

## vSphere Standard Switch (VSS)

Components, Terminology and Data Plane Details

VSPHERE NETWORKING COMPONENTS VSS 1

VSPHERE NETWORK TERMINOLOGY VSS 2

**Virtual Ethernet Adapters**

- Virtual network interface card (vNIC) - virtual machine's interface to the network
- Virtual network kernel network interface card (vmxnic) - vSphere Hypervisor's interface to network (Mgmt, NFS, iSCSI, vMotion, FT)

**Physical Ethernet Adapter**

- Physical network interface card (NIC) - HESOs communicate with outside ESXi host. This is also called as vmnic.

**vSphere Standard Switch (VSS)**

- Forwards packets between vnic, vmxnic, and vmnics

**Port Group (PG)**

- Group of ports sharing the same configuration (e.g. VLAN)

**Virtual Port:** These are the ports where virtual machines or vmnics are connected to the virtual switch. Type of port binding available - No binding. Also called as Ephemeral port binding.

**Uplink:** Connections to physical switches

**NIC Team:** A group of vmnics connected to the same physical network

VIRTUAL ETHERNET ADAPTER TYPES VSS 3

**vlannc** - Emulates AMD 78C970 PCnet32 LANCE 10 Mbps NIC

**vmmnet** - virtual adapter with performance optimization, requires VMware Tools

**flexible** - Behaves as either vlannc or vmmnet. If VMware Tools are installed it will act as vmmnet adapter

**vmxnic000** - Emulates Intel 82545EM Gigabit Ethernet NIC

**vmxnic2 (enhanced)** - Based on vmmnet, with jumbo frame and hardware offloads

**vmxnic3** - New generation of virtual adapter, supports RSS, IPv6 offloads, and VSI/VSI-X interrupts. Available on a limited set of guest operating systems

VSPHERE STANDARD SWITCH - DATA PLANE DETAILS VSS 4

VSS Configuration Parameters

VSS PARAMETERS - VLAN VSS 5

**VLAN Configuration** are done at port group level. The following options are available as part of the **VLAN ID (optional) field**:

- None (0) - No VLAN tagging
- Use this option when you don't want any VLAN tagging (default option). VLAN tagging is left to external switch. Also called as External Switch Tagging - EST
- All (4095) - Trunk configuration
- Use 4095 when linking the Guest virtual machines to the VLAN tagging. This is equivalent to the trunk port configuration on the physical switch port. Also called as Virtual Guest Tagging - VGT
- Type any number as VLAN ID in the field - available numbers 1 to 4094
- In this mode, based on the VLAN number configured on the port group the traffic going out of the virtual switch will be tagged accordingly. Also called as Virtual Switch Tagging - VST

**Best Practices**

- Use VLANs to isolate different traffic types (i.e. use Virtual Switch Tagging (VST))
- While using VST, make sure the physical switch ports, where the host's vmnics are connected, are configured as trunk ports. These trunk ports then carry the tagged traffic from virtual switch.

VSS PARAMETERS - SECURITY VSS 6

**Security Configuration at Port Group level**

- Promiscuous mode
  - Reject - Only forwards traffic that is destined for the VM
  - Accept - All packets received on a particular VLAN of the port group are forwarded to all the VMs connected to the port group
- MAC address Change
  - Reject - If the VM changes its MAC address from the one that was configured in vms file, the inbound traffic with new MAC is dropped
  - Accept - Any change in MAC address is accepted and frames are continued to receive
- Forged Transmits
  - Reject - Outbound traffic is dropped if sent with different MAC address
  - Accept - No checks for MAC address is performed

**Best Practices**

- Select Promiscuous mode - Reject (default)
- Average Bandwidth - Specifies the allowed average bandwidth in kbps
- Peak Bandwidth - The maximum bandwidth allowed in kbps
- Burst Size - Establishes the maximum number of kilo bytes to allow in a burst

VSS PARAMETERS - TRAFFIC SHAPING VSS 7

**Traffic Shaping Policy (Egress traffic management only)**

- A virtual machine's network bandwidth can be controlled by enabling the network traffic shaper.

The following three parameters dictate the traffic-shaping policy on a virtual port

- Average Bandwidth - Specifies the allowed average bandwidth in kbps
- Peak Bandwidth - The maximum bandwidth allowed in kbps
- Burst Size - Establishes the maximum number of kilo bytes to allow in a burst

**Best Practices**

- Use this feature in the following scenarios
  - You want to control vMotion traffic going out of one physical interface such that it doesn't impact virtual machine traffic flowing through the same interface
  - You want to control a tenant's virtual machine bandwidth usage

VSS PARAMETERS - TEAMING VSS 8

**Teaming Configuration on a Port Group depends on the way uplinks of a virtual switch are connected to the physical switches.**

Based on the physical switch connection option (shown in VSS/VDS 3) choose the load balancing type

- Load balancing types
  - Route based on originating virtual port (Port ID Hash)
  - Route based on source MAC hash (MAC Hash)
  - Route based on IP hash (IP Hash)
  - Use Explicit Failover order (Explicit Failover)
- After choosing the Load balancing type
  - Make sure that the Active, Standby and Unused adapters are identified for the port group

**Best Practices**

- Always connect two or more vmnics to a virtual switch and configure teaming for redundancy and higher bandwidth.
- If possible terminate the vmnics on two separate physical access switch

VSS PARAMETERS - TEAMING VSS 8

**Teaming Configuration on a Port Group depends on the way uplinks of a virtual switch are connected to the physical switches.**

Based on the physical switch connection option (shown in VSS/VDS 3) choose the load balancing type

- Load balancing types
  - Route based on originating virtual port (Port ID Hash)
  - Route based on source MAC hash (MAC Hash)
  - Route based on IP hash (IP Hash)
  - Use Explicit Failover order (Explicit Failover)
- After choosing the Load balancing type
  - Make sure that the Active, Standby and Unused adapters are identified for the port group

**Best Practices**

- Always connect two or more vmnics to a virtual switch and configure teaming for redundancy and higher bandwidth.
- If possible terminate the vmnics on two separate physical access switch

VSS PARAMETERS - FAILOVER DETECTION VSS 9

**Network Failover Configuration**

- Link State Only (default)
  - Use when supported by physical switches. For Example - Cisco switches.
  - Beacon probing
  - Use when no link state tracking support is available on physical switches and you don't have redundant connection between access to distribution physical switches.
- Beacon probing should be used when you have more than 2 uplinks in a team. This software approach of detecting failure between the vmmnics is very useful when there is no link state tracking support on the physical switches.

**Best Practices**

- Cisco switches provides Link State tracking feature. By enabling this feature any link failure between access and distribution (upstream) switch is indicated to the ESXi host such that the traffic gets moved to another working uplink. This feature should be used if you don't have redundant connection from the access switch to the distribution switches.
- Beacon probing should be used when you have more than 2 uplinks in a team. This software approach of detecting failure between the vmmnics is very useful when there is no link state tracking support on the physical switches.

VSS PARAMETERS - FAILOVER OPTION AFTER LINK COMES BACK VSS 10

**Failback Configuration**

- Use this feature when using Explicit Failover (EF) teaming configuration
- Configure "No" when selecting EF teaming
- Default Failback is "yes". That means after the link comes back after a failure the traffic will be moved back on the link.

**Best Practices**

- Failback configuration works with Active-Standby teaming option and not Active-Active teaming. So make sure when configuring Explicit failover the failback option is changed from default "yes" to "no". This change makes sure that you don't keep moving the traffic from one link to another if one link is flapping.

## vSphere Switch (VSS/VDS) Operation and Connection Options

VSPHERE SWITCH OPERATION-MAC LEARNING VSS/VDS 1

Virtual switches (VSS/VDS) are not learning switches. They authoritatively configure the MAC forwarding table.

**Forwarding Table Example**

| VMI MAC | Port |
|---------|------|
| VN1 MAC | 1    |
| VN2 MAC | 4    |

**Teaming configuration** decides which VM traffic will be sent over which uplink

For Example: Port ID hash teaming

- Based on the Port ID hash teaming algorithm VM1 traffic will be sent over uplink1 (virtual port 1)
- Based on the Port ID hash teaming algorithm VM2 traffic will be sent over uplink2 (virtual port 1)

VSPHERE SWITCH OPERATION-HANDLING OF BROADCAST PACKETS VSS/VDS 2

VSS/VDS doesn't support Spanning Tree Protocol (STP) because they don't create loops

When VMI sends a broadcast packet that packet is sent to all virtual ports and only one uplink port

When a broadcast packet is received from uplink port that packet is copied only to the VMs and not sent back to the other uplink port

The unique handling of Broadcast packets along with the Authoritative learning prevents loop in the network.

VIRTUAL SWITCH TO PHYSICAL SWITCH CONNECTION OPTIONS VSS/VDS 3

VSS TO VDS MIGRATION VSS/VDS 4

**Three methods of Migration**

- Using VDS wizard available in vCenter Server
- Using combination of the VDS and Host Profiles
- Using PowerCLI or vCLI commands

**Key Steps while migrating using VDS wizard in vCenter Server**

- Create a VDS first. Don't add the hosts to the VDS yet
- Next create the distributed virtual port groups to match the VSS port group configurations
- Add hosts to VDS
- Select the hosts you wish to migrate, along with physical adapters per host. You can either decide to migrate all physical adapters to VDS at once or few at a time
- Next step is to migrate the vmmnics of virtual adapters from the existing VSS port groups to the VDS distributed port groups
- After migrating vmmnics interfaces you can choose to migrate virtual machine networking from VSS to VDS

**Best Practices**

- Using Ether Channel configuration on physical switches make sure to disable it before going through migration steps

**Migration Guide**

- http://www.vmware.com/files/pdf/vsphere-network-vm-migration-configuration-wp.pdf

Components, Terminology and Data Plane Details

VSPHERE DISTRIBUTED SWITCH (VDS) COMPONENTS VDS 1

**Key differences between VSS and VDS**

- Management Plane and Data Plane part of the host in VSS (as shown in VSS 1)
- In VDS, the data plane remains local to each host, but the management plane is centralized with vCenter Server acting as the central control plane for all parameter configurations and virtual network management (as shown in VDS 1)

**Key advantages of VDS**

- Centralized management minimizes the configuration errors that could happen while managing VSS on individual hosts with separate management plane.
- Each vCenter Server instance can support up to 128 VDSs and each VDS can connect up to 500 hosts.
- Along with centralized and simplified management, VDS provides advanced virtual network capabilities.

**One Management Plane - Allow to configure various parameters of the distributed switch Data Plane. Handle the packet switching function on individual hosts**

VDS SPECIFIC TERMINOLOGY VDS 2

**dvUplink**

dvUplink provides a level of abstraction for the physical NICs (vmmnic) on each host. NIC teaming, load balancing, and failover policies on distributed port groups are applied to the dvUplinks and not to the vmmnics on individual hosts.

**dvUplink Port group**

While creating a new distributed switch the number of dvUplinks are defined as part of the dvUplink port group configuration.

**Distributed Port Group (dvPG)**

Distributed Port Group are port groups associated with a VDS and specify configuration that is common across a group of distributed virtual ports. Distributed Port Groups also define how a connection is made through the VDS to the Network.

**Distributed Virtual Port (dvPort)**

These are the ports where virtual machines or vmnics are connected to the virtual switch. Type of port binding is Static, Ephemeral and Dynamic (deprecated option and won't be available in future release).

VDS PARAMETERS - DVUPLINK VDS 4

**dvUplink Configuration** is done at dvUplink port group level

Depending on how many physical NICs are on hosts you can determine how many dvUplinks you should configure.

The following are some examples of dvUplink configurations

- If there are 4 hosts with 4 NICs each - Configure dvUplink as 4 (default)
- In a heterogeneous environment where there are 2 hosts with 6 NICs and 2 other hosts with 8 NICs - Configure dvUplink as 8 (Highest common denominator)

**Best Practices**

- dvUplinks have to be mapped to vmnics on the hosts. This mapping process happens when the hosts are added to the distributed switch. You should do consistent mapping of dvUplink to vmmic across all the hosts.
- In case of host that does not have enough vmnics when compared to uplinks, you can leave the higher order dvUplinks not mapped. In this deployment you should make sure that the dvUplink group on the host don't have teaming property that includes unmapped dvUplinks.

## vSphere Distributed Switch (VDS)

VDS PARAMETERS - DVPORT VDS 5

**dvPort configuration** is done at dvPort group level. You have following three options under the general dvPort group properties

- Port binding type
  - Static binding (Default)
  - Dynamic binding (deprecated and won't be supported in future release)
  - Ephemeral - No binding (Equivalent to Standard switch option)
- Port allocation
  - Static (Default)
  - Elastic
- Number of ports - Required if using Fixed port allocation option

**Best Practices**

- Use Static binding and Elastic port allocation - The static port-binding configuration on a dvPort group helps users to do detailed monitoring of dvPorts. This is not possible with ephemeral and dynamic port-binding configurations, where users lose the visibility and troubleshooting capability.
- With elastic port allocation you don't have to manually manage the number of dvPorts.

VDS PARAMETERS - VLAN VDS 6

**VLAN Configuration** is done at distributed port group level

- No VLAN tagging
- VLAN - Allow you to specify the VLAN ID to use for tagging. This is also called as Virtual Switch Tagging (VST) mode
- VLAN trunking - VLAN trunk range is configure from 0-4094. In this mode Guest virtual machines are allowed to do the tagging and those tags are carried through VDS. This is also called as Virtual Guest Tagging (VGT)
- PVLAN - If PVLANs are configured at the VDS level this option allows to choose the PVLAN for the port group.

**Best Practices**

- Use VLANs to isolate different traffic types (i.e. use Virtual Switch Tagging (VST))
- Make sure on the physical switch port, where the host is connected, all the configured VLANs are trunked
- PVLAN configuration can be used in creating Demilitarized Zones (DMZ) deployments.

VDS PARAMETERS - TEAMING VDS 7

**Teaming Configuration on a distributed Port Group depends on the way uplinks of a virtual switch are connected to the physical switch.**

Based on the physical switch connection option (shown in VSS/VDS 3) choose the load balancing type

- Load balancing types
  - Route based on originating virtual port (Port ID Hash)
  - Route based on source MAC hash (MAC Hash)
  - Route based on IP hash (IP Hash)
  - Use Explicit Failover order (Explicit Failover)
- After choosing the Load balancing type
  - Make sure that the Active, Standby and Unused adapters are identified for the port group

**Best Practices**

- Route based on physical NIC load (LBT) is the recommended option for teaming
- It doesn't need any special configuration on the physical switch side and utilizes the vmmics efficiently.

VDS PARAMETERS - TEaming VDS 7

**Teaming Configuration on a distributed Port Group depends on the way uplinks of a virtual switch are connected to the physical switch.**

Based on the physical switch connection option (shown in VSS/VDS 3) choose the load balancing type

- Load balancing types
  - Route based on originating virtual port (Port ID Hash)
  - Route based on source MAC hash (MAC Hash)
  - Route based on IP hash (IP Hash)
  - Use Explicit Failover order (Explicit Failover)
- After choosing the Load balancing type
  - Make sure that the Active, Standby and Unused adapters are identified for the port group

**Best Practices**

- Route based on physical NIC load (LBT) is the recommended option for teaming
- It doesn't need any special configuration on the physical switch side and utilizes the vmmics efficiently.

## vSphere Distributed Switch Features

VDS FEATURE - PRIVATE VLAN (PVLAN) VDS 8

**Private VLAN feature** provides isolation between the ports in the same broadcast domain. If you have limited VLANs or want to provide further segmentation in a broadcast domain then you can make use of this feature.

A virtual machine's Ingress and Egress network bandwidth can be controlled by enabling the network traffic shaper.

The following three parameters dictate the traffic-shaping policy on a virtual port

- Average Bandwidth - Specifies the allowed average bandwidth in kbps
- Peak Bandwidth - The maximum bandwidth allowed in kbps
- Burst Size - Establishes the maximum number of kilo bytes to allow in a burst

**Best Practices**

- Use this feature in the following scenarios
  - When you want more granular control over traffic type that can't be achieved through the NICOC feature.
  - You want to limit input traffic to a virtual machine or to a vmmnics NIC.

VDS FEATURE - TRAFFIC MANAGEMENT - PORT LEVEL TRAFFIC MANAGEMENT VDS 11

**NetFlow** is a networking protocol that collects IP traffic information as records and sends them to a collector tools for traffic flow analysis.

NetFlow features on VDS allow you to monitor the virtual machine to virtual machine flows that are not seen on the physical switch uplinks.

VDS support NetFlow V2 also called as IPFIX

**NetFlow Use cases**

- Help to help monitor application flows and measures flow performance over time. It also helps in network capacity planning exercises.
- Make use of NetFlow to measure the bandwidth requirement for different traffic types and then use that information to configure the NICOC shares and limits parameter.

The following diagram shows an example how VDS sends the network flow information across to a Collector Tool placed centrally in the network.

## VDS Configuration Parameters

VDS PARAMETERS - DVUPLINK VDS 4

**dvUplink Configuration** is done at dvUplink port group level

Depending on how many physical NICs are on hosts you can determine how many dvUplinks you should configure.

The following are some examples of dvUplink configurations

- If there are 4 hosts with 4 NICs each - Configure dvUplink as 4 (default)
- In a heterogeneous environment where there are 2 hosts with 6 NICs and 2 other hosts with 8 NICs - Configure dvUplink as 8 (Highest common denominator)

**Best Practices**

- dvUplinks have to be mapped to vmnics on the hosts. This mapping process happens when the hosts are added to the distributed switch. You should do consistent mapping of dvUplink to vmmic across all the hosts.
- In case of host that does not have enough vmnics when compared to uplinks, you can leave the higher order dvUplinks not mapped. In this deployment you should make sure that the dvUplink group on the host don't have teaming property that includes unmapped dvUplinks.

## Virtual Extensible Local Area Network (VXLAN)

Components, Terminology and Data Plane Details

VXLAN VXLAN 1

The following diagram shows VXLAN deployment within a data center where two clusters are deployed on two different racks and are backed by different VLANs.

VXLAN SPECIFIC TERMINOLOGY VXLAN 2

**Overlay Network**

It is a network built over another network. For Example, VXLAN network is built over IP network.

**Encapsulation and Decapsulation**

Process of adding and removing packet header is called as Encap and Decap respectively

**VTEP**

"Virtual Tunnel Endpoints" refers to the encapsulation/decapsulation endpoints of a VXLAN tunnel. In the vSphere environment the vmmnics module of VDS on each host acts as VTEP.

**Virtual Wire or VXLAN Segment ID**

A logical layer-2 network identified by a unique 24 bit segment ID.

**MTU**

"Maximum Transmission Unit" of a communications protocol of a layer is the size (in bytes) of the largest protocol data unit that the layer can pass onwards.

**IGMP**

"Internet Group Management Protocol" communications protocol is used by hosts and adjacent routers to establish multicast group memberships. It is analogous to ICMP for used by connections.

**IGMP Snooping**

A network switch keeps a map/table of IGMP conversation between hosts and routers to allow filtering of unneeded Multicasts for particular lines.

**LACP**

Link Aggregation Control Protocol provides a method to control the bundling of several physical ports together to form a single channel.

VXLAN PACKET HEADER DETAILS VXLAN 3

**Ethernet II IP - overlay network**

- Entire L2 frame encapsulated in UDP frame
- 50 bytes of overhead

**Included 24 bit VXLAN Identifier**

- 24-bit VNI/24-bit logical identifiers

**VXLAN can cross Layer 3**

Tunnels between vSphere hosts

- VNIs do NOT see VXLAN ID
- External switch don't see VMs IP and MAC address

**IP multicast used for L2 broadcast, unknown unicast technology submitted to IETF for standardization**

- With Cisco, Dell, Red Hat, Broadcom, Arista and Others

VXLAN Traffic Flows

VXLAN COMMUNICATION PATHS - TWO VMs ON SAME VXLAN VIRTUAL WIRE VXLAN 4

The following diagram shows VXLAN traffic flow between virtual machines and to the external world

VXLAN COMMUNICATION PATHS - TWO VMs ON DIFFERENT VXLAN VIRTUAL WIRES VXLAN 5

The following diagram shows VXLAN traffic flow between virtual machines on two different virtual wires

VXLAN ADVANTAGES VXLAN 6

Provides the ability to provision on-demand logical layer 2 isolated networks

- More multitenant groups are better
- Multiple segments can be mapped to a single multitenant group
- If VXLAN transport is contained in a single VLAN, "IGMP Querier" must be enabled for the VXLAN on the switch
- Multiple segments can be mapped to a single multitenant group
- Multiple segments can be mapped to a single multitenant group
- Multiple segments can be mapped to a single multitenant group

**Allows customer the flexibility to provision compute resources across layer 2 boundaries.**

For example, if there is a rack in the datacenter that has run out of compute capacity and there is another rack that is backed by different VLANs with some free compute resources, you can now easily increase the compute resources for your application by extending your logical layer 2 network over the rack where compute capacity is available.

**On demand networks without physical network re-configuration**

If customers want to deploy a new application they don't have to request the network admins to plumb new VLANs through their physical network infrastructure.

VXLAN INFRASTRUCTURE PREREQUISITES VXLAN 7

**IP Multicast forwarding is required**

- If VXLAN transport is contained in a single VLAN, "IGMP Querier" must be enabled for the VXLAN on the switch
- Multiple segments can be mapped to a single multitenant group
- Multiple segments can be mapped to a single multitenant group

**Increased MTU needed to accommodate VLAN encapsulation overhead**

- Physical infrastructure must carry 50 bytes more than the VM NIC MTU size.
- E.g. 1500 MTU on VMs - 50 + 1500 MTU on switches and routers.

**Leverage 5-tuple hash distribution for uplink and inter-switch LACP**

- Encapsulation will generate a source UDP port based on a hash of the inner packet 5-tuple (header information)

**If VXLAN traffic is traversing a router, proxy ARP must be enabled on first hop router**