



VMware Cloud on AWS: vCenter Architecture

VMware Architecture

Table of contents

VMware Cloud on AWS: vCenter Architecture	3
Overview	3
vCenter Inventory	4
Datacenters	4
Clusters	4
Hosts	4
Resource Pools	5
Datastores	6
Networks	7
SDDC Group	8
Multi SDDC Management	10
Hybrid Linked Mode	10
vCenter Linking	10
User Access and Roles	11
Cloud Admin Role	11
Cloud Admin Group	11
Cloud Admin User	11
Authenticating Users with LDAP or Enterprise SSO (single sign-on)	11
On-premises LDAP domain as an identity source for the SDDC vCenter Server	11
Configuring Enterprise Federation (SSO)	11

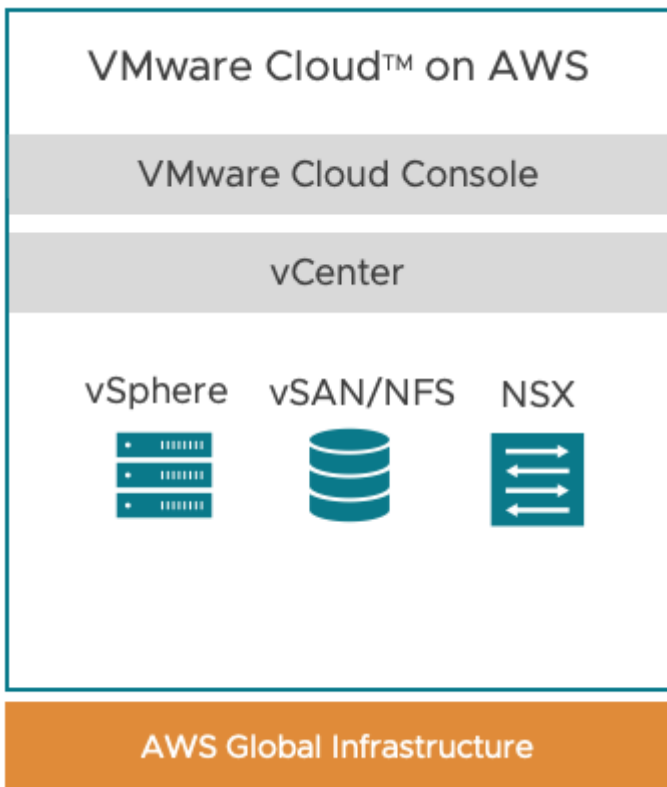
VMware Cloud on AWS: vCenter Architecture

Overview

The computing assets within a VMware Cloud on AWS SDDC (Software-Defined Data Center) consist of a collection of bare-metal AWS EC2 servers, each installed with ESXi and managed by a vCenter Server instance.

VMware vCenter Server offers centralized administration for vSphere virtual infrastructure, empowering essential features such as VMware vSphere vMotion®, VMware vSphere Elastic Distributed Resource Scheduler™, VMware vSphere High Availability (HA), and Hybrid Linked Mode.

The following image shows VMware Cloud on AWS SDDC components.



In the following sections, we will delve into the specifics of the vCenter architecture.

vCenter Inventory

The vCenter inventory provides a comprehensive overview of your virtual infrastructure, encompassing essential components such as datacenters, clusters, resource pools, datastores, and networks. This centralized management platform enables efficient administration and organization of your virtualized environment, ensuring optimal resource utilization, performance, and availability.

Datacenters

A VMware Cloud on AWS SDDC supports a single logical datacenter which is named "SDDC-Datacenter". All resources for the SDDC reside within this datacenter. Datacenters within an SDDC cannot be renamed.

Clusters

In VMware Cloud on AWS SDDC, a cluster represents a group of ESXi hosts that work together to provide compute resources for running virtual machines (VMs) and supporting infrastructure services. The first cluster created when SDDC is created is called a base cluster.

Clusters enable features like vSphere High Availability (HA) and vSphere Distributed Resource Scheduler (DRS) functionalities within an SDDC. Clusters follow the naming convention "Cluster-n," where "n" represents the sequential number assigned to the cluster. [Renaming a cluster](#) is a supported operation and can be executed from the VMware cloud console only.

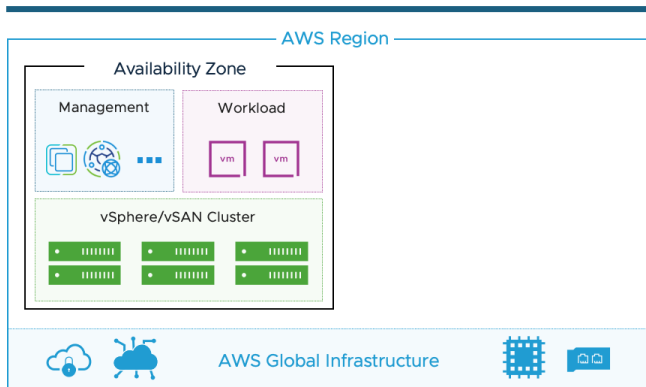
Every VMware Cloud on AWS SDDC is created with a cluster denoted as "Cluster-1," serving both management appliances and end-user workloads. Often referred to as the base or primary cluster of the SDDC. Cluster-1 always hosts the management appliance for the SDDC, it forms the foundation. Additional clusters can be created as necessary, within the limits [configuration maximum](#) defined for the SDDC.

By default, clusters are associated with a single AWS Availability Zone (AZ), referred to as a standard or non-stretched cluster. However, VMware Cloud on AWS enables deploying an SDDC with a [stretched-cluster](#) for customers requiring higher availability to meet their business demands. Stretched clustering enables the distribution of cluster hosts across two Availability Zones (AZs) within an AWS region.

The vCenter, NSX Managers, and NSX edge appliances are typically deployed in the same AZ on standard and stretched clusters. In the event of an AZ failure, if vCenter, NSX Managers, and NSX edge appliances are in the impacted AZ, vSphere HA will restart these appliances and customers' workloads in the remaining AZ.

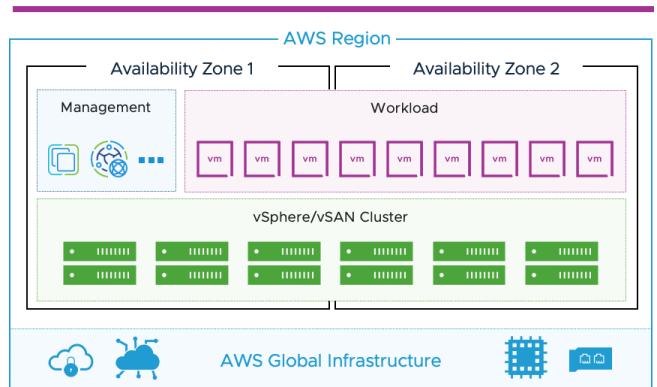
The decision between a stretched cluster and a standard SDDC is critical and must be determined during the **initial deployment** of the SDDC. Once deployed, switching from a standard SDDC to a stretched cluster SDDC (or vice-versa) is not supported. Additionally, all clusters within the SDDC will be uniformly configured as either stretched or standard.

Standard Cluster



99.9% Availability guarantee

Stretched Cluster



99.99% Availability guarantee

Entry level stretch clusters (2/4 hosts with 99.9%)

The image shares a comparison of workload data when in a standard and a stretched cluster.

Note: The north-south(N-S) traffic of workloads in a stretched cluster traverses from the edge node and via the uplink of AZ where the edge resides. Please refer to the [Network Architecture](#) document for traffic flow information.

Hosts

In the context of VMware vCenter, a "Host" refers to a physical server that runs the VMware ESXi hypervisor. For production

deployments, a primary cluster within an SDDC necessitates a minimum of 2 hosts.

vSphere High Availability (HA) ensures the availability of the virtual machines in your SDDC. If a host fails, vSphere HA restarts its VMs on a different host. All clusters in a VMware Cloud on AWS SDDC are configured to use vSphere HA. This setting **cannot be reconfigured**.

To ensure the availability of all workload and management VMs in your SDDC, VMware Cloud on AWS must maintain sufficient capacity to power them on in the event of host failure. HA admission control is the primary mechanism for capacity maintenance. Admission control imposes constraints on resource usage and can prevent any action that consumes more resources than the cluster can support during a failover. These constraints apply to actions like powering on or migrating a VM, or reserving additional CPU or memory resources for a VM, and effectively limit the availability of host resources as shown here:

- In a two-host i3en.metal SDDC cluster, admission control prevents you from powering-on more than 36 VMs or assigning more than 1152 MHz CPU reservation to a single VM.
- In SDDC clusters with three to five hosts, admission control reserves one host for failover.

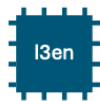
Details regarding VMware Cloud on AWS instance types can be found in the "[Feature Brief: SDDC Host Types](#)".



M7i

- 4th gen Intel Xeon Sapphire Rapids
- Turbo frequency up to 3.8 GHz
- 48 physical cores (96 logical cores with HT)
- 384 GiB Memory
- Up to 37.5 Gbps Networking Speed

CPU Intensive Workloads
AI/ML Workloads
Ransomware & Disaster Recovery



I3en

- Intel 2nd Generation Xeon Scalable
- Cascade Lake Microarchitecture
- 2.5 GHz x 96 logical cores (48+HT)
- 768 GB RAM
- ~45 TiB NVMe

Storage-optimized clusters for databases, large file systems, analytics, and high random I/O



I4i

- Intel 3rd Generation Xeon Scalable
- Ice Lake Microarchitecture
- 2.9 GHz x 128 logical cores (64+HT)
- 1024 GB RAM
- ~20 TiB NVMe

General purpose computing, resource-optimized for all applications with 50% more memory and 2x storage of i3

VMware Cloud on AWS also allows customers to select a reduced number of CPU cores to run per host as the default number of cores for the host type. This feature is designed to avoid running servers with underutilized CPU capacity. This custom core count cluster feature is configurable only for the additional cluster created on the SDDC and does not apply to the base cluster. Refer to [Feature Brief: Custom CPU Core Count](#) for more information.

Hosts within a cluster share the same instance type; the creation of "mixed-host" clusters is not supported. However, hosting multiple clusters with different instance types within a single SDDC is permissible.

New hosts can be incorporated into one cluster within the SDDC and will allocate their resources exclusively to that cluster. Also, the number of cores per CPU that you select applies to all the hosts in that cluster at the time of deployment.

Additionally, there are prescribed maximums not only on the number of hosts per cluster but also on the total number of hosts permissible within the entire SDDC. Refer to the [configuration maximums](#) for specific details.

Resource Pools

A "resource pool" is a logical abstraction that aggregates and partitions CPU and memory resources under a cluster. Resource pools within an SDDC primarily serve to safeguard the management appliances and perform two key functions:

1. They provide a mechanism to reserve compute resources for management appliances.
2. They serve as the designated objects for assigning permissions to management appliances.

Upon creation, the SDDC is configured with two resource pools within the base cluster:

- Mgmt-ResourcePool - This Resource Pool contains the management appliances
- Compute-ResourcePool - This Resource Pool is Intended for hosting end-user workloads.

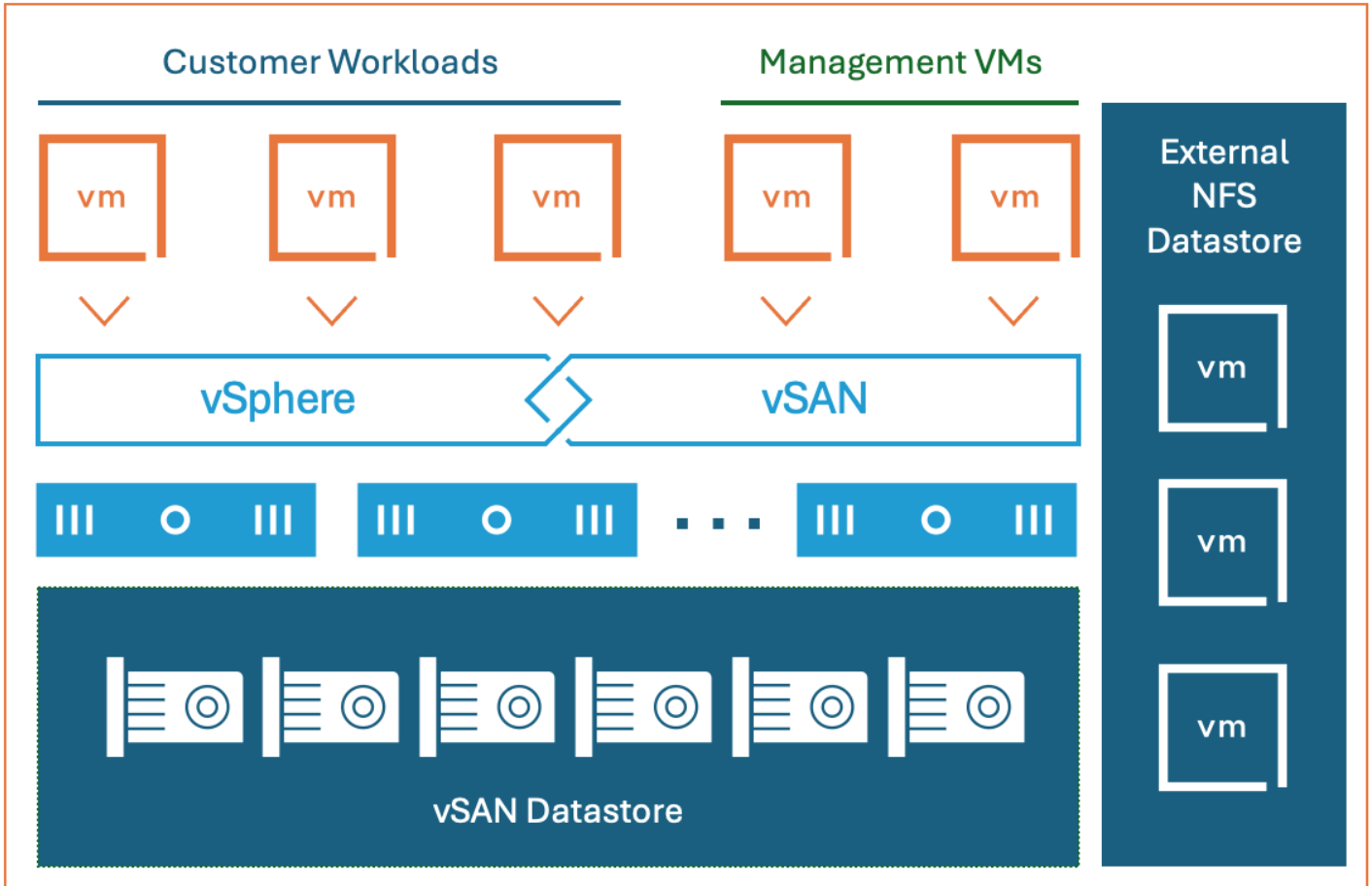
Additional resource pools can be created under cluster or as child resource pools under 'Compute-ResourcePool'. A Simple use

case can be for hosting a set of workloads with specific resource requirements.

Because of vCenter's permissions model, it is possible to create workloads outside of Compute-ResourcePool. However, this approach is highly discouraged due to the risk of resource conflicts between workloads inside resource pools and those outside. As a general practice, it is recommended to place all workloads within the Compute-ResourcePool to mitigate this issue.

Datstores

In VMware vCenter, a "datastore" refers to a storage repository where virtual machine (VM) files, such as virtual disks, configuration files, and snapshots, are stored.



VMware Cloud on AWS offers two distinct types of datastore storage technologies: vSAN and NFS.

vSAN and NFS offer further choices:

vSAN provides:

- ESA (Express Storage Architecture)
- OSA (Original Storage Architecture)

Note - vSAN ESA requires a VMware cloud on AWS SDDC version 1.24 a minimum of three i4i nodes and deployed single AZ.

NFS Storage includes:

- Amazon FSx for NetApp ONTAP
- VMware Cloud Flex Storage

The base cluster of the SDDC with a datastore may contain the following:

1. vsanDatastore - This datastore is reserved for SDDC management appliances with vSAN datastore.
2. WorkloadDatastore - This datastore is available for end-user workloads.
3. managementDatastore- This datastore is reserved for SDDC management appliances that are running M7i hosts(using NFS storage only).

Storage pools within the SDDC 'vsanDatastore' and 'managementDatastore' are protected and cannot be modified.

Since management appliances reside exclusively within the base cluster, additional clusters that are added to the SDDC will contain only a single datastore for end-user workloads. The naming convention for vSAN datastores is "WorkloadDatastore (n)", where "n" is the sequence number of the datastore. Note that the sequence begins with an un-numbered datastore and is inconsistent with the cluster naming convention. As an example:

- Cluster-1 owns WorkloadDatastore
- Cluster-2 owns WorkloadDatastore (1) etc

Note:

- Clusters with ESA are supported with i4i.metal hosts only. For more information refer to [vSAN ESA with VMware Cloud on AWS: Technical Deep Dive](#)
- Custom workload datastore name is supported for NFS storage only.

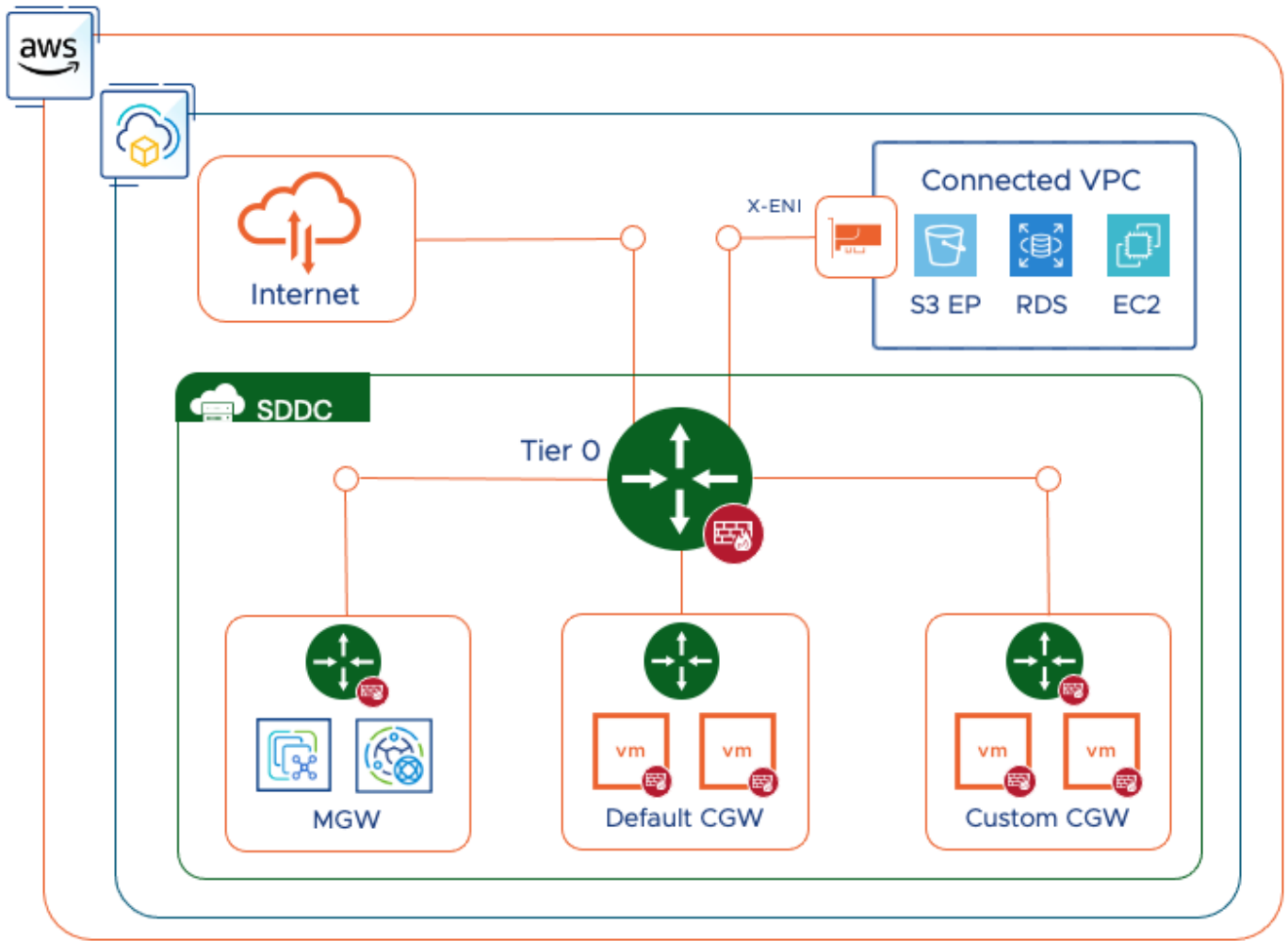
Networks

The hosts within the VMware Cloud on AWS SDDC utilize an NSX Virtual Distributed Switch (VDS). Management of the VDS is handled by the local clustered instance of NSX deployed within the management resource pool of the SDDC. Standard NSX user interface (UI) is provided for customers to simplify administration.

The native UI for NSX allows customers to perform operations such as segment creation, configuring DFW (Distributed Firewall), VPN, additional Compute gateway(CGW), and more.

An SDDC network has two notional tiers:

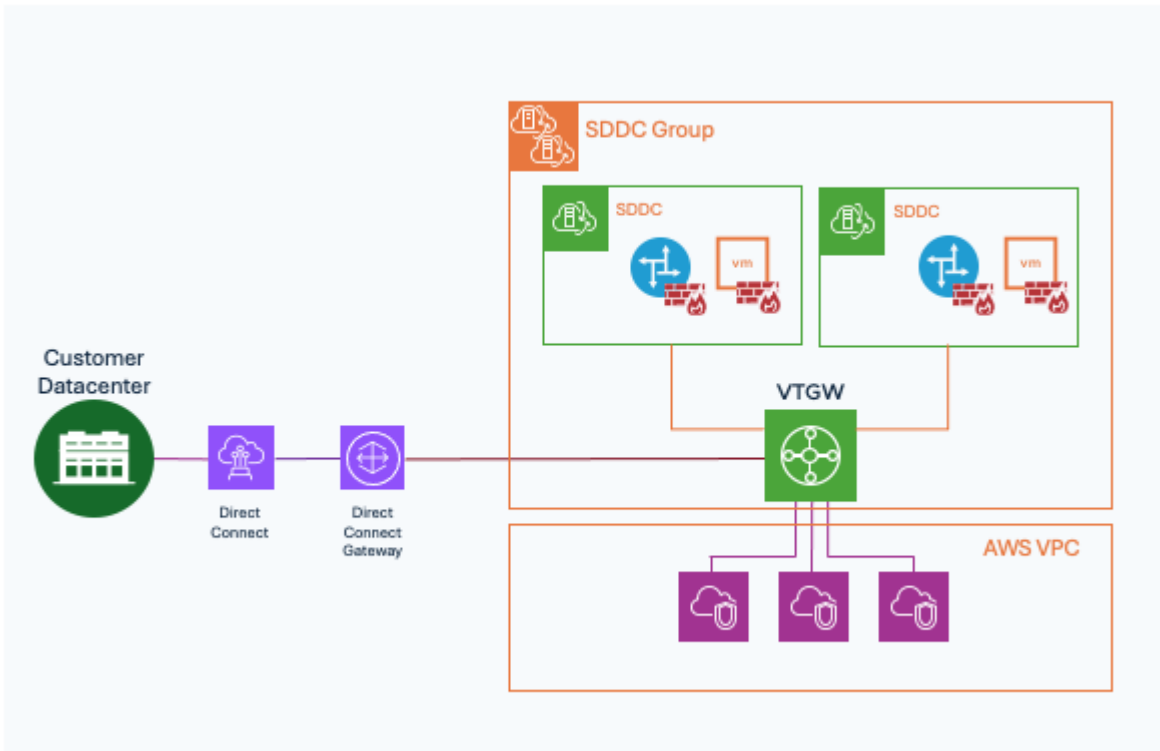
- Tier-0 handles north-south traffic (traffic leaving or entering the SDDC, or between the Management and Compute gateways). In the default configuration, each SDDC has a single Tier-0 router.
- Tier-1 handles east-west traffic (traffic between routed network segments within the SDDC). In the default configuration, each SDDC has a single Tier-1 router. You can create and configure additional Tier-1 gateways if you need them. See [Add a Custom Tier-1 Gateway to a VMware Cloud on AWS SDDC](#).



The Image highlights the NSX networking on VMC on AWS environment. Please refer to the Network Architecture document for more information.

SDDC Group

An SDDC group uses VMware Transit Connect to provide high-bandwidth, low-latency connections between SDDCs in the group. An SDDC group can include customer’s VPCs. Customers can also add an AWS Direct Connect Gateway (DXGW) to provide connectivity between SDDC group to on-premises SDDCs.

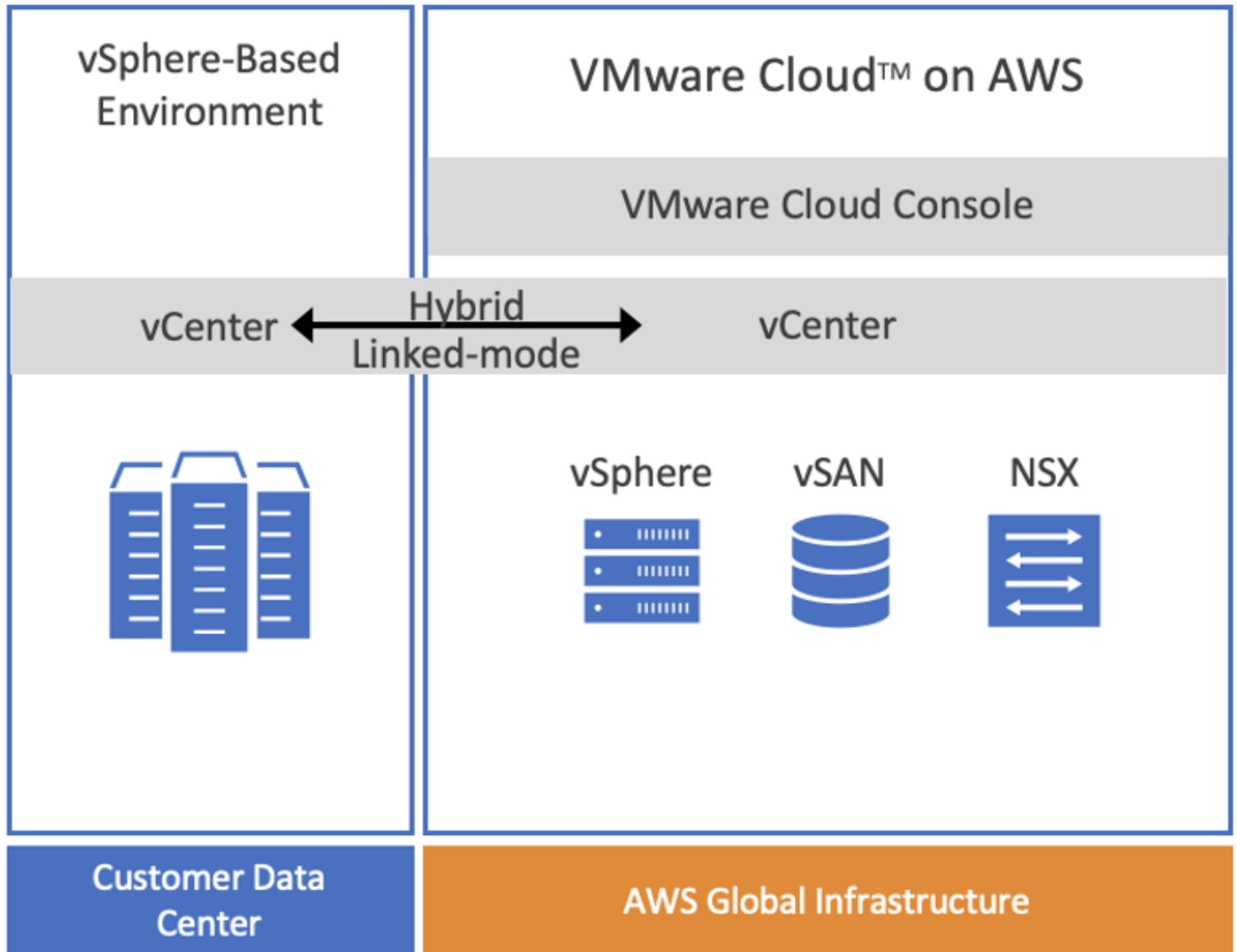


Multi SDDC Management

VMware Cloud on AWS enables multiple SDDC management under a single pane with features like Hybrid linked mode and vCenter linking.

Hybrid Linked Mode

Hybrid Linked Mode allows you to link your VMware cloud on AWS vCenter instance with an on-premises vCenter Single Sign-On domain. If you link your cloud vCenter to a domain that contains multiple vCenter instances linked using Enhanced Linked Mode, all those instances are linked to your cloud SDDC.



In order to link your cloud SDDC to your on-premises vCenter, you must add an identity source to the SDDC LDAP domain.

Refer to [Configuring Hybrid Linked Mode](#) for more information.

vCenter Linking

An organization that contains an [SDDC deployment group](#) can link the vCenter systems in those SDDCs to enable an administrator to manage their combined inventories in the same vSphere Client view.

When you enable vCenter linking in an SDDC group, a cloud administrator can log in as `cloudadmin@vmc.local` and use the vSphere Client to manage all the vCenter systems in the group. If the `cloudadmin@vmc.local` account configures these systems to use single sign-on, then users with accounts in that single sign-on domain can access all the linked systems in the group.

After vCenter linking has been enabled in an SDDC group, the vCenter systems in SDDCs added to the group are linked automatically, and vCenter systems in SDDCs that are removed from the group are unlinked automatically.

User Access and Roles

The permissions model of a VMware cloud on AWS SDDC is designed to allow the service provider (VMware) and the tenant (customer) to manage the environment in a [shared responsibility model](#). The high-level goals of the shared permissions model are:

- Permit VMware and customer joint access to vCenter
- Enable customers to manage their workloads, users/groups (via AD/LDAP), tags, permissions (on their inventory objects), roles (from a subset of their permissions), etc.
- Protect VMware-managed objects (management appliances, users/groups, global policies, roles, permissions, hosts, storage, etc..).

Refer to [Feature Brief: Understanding Accounts, Roles and Privileges](#) for more information.

Cloud Admin Role

The Cloud Admin Role defines the permissions for customers within the SDDC. The CloudAdmin role has the privileges necessary to create and manage SDDC workloads and related objects such as storage policies, content libraries, vSphere tags, and resource pools. This role cannot access or configure objects that are supported and managed by VMware, such as hosts, clusters, and management virtual machines.

Note:

The CloudAdmin user can grant other users or groups read-only access to VMware Cloud on AWS vCenter management objects such as the **Mgmt-ResourcePool**, **Management VMs** folder, **Discovered Virtual Machines** folder, **vmc-hostswitch**, and **vsanDatastore**. Because this read-only access does not propagate to management objects, you cannot grant it as a Global Permission and instead must explicitly grant it for each management object.

Refer to [vSphere Permissions and Privileges in VMware Cloud on AWS](#) for more information.

Cloud Admin Group

The Cloud Admin Group is used for defining access to objects within the vCenter inventory. This group has been granted a Cloud Admin Role within Global Permissions as well as on the datacenter object of vCenter. It also has read-only permission for the management resources within vCenter (storage, networks, Resource Pools, VMs) as well as on the “Discovered virtual machines” folder.

Users or groups managed within customers Active Directory (AD) or Lightweight Directory Access Protocol (LDAP) will be allocated to the Cloud Admin Group or with tailored permissions and roles within vCenter.

Cloud Admin User

Customers can access vCenter using the cloudadmin@vmc.local user account. The Cloud Admin User is initially provisioned with a randomized password which is available from within the SDDC view of the VMC console. It should be noted that if this password is changed from within vCenter, the password will not propagate to the VMC console (which will continue to reflect the old password).

Refer to [Cloud Admin privileges](#) documentation for more information.

Note: Access to the vCenter user interface is blocked by default and can only be accessed by creating a management firewall rule for a trusted IP or VPN interface. Refer to the [Add or Modify Management Gateway Firewall Rules](#) for more information.

Authenticating Users with LDAP or Enterprise SSO (single sign-on)

VMware Cloud on AWS enables integration for user authentication and authorization in below ways.

On-premises LDAP domain as an identity source for the SDDC vCenter Server

You can configure Hybrid Linked Mode from your SDDC if your on-premises LDAP service is provided by a native Active Directory (Integrated Windows Authentication) domain or an OpenLDAP directory service.

This step is optional when configuring Hybrid Linked Mode from the VMware Cloud Gateway, but adding an identity source does allow you to configure users or groups with a lesser level of access than the Cloud Administrator

Refer to [Add an Identity Source to the SDDC LDAP Domain](#) for more information.

Configuring Enterprise Federation (SSO)

Customers are highly encouraged to secure their environment by configuring federated authentication backed by an enterprise

managed system (typically Active Directory). VMware Cloud services users with a federated domain use their corporate credentials to log in to the Cloud Services Console across Organizations.

Refer to the [Feature brief](#) and [what is Enterprise Federation](#) for more information.

Additionally, workload access to AWS services can also be secured by using AWS IAM roles. Refer to [Using IAM Roles Anywhere to Help Secure VMware Cloud on AWS Workloads](#) for more information.

