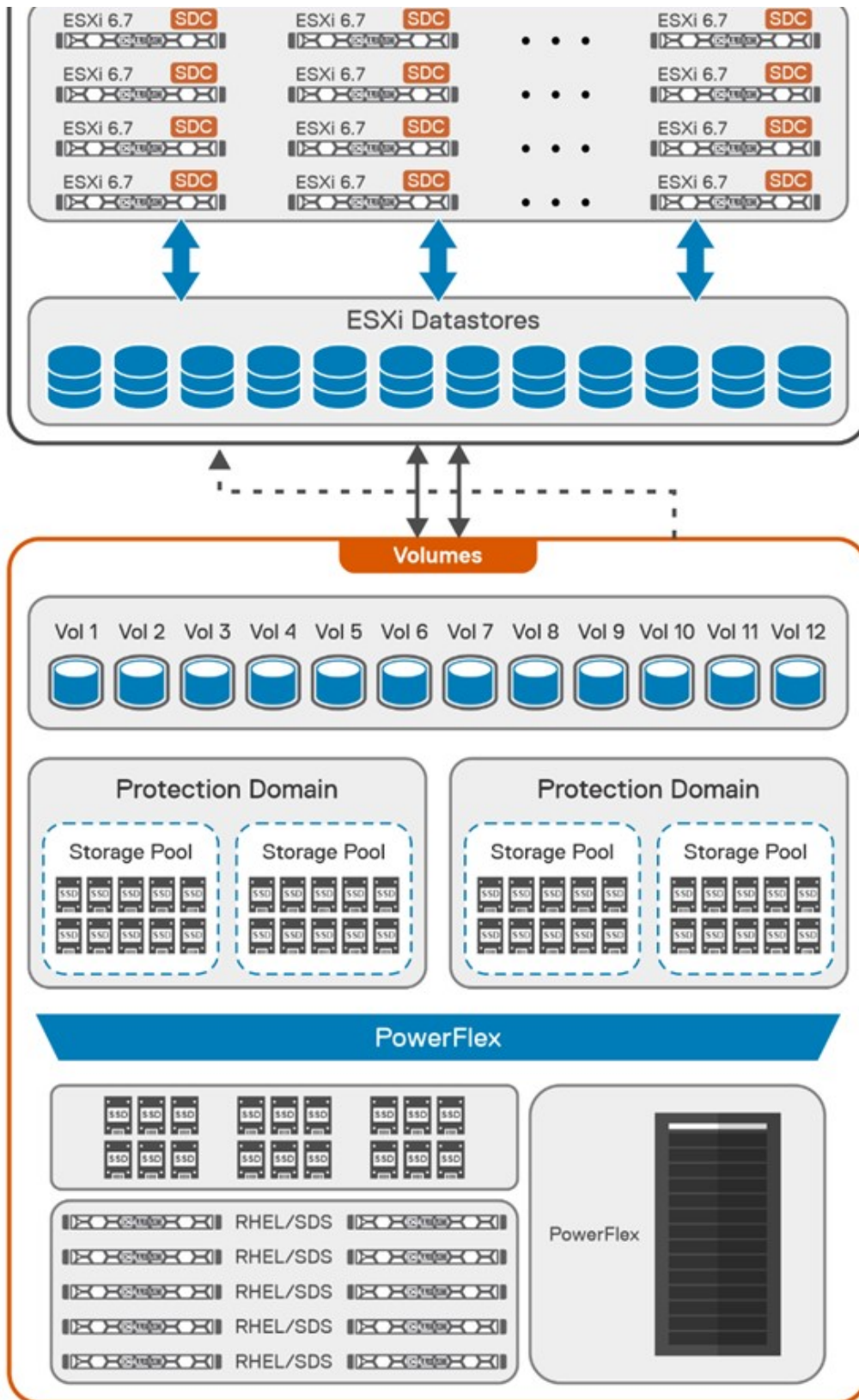
If you have read my other blogs or are familiar with PowerFlex, you know it has powerful transmorphic properties. For example, PowerFlex nodes sometimes function as both storage and compute, like hyperconverged infrastructure (HCI). At other times, PowerFlex functions as a storage-only (SO) node or a compute-only (CO) node. Even more interesting, these node types can be mixed and matched in the same environment to meet the needs of the organization and the workloads that they run.

This transmorphic property of PowerFlex is helpful in a Greenplum deployment, especially with the configuration described in the solutions guide. Because the deployment is built on open-source PostgreSQL, it is optimized for the needs of an MPP database, like Greenplum. PowerFlex can deliver the compute performance necessary to support massive data IO with its CO nodes. The PowerFlex infrastructure can also support workloads running on CO nodes or nodes that combine compute and storage (hybrid nodes). By leveraging the malleable nature of PowerFlex, no additional silos are needed in the data center, and it may even help remove existing ones.

The architecture used in the solutions guide consists of 12 CO nodes and 10 SO nodes. The CO nodes have VMware ESXi installed on them, with Greenplum instances deployed on top. There are 10 segments and one director deployed for the Greenplum environment. The 12th CO node is used for redundancy.

The storage tier uses the 10 SO nodes to deliver 12 volumes backed by SSDs. This configuration creates a high speed, highly redundant storage system that is needed for Greenplum. Also, two protection domains are used to provide both primary and mirror storage for the Greenplum instances. Greenplum mirrors the volumes between those protection domains, adding an additional level of protection to the environment, as shown in the following figure:





By using this fluid and composable architecture, the components can be scaled

Testing and validation with Greenplum: we have you covered

The solutions guide not only describes how to build a Greenplum environment, it also addresses testing, which many administrators want to perform before they finish a build. The guide covers performing basic validations with FIO (https://fio.readthedocs.io/en/latest/fio_doc.html) and gpcheckperf (https://greenplum.docs.pivotal.io/6-20/utility_guide/ref/gpcheckperf.html). In the simplest terms, these tools ensure that IO, memory, and network performance are acceptable. The FIO tests that were run for the guide showed that the HBA was fully saturated, maximizing both read and write operations. The gpcheckperf testing showed a performance of 14,283.62 MB/sec for write workloads.

Wouldn't you feel better if a Greenplum environment was tested with a real-world dataset? That is, taking it beyond just the minimum, maximum, and average numbers? The great news is that the architecture was tested that way! Our Dell Digital team has developed an internal test suite running static benchmarked data. This test suite is used at Dell Technologies across new Greenplum environments as the gold standard for new deployments.

In this test design, all the datasets and queries are static. This scenario allows for a consistent measurement of the environment from one run to the next. It also provides a baseline of an environment that can be used over time to see how its performance has changed -- for example, if the environment sped up or slowed down following a software update.

Massive performance with real data

So how did the architecture fare? It did very well! When 182 parallel complex queries were run simultaneously to stress the system, it took just under 12 minutes for the test to run. In that time, the environment had a read bandwidth of 40 GB/s and a write bandwidth of 10 GB/s. These results are using actual production-based queries from the Dell Digital team workload. These results are close to saturating the network bandwidth for the environment, which indicates that there are no storage bottlenecks.

The design covered in this solution guide goes beyond simply verifying that the

One of the key areas that we tested was the impact of snapshots on performance. Snapshots are a frequent operation in data centers and are used to create test copies of data as well as a source for backups. For this reason, consider the impact of snapshots on MPP databases when looking at an environment, not just how fast the database performs when it is first deployed.

In our testing, we used the native snapshot capabilities of PowerFlex to measure the impact that snapshots have on performance. Using PowerFlex snapshots provides significant flexibility in data protection and cloning operations that are commonly performed in data centers.

We found that when the first storage-consistent snapshot of the database volumes was taken, the test took 45 seconds longer to complete than initial tests. This result was because it was the first snapshot of the volumes. Follow-on snapshots during testing resulted in minimal impact to the environment. This minimal impact is significant for MPP databases in which performance is important. (Of course, performance can vary with each deployment.)

We hope that these findings help administrators who are building a Greenplum environment feel more at ease. You not only have a solution guide to refer to as you architect the environment, you can be confident that it was built on best-in-class infrastructure and validated using common testing tools and real-world queries.

The bottom line

Now that you know the assignment is coming to build an MPP database using VMware Greenplum -- are you up to the challenge?

If you are, be sure to read the solution guide (<https://infohub.delltechnologies.com/t/vmware-tanzu-greenplum-on-dell-emc-powerflex/>). If you need additional guidance on building your Greenplum environment on PowerFlex, be sure to reach out to your Dell representative.

Resources

- Dell PowerFlex (<https://www.dell.com/en-us/dt/storage/powerflex.htm>)

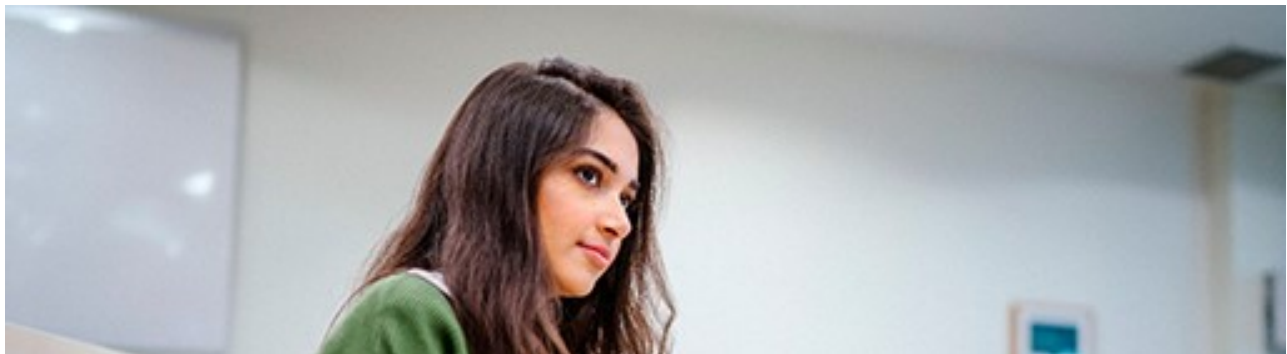
/wonder_nerd)

LinkedIn (https://linkedin.com/in/wondernerd/)

- Sue Mosovich – VMware

Tags: VMware PowerFlex Greenplum

Related Blog Posts



(/p/introducing-nvme-over-tcp-nvme-tcp-in-powerflex-4-0/)

VMware PowerFlex NVMe/TCP

Introducing NVMe over TCP (NVMe/TCP) in PowerFlex 4.0

(/p/introducing-nvme-over-tcp-nvme-tcp-in-powerflex-4-0/)



Kevin M. Jones, Tony Foster

Fri, 26 Aug 2022 18:59:38 -0000 | Read Time: 0 minutes

(/p/introducing-nvme-over-tcp-nvme-tcp-in-powerflex-4-0/)

Anyone who has used or managed PowerFlex knows that an environment is built from three lightweight software components: the MDM, the SDS, and the SDC. To deploy a PowerFlex environment, the typical steps are:

1. Deploy the MDM...



(/p/powerflex-and-amazon-destination-eks-anywhere/)

VMware vSphere Kubernetes PowerFlex Amazon EKS

PowerFlex and Amazon: Destination EKS Anywhere

(/p/powerflex-and-amazon-destination-eks-anywhere/)

 Tony Foster

Wed, 19 Jan 2022 17:09:54 -0000 | Read Time: 0 minutes

(/p/powerflex-and-amazon-destination-eks-anywhere/)

Welcome to your destination. Today Dell Technologies is pleased to share that Amazon Elastic Kubernetes Service (Amazon EKS) Anywhere has been validated (<https://aws.amazon.com/eks/eks-anywhere/partners>) on Dell PowerFlex ([https://www.delltechnologies.com/en-](https://www.delltechnologies.com/en-us/whitepapers/PowerFlex-Cloud-Managed-Infrastructure-Partners-Program-Announces-Amazon-EKS-Anywhere-Validation)

logo

© 2023 Dell Inc. (<https://www.dell.com/learn/us/en/uscorp1/site-terms-of-use-copyright>)

Privacy (<https://www.dell.com/learn/us/en/uscorp1/policies-privacy>)

Terms Of Use (<https://www.dell.com/learn/us/en/uscorp1/site-terms-of-use>)

Legal (<https://www.dellemc.com/en-us/customer-services/product-warranty-and-service-descriptions.htm>)

Anti-Slavery and Human Trafficking (<https://i.dell.com/sites/doccontent/corporate/corp-comm/en/Documents/dell-california-trafficking.pdf>)