# FlexPod Datacenter with VMware vSphere 6.5 U1 and Cisco ACI 3.1 Design Guide

Last Updated: September 10, 2018

# About the Cisco Validated Design Program

The Cisco Validated Design (CVD) program consists of systems and solutions designed, tested, and documented to facilitate faster, more reliable, and more predictable customer deployments. For more information, visit:

http://www.cisco.com/go/designzone.

ALL DESIGNS, SPECIFICATIONS, STATEMENTS, INFORMATION, AND RECOMMENDATIONS (COLLECTIVELY, "DESIGNS") IN THIS MANUAL ARE PRESENTED "AS IS," WITH ALL FAULTS.  CISCO AND ITS SUPPLIERS DISCLAIM ALL WARRANTIES, INCLUDING, WITHOUT LIMITATION, THE WARRANTY OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT OR ARISING FROM A COURSE OF DEALING, USAGE, OR TRADE PRACTICE.  IN NO EVENT SHALL CISCO OR ITS SUPPLIERS BE LIABLE FOR ANY INDIRECT, SPECIAL, CONSEQUENTIAL, OR INCIDENTAL DAMAGES, INCLUDING, WITHOUT LIMITATION, LOST PROFITS OR LOSS OR DAMAGE TO DATA ARISING OUT OF THE USE OR INABILITY TO USE THE DESIGNS, EVEN IF CISCO OR ITS SUPPLIERS HAVE BEEN ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

THE DESIGNS ARE SUBJECT TO CHANGE WITHOUT NOTICE.  USERS ARE SOLELY RESPONSIBLE FOR THEIR APPLICATION OF THE DESIGNS.  THE DESIGNS DO NOT CONSTITUTE THE TECHNICAL OR OTHER PROFESSIONAL ADVICE OF CISCO, ITS SUPPLIERS OR PARTNERS.  USERS SHOULD CONSULT THEIR OWN TECHNICAL ADVISORS BEFORE IMPLEMENTING THE DESIGNS.  RESULTS MAY VARY DEPENDING ON FACTORS NOT TESTED BY CISCO.

CCDE, CCENT, Cisco Eos, Cisco Lumin, Cisco Nexus, Cisco StadiumVision, Cisco TelePresence, Cisco WebEx, the Cisco logo, DCE, and Welcome to the Human Network are trademarks; Changing the Way We Work, Live, Play, and Learn and Cisco Store are service marks; and Access Registrar, Aironet, AsyncOS, Bringing the Meeting To You, Catalyst, CCDA, CCDP, CCIE, CCIP, CCNA, CCNP, CCSP, CCVP, Cisco, the Cisco Certified Internetwork Expert logo, Cisco IOS, Cisco Press, Cisco Systems, Cisco Systems Capital, the Cisco Systems logo, Cisco Unified Computing System (Cisco UCS), Cisco UCS B-Series Blade Servers, Cisco UCS C-Series Rack Servers, Cisco UCS S-Series Storage Servers, Cisco UCS Manager, Cisco UCS Management Software, Cisco Unified Fabric, Cisco Application Centric Infrastructure, Cisco Nexus 9000 Series, Cisco Nexus 7000 Series. Cisco Prime Data Center Network Manager, Cisco NX-OS Software, Cisco MDS Series, Cisco Unity, Collaboration Without Limitation, EtherFast, EtherSwitch, Event Center, Fast Step, Follow Me Browsing, FormShare, GigaDrive, HomeLink, Internet Quotient, IOS, iPhone, iQuick Study, LightStream, Linksys, MediaTone, MeetingPlace, MeetingPlace Chime Sound, MGX, Networkers, Networking Academy, Network Registrar, PCNow, PIX, PowerPanels, ProConnect, ScriptShare, SenderBase, SMARTnet, Spectrum Expert, StackWise, The Fastest Way to Increase Your Internet Quotient, TransPath, WebEx, and the WebEx logo are registered trademarks of Cisco Systems, Inc. and/or its affiliates in the United States and certain other countries.

All other trademarks mentioned in this document or website are the property of their respective owners. The use of the word partner does not imply a partnership relationship between Cisco and any other company. (0809R)

# Table of Contents

# Executive Summary

Cisco Validated Designs consist of systems and solutions that are designed, tested, and documented to facilitate and improve customer deployments. These designs incorporate a wide range of technologies and products into a portfolio of solutions that have been developed to address the business needs of our customers.

This document describes the Cisco and NetApp® FlexPod® solution, which is a validated approach for deploying Cisco and NetApp technologies as a shared cloud infrastructure. FlexPod is a leading integrated infrastructure supporting a broad range of enterprise workloads and use cases. This validated design provides a framework for deploying a VMware vSphere virtualization platform based on FlexPod in enterprise class datacenters. This solution enables customers to deploy VMware vSphere based private cloud on integrated infrastructure quickly and reliably.

The solution architecture is based on Cisco UCS running a software release that supports all Cisco UCS hardware platforms including Cisco UCS B-Series Blade Servers and Cisco UCS C-Series Rack-Mount Servers, Cisco UCS 6300 or 6200 Fabric Interconnects, Cisco Nexus 9000 Series Switches, and NetApp All Flash series storage arrays. In addition to that, it includes VMware vSphere 6.5 U1, which provides a number of new features for optimizing storage utilization and facilitating private cloud.

# FlexPod Program

Cisco and NetApp® have carefully validated and verified the FlexPod solution architecture and its many use cases while creating a portfolio of detailed documentation, information, and references to assist customers in transforming their data centers to this shared infrastructure model. This portfolio includes, but is not limited to the following items:

- Best practice architectural design

- Workload sizing and scaling guidance

- Implementation and deployment instructions

- Technical specifications (rules for what is a FlexPod® configuration)

- Frequently asked questions and answers (FAQs)

- Cisco Validated Designs (CVD) and NetApp Validated Architectures (NVA) that cover different use cases

Cisco and NetApp have also built a robust and experienced support team focused on FlexPod solutions, from customer account and technical sales representatives to professional services and technical support engineers. The support alliance between NetApp and Cisco gives customers and channel services partners direct access to technical experts who collaborate with cross vendors and have access to shared lab resources to resolve potential issues.

FlexPod supports tight integration with virtualized and cloud infrastructures, making it the logical choice for long-term investment. FlexPod also provides a uniform approach to IT architecture, offering a well-characterized and documented shared pool of resources for application workloads. FlexPod delivers operational efficiency and consistency with the versatility to meet a variety of SLAs and IT initiatives, including:

- Application rollouts or application migrations

- Business continuity and disaster recovery

- Desktop virtualization

- Cloud delivery models (public, private, hybrid) and service models (IaaS, PaaS, SaaS)

- Asset consolidation and virtualization

# Solution Overview

## Introduction

Industry trends indicate a vast data center transformation toward shared infrastructure and cloud computing. Business agility requires application agility, so IT teams need to provision applications quickly and resources need to be able to scale up (or down) in minutes.

FlexPod Datacenter is a best practice data center architecture, designed and validated by Cisco and NetApp to meet the needs of enterprise customers and service providers. FlexPod Datacenter is built on NetApp All Flash FAS, Cisco Unified Computing System (UCS), and the Cisco Nexus family of switches. These components combine to enable management synergies across all IT infrastructure in a business. FlexPod Datacenter has been proven to be the optimal platform for virtualization and workload consolidation, enabling enterprises to standardize their entire IT infrastructure.

To simplify the evolution to a shared cloud infrastructure based on an application driven policy model, Cisco and NetApp have developed this FlexPod with Cisco Application Centric Infrastructure (ACI) solution for VMware vSphere environments. Cisco ACI in the data center is a holistic architecture with centralized automation and policy-driven application profiles that combines software flexibility with hardware performance.

## Audience

The audience for this document includes, but is not limited to; sales engineers, field consultants, professional services, IT managers, partner engineers, and customers who want to take advantage of an infrastructure built to deliver IT efficiency and enable IT innovation.

## What's New?

The following design elements distinguish this version of FlexPod from previous models:

- Validation of the Cisco ACI 3.1 Release on Cisco Nexus 9000 Series Switches

- Support for the Cisco UCS 3.2 release and Cisco UCS B200-M5 servers

- Support for the latest release of NetApp Data ONTAP® 9.3

- Support for VMware vSphere 6.5 Update 1 on the above infrastructure

- A storage design supporting NFS, iSCSI and Fiber Channel SAN

- Support for FC/FCoE storage by directly connecting to Cisco UCS Fabric Interconnects

- Application design guidance for multi-tiered applications using Cisco ACI application profiles and policies

## FlexPod Datacenter Overview

FlexPod Datacenter is a validated, converged infrastructure platform that is designed and built for Enterprise datacenters and cloud deployments. FlexPod Datacenter solutions use the following family of components for the compute, networking and storage layers of the design:

- Cisco Unified Computing System (Cisco UCS)

- Cisco Nexus and MDS Family of Switches

- NetApp All Flash FAS (AFF) Storage Systems

These components are connected, configured and integrated based on technology and product best practices recommended by both Cisco and NetApp to deliver a highly available, flexible and scalable platform for running a variety of enterprise workloads with confidence. FlexPod can scale up for greater performance and capacity (adding compute, network, or storage resources individually as needed), or it can scale out for environments that require multiple consistent deployments (such as rolling out of additional FlexPod stacks). FlexPod solutions address the four primary architectural goals of scalability, flexibility, availability, and manageability.

Performance is a key design criterion that is not addressed in this document. It is addressed through other collateral, benchmarking and solution testing efforts; this design guide focusses on the functional design.

Figure 1    FlexPod Component Families



One of the key benefits of FlexPod is its ability to maintain consistency during scale. Each of the component families shown (Cisco UCS, Cisco Nexus, and NetApp AFF) offers platform and resource options to scale the infrastructure up or down, while supporting the same features and functionality that are required under the configuration and connectivity best practices of FlexPod.

Specific models from the above component families are used to validate the reference architecture covered in this document. See the Solution Validation section of this document for more details.

FlexPod Datacenter is a flexible solution that can support other models to meet customer require-ments. Confirm support using Cisco's Hardware Compatibility List (HCL) and NetApp's Interoperability Matrix (IMT).

# Solution Components Overview

This section provides a technical overview of the compute, network, storage and management components in this FlexPod Datacenter solution. For additional information on any of the components covered in this section, see Solution References section.

## Cisco Unified Computing System

The Cisco Unified Computing System™ (Cisco UCS) is a next-generation data center platform that integrates computing, networking, storage access, and virtualization resources into a cohesive system designed to reduce total cost of ownership and increase business agility. The system integrates a low-latency, lossless 10 or 40 Gigabit Ethernet unified network fabric with enterprise-class, x86-architecture servers. The system is an integrated, scalable, multi-chassis platform with a unified management domain for managing all resources.

The Cisco Unified Computing System consists of the following subsystems:

- Compute - **The compute piece of the system incorporates servers based on latest Intel's x86** processors. Servers are available in blade and rack form factor, managed by Cisco UCS Manager.

- Network - The integrated network fabric in the system provides a low-latency, lossless, 10/40 Gbps Ethernet fabric. Networks for LAN, SAN and management access are consolidated within the fabric. The unified fabric uses the innovative Single Connect technology to lowers costs by reducing the number of network adapters, switches, and cables. This in turn lowers the power and cooling needs of the system.

- Virtualization - The system unleashes the full potential of virtualization by enhancing the scalability, performance, and operational control of virtual environments. Cisco security, policy enforcement, and diagnostic features are now extended into virtual environments to support evolving business needs.

- Storage access – Cisco UCS system provides consolidated access to both SAN storage and Network Attached Storage over the unified fabric. This provides customers with storage choices and investment protection. Also, the server administrators can pre-assign storage-access policies to storage resources, for simplified storage connectivity and management leading to increased productivity.

- Management: The system uniquely integrates compute, network and storage access subsystems, enabling it to be managed as a single entity through Cisco UCS Manager software. Cisco UCS Manager increases IT staff productivity by enabling storage, network, and server administrators to collaborate on Service Profiles that define the desired physical configurations and infrastructure policies for applications. Service Profiles increase business agility by enabling IT to automate and provision resources in minutes instead of days.

### Cisco UCS Differentiators

**Cisco's Unified Compute System** has revolutionized the way servers are managed in data-center and provides a number of unique differentiators that are outlined below:

- Embedded Management – Servers in Cisco UCS are managed by embedded software in the Fabric Interconnects, eliminating need for any external physical or virtual devices to manage the servers.

- Unified Fabric – Cisco UCS uses a wire-once architecture, where a single Ethernet cable is used from the FI from the server chassis for LAN, SAN and management traffic. Adding compute capacity does not require additional connections. This converged I/O reduces overall capital and operational expenses.

- Auto Discovery – By simply inserting a blade server into the chassis or a rack server to the fabric interconnect, discovery of the compute resource occurs automatically without any intervention.

- Policy Based Resource Classification – Once a compute resource is discovered, it can be automatically classified to a resource pool based on policies defined which is particularly useful in cloud computing.

- Combined Rack and Blade Server Management – Cisco UCS Manager is hardware form factor agnostic and can manage both blade and rack servers under the same management domain.

- Model based Management Architecture – Cisco UCS Manager architecture and management database is model based and data driven. An open XML API is provided to operate on the management model which enables easy and scalable integration of Cisco UCS Manager with other management systems.

- Policies, Pools, and Templates – The management approach in Cisco UCS Manager is based on defining policies, pools and templates, instead of cluttered configuration, which enables a simple, loosely coupled, data driven approach in managing compute, network and storage resources.

- Policy Resolution – In Cisco UCS Manager, a tree structure of organizational unit hierarchy can be created that mimics the real-life tenants and/or organization relationships. Various policies, pools and templates can be defined at different levels of organization hierarchy.

- Service Profiles and Stateless Computing – A service profile is a logical representation of a server, carrying its various identities and policies. This logical server can be assigned to any physical compute resource as far as it meets the resource requirements. Stateless computing enables procurement of a server within minutes, which used to take days in legacy server management systems.

- Built-in Multi-Tenancy Support – The combination of a profiles-based approach using policies, pools and templates and policy resolution with organizational hierarchy to manage compute resources makes Cisco UCS Manager inherently suitable for multi-tenant environments, in both private and public clouds.

## Cisco UCS Manager

Cisco UCS Manager (UCSM) provides unified, integrated management for all software and hardware components in Cisco UCS. Using Cisco Single Connect technology, it manages, controls, and administers multiple chassis for thousands of virtual machines. Administrators use the software to manage the entire Cisco Unified Computing System as a single logical entity through an intuitive graphical user interface (GUI), a command-line interface (CLI), or a through a robust application programming interface (API).

Cisco UCS Manager is embedded into the Cisco UCS Fabric Interconnects and provides a unified management interface that integrates server, network, and storage. Cisco UCS Manger performs auto-discovery to detect inventory, manage, and provision system components that are added or changed. It offers comprehensive set of XML API for third party integration, exposes thousands of integration points and facilitates custom development for automation, orchestration, and to achieve new levels of system visibility and control.

## Cisco UCS Fabric Interconnects

The Cisco UCS Fabric interconnects (FIs) provide a single point for connectivity and management for the entire Cisco UCS system. Typically deployed as an active-active pair, the system's fabric interconnects integrate all components into a single, highly-available management domain controlled by the Cisco UCS Manager. Cisco UCS FIs provide a single unified fabric for the system, with low-latency, lossless, cut-through switching that supports LAN, SAN and management traffic using a single set of cables.

### Cisco UCS 6200 Series Fabric Interconnects

Cisco UCS 6200 Series FIs are a family of line-rate, low-latency, lossless switches that can support 1/10 Gigabit Ethernet and Fibre Channel over Ethernet (FCoE), or native (4/2/1 or 8/4/2) Fibre Channel (FC) connectivity.

### Cisco UCS 6300 Series Fabric Interconnects

Cisco UCS 6300 Series Fabric Interconnects provide a 40Gbps unified fabric with higher workload density and double the switching capacity of Cisco UCS 6200 FIs. When combined with the newer 40Gbps Cisco UCS 2300 Series Fabric Extenders, they provide 40GbE / FCoE port connectivity to enable end-to-end 40GbE / FCoE solution. The unified ports also support 16G FC for high speed FC connectivity to SAN.

The two 6300 Fabric Interconnect models currently available are:

- Cisco UCS 6332 Fabric Interconnect is a 1RU 40GbE and FCoE switch offering up to 2.56 Tbps of throughput. The switch has a density of 32 x 40Gbps Ethernet, and FCoE ports. This model is aimed at IP storage deployments requiring high-performance 40Gbps FCoE connectivity to Cisco MDS switches.



- Cisco UCS 6332-16UP Fabric Interconnect is a 1RU 40GbE/FCoE switch and 1/10 Gigabit Ethernet, FCoE and FC switch offering up to 2.24 Tbps throughput. The switch has 24x40Gbps fixed Ethernet/FCoE ports and 16x1/10Gbps Ethernet/FCoE or 4/8/16Gbps FC ports. This model is aimed at FC storage deployments requiring high performance 16Gbps FC connectivity to Cisco MDS switches.



Table 1  provides a comparison of the port capabilities of the different Fabric Interconnect models.

Table 1   Cisco UCS 6200 and 6300 Series Fabric Interconnects

| Features | 6248 | 6296 | 6332 | 6332-16UP |
|---|---|---|---|---|

| Features | 6248 | 6296 | 6332 | 6332-16UP |
|---|---|---|---|---|
| Max 10G ports | 48 | 96 | 96* + 2** | 72* + 16 |
| Max 40G ports | - | - | 32 | 24 |
| Max unified ports | 48 | 96 | - | 16 |
| Max FC ports | 48 x 2/4/8G FC | 96 x 2/4/8G FC | - | 16 x 4/8/16G FC |

*Using 40G to 4x10G breakout cables*        **Requires QSA module*

See the Solution References section for additional information on Fabric Interconnects.

## Cisco UCS 5108 Blade Server Chassis

The Cisco UCS 5108 Blade Server Chassis is a fundamental building block of the Cisco Unified Computing System, delivering a scalable and flexible blade server architecture. The Cisco UCS blade server chassis uses an innovative unified fabric with fabric-extender technology to lower TCO by reducing the number of network interface cards (NICs), host bus adapters (HBAs), switches, and cables that need to be managed, cooled, and powered. It is a 6-RU chassis that can house up to 8 x half-width or 4 x full-width Cisco UCS B-series blade servers. A passive mid-plane provides up to 80Gbps of I/O bandwidth per server slot and up to 160Gbps for two slots (full-width). The rear of the chassis contains two I/O bays to house a pair of Cisco UCS 2000 Series Fabric Extenders to enable uplink connectivity to FIs for both redundancy and bandwidth aggregation.

**Figure 2    Cisco UCS 5108 Blade Server Chassis**



See Solution References section for additional information on Cisco UCS Blade Server Chassis.

## Cisco UCS Fabric Extenders

The Cisco UCS Fabric extenders (FEX) or I/O Modules (IOMs) multiplexes and forwards all traffic from servers in a blade server chassis to a pair of Cisco UCS Fabric Interconnects over a 10Gbps or 40Gbps unified fabric links. All traffic, including traffic between servers on the same chassis, or between virtual machines on the same server, is forwarded to the parent fabric interconnect where Cisco UCS Manager runs, managing the profiles and polices for the servers. FEX technology was developed by Cisco. Up to two FEXs can be deployed in a chassis.

For more information about the benefits of FEX, see: http://www.cisco.com/c/en/us/solutions/data-center-virtualization/fabric-extender-technology-fex-technology/index.html

## Cisco UCS B200 M5 Servers

The enterprise-class Cisco UCS B200 M5 Blade Server extends the Cisco UCS portfolio in a half-width blade form-factor. This M5 server uses the latest Intel® Xeon® Scalable processors with up to 28 cores per processor, 3TB of RAM (using 24 x128GB DIMMs), 2 drives (SSD, HDD or NVMe), 2 GPUs and 80Gbps of total I/O to each server. It supports the Cisco VIC 1340 adapter to provide 40Gb FCoE connectivity to the unified fabric. For more information on Cisco UCS B-series servers, see Solution References.

## Cisco UCS C220 M5 Servers

The enterprise-class Cisco UCS C220 M5 server extends the Cisco UCS portfolio in a one rack-unit (1RU) form-factor. This M5 server uses the latest Intel® Xeon® Scalable processors with up to 28 cores per processor (1 or 2), 3TB of RAM (with 24 x 128GB DIMMS), with various drive combinations (HDD, SSD or NVMe), 77TB of storage capacity (using 10 x 2.5in NVMe PCIe SSDs), and 80Gbps of throughput connectivity. It supports the Cisco VIC 1385 and 1387 adapters (also supports other NICs and HBAs) to provide 40Gb Ethernet and FCoE connectivity. For more information on Cisco UCS C-series servers, see Solution References.

## Cisco UCS Network Adapters

Cisco UCS supports converged network adapters (CNAs) to provide connectivity to the blade and rack mount servers. CNAs eliminate the need for multiple network interface cards (NICs) and host bus adapters (HBAs) by converging LAN, SAN and Management traffic to a single interface. While Cisco UCS supports a wide variety of interface cards, the recommended adapters currently available for blade and rack servers are outlined below.

## Cisco UCS Virtual Interface Card 1340

Cisco UCS Virtual Interface Card (VIC) 1340 is a 2 x 40Gbps Ethernet or 2 x 4 x 10-Gbps Ethernet and FCoE-capable modular LAN on motherboard (mLOM), designed for Cisco UCS B200 Series Blade Servers. The Cisco UCS VIC 1340 enables a policy-based, stateless, agile server infrastructure that can present over 256 PCIe standards-compliant interfaces to the host, and dynamically configure them as NICs or HBAs.

## Cisco UCS Virtual Interface Card 1385 and 1387

Cisco UCS VIC 1385 and VIC 1387 are different form factors (PCIe, MLOM) of the same adapter, designed for Cisco UCS C-series servers and support 2 x 40Gbps Ethernet and FCoE. Similar to Cisco UCS VIC 1340, these adapters incorporate **Cisco's CNA technology to enable a policy**-based, stateless, agile server that can present over 256 PCIe standards-compliant interfaces to the host, and can be dynamically configured as NICs or HBAs. For more information on Cisco UCS Network Adapters, see Solution References.

# Cisco Application Centric Infrastructure and Nexus Switching

**Cisco ACI is an evolutionary leap from SDN's initial vision of operational efficiency through network agility** and programmability. Cisco ACI has industry leading innovations in management automation, programmatic policies, and dynamic workload provisioning. The ACI fabric accomplishes this with a combination of hardware, policy-based control systems, and closely coupled software to provide advantages not possible in other architectures.

Cisco ACI takes a policy-based, systems approach to operationalizing the data center network. The policy is centered around the needs (reachability, access to services, security policies) of the applications. Cisco ACI delivers a resilient fabric to satisfy today's dynamic applications.

## Cisco ACI Architecture

The Cisco ACI fabric is a leaf-and-spine architecture where every leaf connects to every spine using high-speed 40/100-Gbps Ethernet links, with no direct connections between spine nodes or leaf nodes. The ACI fabric is a routed fabric with a VXLAN overlay network, where every leaf is VXLAN Tunnel Endpoint (VTEP). Cisco ACI provides both Layer 2 (L2) and Layer 3 (L3) forwarding across this routed fabric infrastructure.

Figure 3    Cisco ACI Fabric Architecture



## Architectural Building Blocks

The key architectural buildings blocks of the Cisco ACI fabric are:

- Application Policy Infrastructure Controller (APIC) - Cisco APIC is the unifying point of control in Cisco ACI for automating and managing the end-to-end data center fabric. The Cisco ACI fabric is built on a network of individual components that are provisioned and managed as a single entity. The APIC is a physical appliance that serves as a software controller for the overall fabric. It is based on Cisco UCS C-series rack mount servers with 2x10Gbps links for dual-homed connectivity to a pair of leaf switches and 1Gbps interfaces for out-of-band management. Cisco APIC optimizes the application lifecycle for scale and performance, and supports flexible application provisioning across physical and virtual resources. The Cisco APIC exposes northbound APIs through XML and JSON and provides both a command-line interface (CLI) and GUI that manage the fabric through the same APIs available programmatically. For control plane resiliency, a cluster of three APICs, dual-homed to a pair of leaf switches in the fabric are typically deployed.

- Cisco Nexus 9000 Series Switches – The Cisco ACI fabric is built on a network of Cisco Nexus 9000 series switches that provide low-latency, high-bandwidth connectivity with industry proven protocols and innovative technologies to create a flexible, scalable, and highly available architecture. ACI is supported on several models of Nexus 9000 series switches and line cards. The selection of a

Nexus 9000 series switch as an ACI spine or leaf switch will depend on a number of factors such as physical layer connectivity (1/10/25/40/50/100-Gbps), FEX aggregation support, analytics support in hardware (Cloud ASICs), FCoE support, link-level encryption, support for Multi-Pod, Multi-Site etc.

- Spine Nodes – The spines provide high-speed (40/100-Gbps) connectivity between leaf nodes. The ACI fabric forwards traffic by doing a host lookup in a mapping database that contains information about the leaf node where an endpoint (IP, Mac) resides. All known endpoints are maintained in a hardware database on the spine switches. The number of endpoints or the size of the database is a key factor in the choice of a Nexus 9000 model as a spine switch. Leaf switches also maintain a database but only for those hosts that send/receive traffic through it.

- Leaf Nodes – Leaf switches are essentially Top-of-Rack (ToR) switches that end devices connect into. They provide Ethernet connectivity to devices such as servers, firewalls, storage and other network elements. Leaf switches provide access layer functions such as traffic classification, policy enforcement, L2/L3 forwarding of edge traffic etc. The criteria for selecting a specific Nexus 9000 model as a leaf switch will be different from that of a spine switch.

There are currently two generations of ACI-capable Nexus 9000 switches. A mixed deployment is supported but it can impact the design and features available in the ACI fabric.
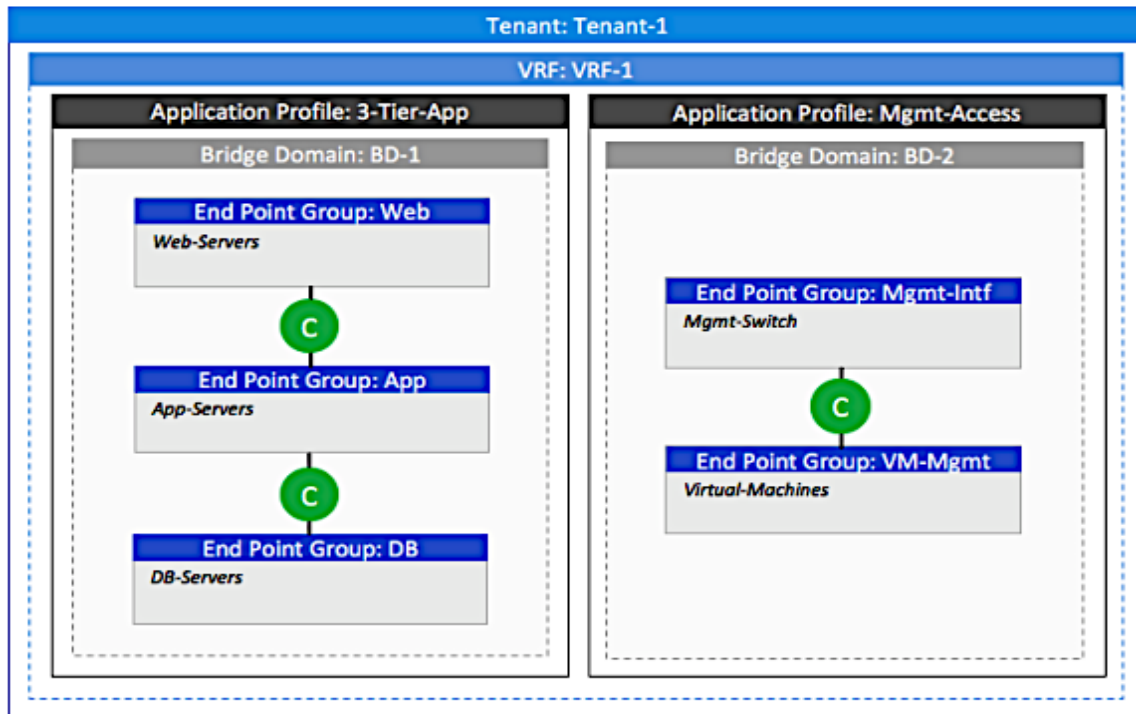
## Cisco ACI Fabric Design Constructs

Cisco ACI architecture uses a number of design constructs that are critical to an Cisco ACI based design and are summarized below. Figure 4 illustrates the relationship between various ACI elements.

- Tenant – A tenant is a logical container which can represent an actual tenant, organization, application or a construct for grouping. From a policy perspective, a tenant represents a unit of isolation. All application configurations in Cisco ACI are part of a tenant. Within a tenant, multiple VRF contexts, bridge domains, and EPGs can be defined according to application requirements.

- VRF – Tenants can be further divided into Virtual Routing and Forwarding (VRF) instances (separate IP spaces) to further separate the organizational and forwarding requirements for a given tenant. A Tenant can have multiple VRFs. IP addressing can be duplicated across VRFs for multitenancy.

- Bridge Domain – A bridge domain (BD) is a L2 forwarding construct that represents a broadcast domain within the fabric. A bridge domain is associated with a single tenant VRF but a VRF can have multiple bridge domains and endpoints. The endpoints in a BD can be anywhere in the ACI fabric, distributed across multiple leaf switches. To minimize flooding across the fabric, ACI provides several features such as learning of endpoint addresses (Mac/IP/Both), forwarding of ARP Requests directly to the destination leaf node, maintaining a mapping database of active remote conversations, local forwarding, and probing of endpoints before they expire. Subnet(s) can be associated with a BD to provide an L3 gateway to the endpoints in the BD.

- End Point Group – An End Point Group (EPG) is a collection of physical and/or virtual end points that require common services and policies, independent of their location. Endpoints could be physical servers, VMs, storage arrays, etc. For example, a Management EPG could be a collection of endpoints that connect to a common segment for management. Each EPG is associated with a single bridge domain but a bridge domain can have multiple EPGs mapped to it.

- Application Profile – An application profile (AP) models application requirements and contains one or more EPGs as necessary to provide the application capabilities. A Tenant can contain one or more application profiles and an application profile can contain one or more EPGs.

- Contracts – Contracts are rules and policies that define the interaction between EPGs. Contracts determine how applications use the network. Contracts are defined using provider-consumer relationships; one EPG provides a contract and another EPG(s) consumes that contract. Contracts utilize inbound/outbound filters to limit the traffic between EPGs or applications based EtherType, IP protocols, TCP/UDP port numbers and can specify QoS and L4-L7 redirect policies.

Figure 4    ACI Design Constructs - Relationship Between Major Components



As shown in Figure 4, devices in different EPGs talk to each other using contracts and associated filter but no special configuration is necessary within the same EPG. Also, each tenant can have multiple VRFs and bridge domains and each BD can be associated with multiple application profiles and EPGs.

## NetApp A-Series All Flash FAS Storage

With the new A-Series All Flash FAS (AFF) controller lineup, NetApp provides industry leading performance while continuing to provide a full suite of enterprise-grade data management and data protection features. The A-Series lineup offers double the IOPS while decreasing the latency. Table 2  summarizes the A-series controllers and their specifications. For more information on the AFF A-Series, see Solution References.

Table 2   NetApp A-Series Controller Specifications

|  | AFF A200 | AFF A300 | AFF A700 | AFF A700s |
|---|---|---|---|---|

| NAS Scale-out | 2-8 nodes | 2-24 nodes | 2-24 nodes | 2-24 nodes |
|---|---|---|---|---|
| SAN Scale-out | 2-8 nodes | 2-12 nodes | 2-12 nodes | 2-12 nodes |
| Per HA Pair Specifications (Active-Active Dual Controller) | | | | |
| Maximum SSDs | 144 | 384 | 480 | 216 |
| Maximum Raw Capacity | 2.2PB | 5.9PB | 7.3PB | 3.3PB |
| Effective Capacity | 8.8PB | 23.8PB | 29.7PB | 13PB |
| Chassis Form Factor | 2U chassis with two HA controllers and 24 SSD slots | 3U chassis with two HA controllers | 8u chassis with two HA controllers | 4u chassis with two HA controllers and 24 SSD slots |
| ONTAP 9 Base Bundle | ✓ | ✓ | ✓ | ✓ |
| ONTAP 9 Premium Bundle (FlexClone, SnapMirror, SnapVault, SnapCenter & more) | ✓ | ✓ | ✓ | ✓ |

NetApp A-Series all-flash controllers were designed from the ground up with flash storage in mind. They provide industry leading density, scalability, and network connectivity, allowing customers to do more with their flash storage. This controller provides a rich set of data management features as well as industry leading performance in a 2u form factor. ONTAP 9 has many key features that optimize SSD performance and endurance, including the following:

- Coalesced writes to free blocks to maximize flash performance and longevity

- Flash-specific read path optimizations to enable consistent low latency

- Advanced drive partitioning to increase storage efficiency, increasing usable capacity by almost 20%

- Support for multi-stream writes to increase write performance to SSDs

## NetApp AFF A300

The NetApp AFF A300 controller provides the high-performance benefits of 40GbE and all flash SSDs, offering better performance than comparable options, while taking up less space in the rack. Combined with the disk shelf of 3.8TB disks, this solution can provide over ample horsepower and over 90TB of capacity, all while taking up only 5U of valuable rack space. This makes it an ideal controller for a shared workload **converged infrastructure. As an infrastructure's capacity or performance needs grow, the NetApp AFF A300** can increase capacity with additional storage shelves and performance by adding additional controllers to the cluster; a cluster can scale up to 24 nodes.

Figure 5    NetApp A300 Front View



Figure 6    NetApp A300 Rear View



Note: The 40GbE cards are installed in the expansion slot 2 and the ports are e2a, e2e.

## Cluster Storage and HA Pairs

In NetApp storage architecture, each controller (FAS or AFF), its storage, its network connectivity, and the instance of ONTAP running on the controller is called a *node.*

Nodes are paired for high availability (HA). Together these pairs (up to 12 nodes for SAN, up to 24 nodes for NAS) comprise the cluster. Nodes communicate with each other over a private, dedicated cluster interconnect.

Nodes in an HA pair must use the same storage array model. However, a cluster can use any supported combination of controllers. You can scale out for capacity by adding nodes with like storage array models, or for performance by adding nodes with higher-end storage arrays.

An internal HA interconnect allows each node to continually check whether its partner is functioning and to **mirror log data for the other's nonvolatile memory. When a write request is made to a node, it is logged in** NVRAM on both nodes before a response is sent back to the client or host. On failover, the surviving partner commits the failed node's uncommitted write requests to disk, ensuring data consistency.

Depending on the controller model, node storage consists of flash disks, capacity drives, or both. Network ports on the controller provide access to data. Physical storage and network connectivity resources are virtualized, visible to cluster administrators only, not to NAS clients or SAN hosts.

**Connections to the other controller's storage media allow each node to access the other's storage in the** event of a takeover. Network path failover mechanisms ensure that clients and hosts continue to communicate with the surviving node.

Of course, you can scale up in all the traditional ways as well, upgrading disks or controllers as needed. ONTAP's virtualized storage infrastructure makes it easy to move data non-disruptively, meaning that you can scale vertically or horizontally without downtime.

# NetApp ONTAP 9.3

NetApp ONTAP data management software offers unified storage for applications that read and write data over block- or file-access protocols, in storage configurations that range from high-speed flash, to lower-priced spinning media, to cloud-based object storage.

ONTAP implementations run on NetApp-engineered Fabric-Attached Storage (FAS) or AFF appliances, on commodity hardware (ONTAP Select), and in private, public, or hybrid clouds (NetApp Private Storage, ONTAP Cloud). Specialized implementations offer best-in-class converged infrastructure (FlexPod Datacenter) and access to third-party storage arrays (FlexArray Virtualization).

Together these implementations form the basic framework of the *NetApp Data Fabric*, with a common software-defined approach to data management and fast, efficient replication across platforms. ONTAP serves as the foundation for hybrid cloud and virtualization designs.

## NetApp Storage Virtual Machine (SVM)

A cluster serves data through at least one and possibly multiple storage virtual machines (SVMs - formerly called Vservers). An SVM is a logical abstraction that represents the set of physical resources of the cluster. Data volumes and network logical interfaces (LIFs) are created and assigned to an SVM and may reside on any node in the cluster to which the SVM has been given access. An SVM may own resources on multiple nodes concurrently, and those resources can be moved non-disruptively from one node to another. For example, a flexible volume can be non-disruptively moved to a new node and aggregate, or a data LIF can be transparently reassigned to a different physical network port. The SVM abstracts the cluster hardware and it is not tied to any specific physical hardware.

An SVM can support multiple data protocols concurrently. Volumes within the SVM can be joined together to form a single NAS namespace, which makes all of an SVM's data available through a single share or mount point to NFS and CIFS clients. SVMs also support block-based protocols, and LUNs can be created and exported by using iSCSI, FC, or FCoE. Any or all of these data protocols can be configured for use within a given SVM.

Because it is a secure entity, an SVM is only aware of the resources that are assigned to it and has no knowledge of other SVMs and their respective resources. Each SVM operates as a separate and distinct entity with its own security domain. Tenants can manage the resources allocated to them through a delegated SVM administration account. Each SVM can connect to unique authentication zones such as Active Directory, LDAP, or NIS. A NetApp cluster can contain multiple SVMs. If you have multiple SVMs, you can delegate an SVM to a specific application. This allows administrators of the application to access only the dedicated SVMs and associated storage, increasing manageability, and reducing risk.
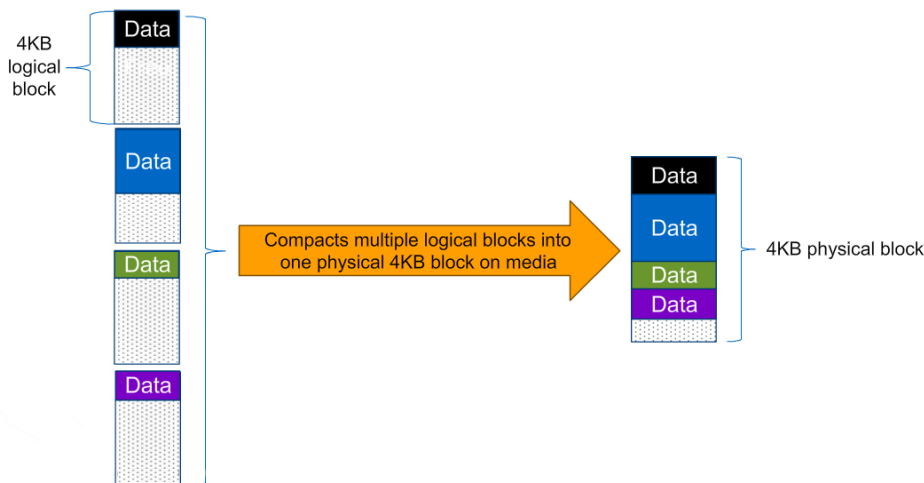
## Storage Efficiency

Storage efficiency has always been a primary architectural design point of ONTAP. A wide array of features allows businesses to store more data using less space. In addition to deduplication and compression,

businesses can store their data more efficiently by using features such as unified storage, multi-tenancy, and NetApp Snapshot® technology. Storage efficiency features with ONTAP 9 include:

- **Thin provisioning**: A *thin-provisioned* volume or LUN is one for which storage is not reserved in advance. Instead, storage is allocated dynamically, as it is needed. Free space is released back to the storage system when data in the volume or LUN is deleted.

- **Deduplication:** *Deduplication* reduces the amount of physical storage required for a volume (or all the volumes in an AFF aggregate) by discarding duplicate blocks and replacing them with references to a single shared block. Reads of deduplicated data typically incur no performance charge. Writes incur a negligible charge except on overloaded nodes.

- **Compression**: *Compression* reduces the amount of physical storage required for a volume by combining data blocks in *compression groups,* each of which is stored as a single block. Reads of compressed data are faster than in traditional compression methods because ONTAP decompresses only the compression groups that contain the requested data, not an entire file or LUN.

- **Compaction**: Compaction, which was introduced in ONTAP 9, is the latest patented storage efficiency technology released by NetApp. In the ONTAP WAFL file system, all I/O takes up 4KB of space, even if it does not actually require 4KB of data. Compaction combines multiple blocks that are not using their full 4KB of space together into one block. This one block can be more efficiently stored on-disk to save space. See Figure 7 for an illustration of compaction.

- **FlexClone volumes, files, and LUNs**: *FlexClone* technology references Snapshot metadata to create writable, point-in-time copies of a volume. Copies share data blocks with their parents, consuming no storage except what is required for metadata until changes are written to the copy. FlexClone files and FlexClone LUNs use identical technology, except that a backing Snapshot copy is not required.

Figure 7    Compaction in ONTAP 9



## NetApp Volume Encryption

Data security continues to be an important consideration for customers purchasing storage systems. NetApp has supported self-encrypting drives in storage clusters prior to ONTAP 9. However, in ONTAP 9, the encryption capabilities of ONTAP are extended by adding an Onboard Key Manager (OKM). The OKM generates and stores the keys for each of the drives in ONTAP, allowing ONTAP to provide all functionality

required for encryption out of the box. Through this functionality, sensitive data stored on disks is secure and can only be accessed by ONTAP.
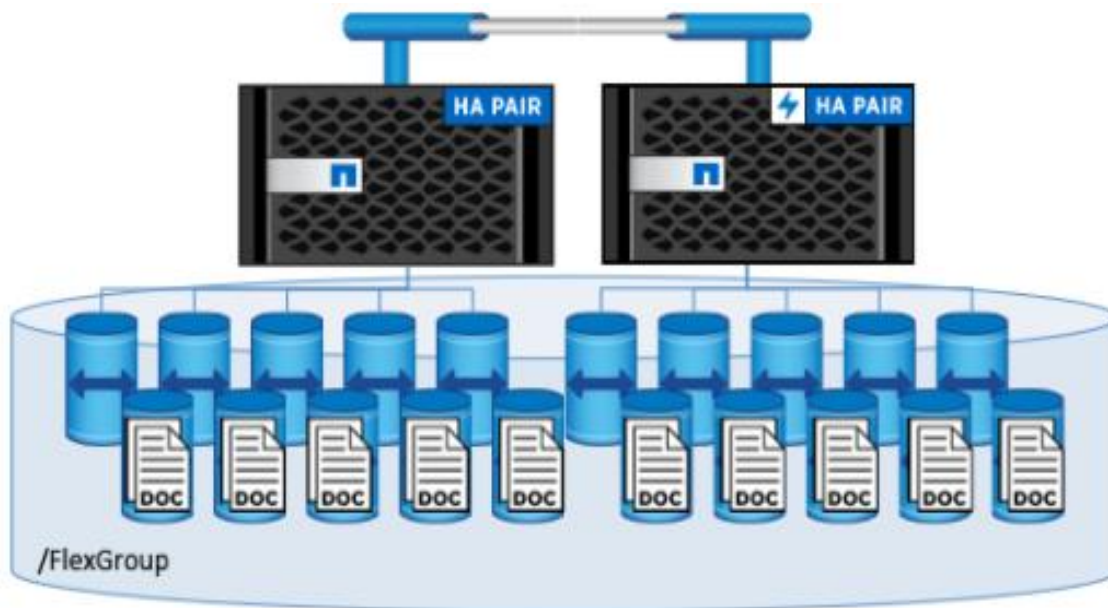
NetApp extended the encryption capabilities further with NetApp Volume Encryption (NVE) in ONTAP 9.1. NVE is a software-based mechanism for encrypting data. It allows a user to encrypt data at the per volume level instead of requiring encryption of all data in the cluster, thereby providing more flexibility and granularity to the ONTAP administrators. Once enabled, this encryption extends to Snapshot copies and FlexClone® volumes created in the cluster. To preserve all storage efficiency benefits, ONTAP executes NVE after all storage efficiency features. This ensures that customers have secure encrypted data while enjoying the benefits of reduced capacity utilization in the cluster.

For more information about encryption in ONTAP 9, see the [NetApp Encryption Power Guide](#).

## NetApp FlexGroup

NetApp introduced the FlexGroup - a scale-out NAS container that provides high performance, automatic load distribution, and scalability – in ONTAP 9. A FlexGroup volume contains several FlexVol volumes that automatically and transparently share traffic in the cluster.

Figure 8    NetApp FlexGroup



Files in a FlexGroup volume are allocated to individual member volumes and are not striped across volumes or nodes.  When a client adds files and sub-directories to a FlexGroup volume, ONTAP automatically determines the best FlexVol member to use for storing each new file and subdirectory. The FlexGroup volume attempts to organize the files, both for fastest accesses and for better throughput performance.

Advantages of NetApp ONTAP FlexGroup

- Massive capacity: FlexGroup volumes can scale up to multiple petabytes and with high file counts (hundreds of billions of files), The only limiting factors being the physical limits of the hardware and total volume limits of ONTAP. For example, a 10-node cluster can have a 20PB FlexGroup volume that can handle 400 billion files

- Predictable low latency for High-metadata workload: A FlexGroup volume utilizes all the cluster resources i.e. multiple aggregates, nodes, CPU cores and other hardware assets, thereby enabling multiple volume affinities to a single storage container for metadata intensive workloads.

- Ease of management: A FlexGroup volume can provision storage across every node and aggregate in a cluster (without any junction path / capacity management overhead) through the FlexGroup tab in NetApp OnCommand® System Manager.

## Backup and Replication

Traditionally, ONTAP replication technologies served the need for disaster recovery (DR) and data archiving. With the advent of cloud services, ONTAP replication has been adapted to data transfer between endpoints in the NetApp Data Fabric. The foundation for all these uses is ONTAP Snapshot technology.

- Snapshot copies: A *Snapshot copy* is a read-only, point-in-time image of a volume. The image consumes minimal storage space and incurs negligible performance overhead because it records only changes to the active file system since the last Snapshot copy. A volume can contain up to 255 Snapshot copies.

- SnapMirror disaster recovery and data transfer: *SnapMirror* is disaster recovery technology, designed for failover from primary storage to secondary storage at a geographically remote site. As its name implies, SnapMirror creates a replica, or *mirror,* of your working data in secondary storage from which you can continue to serve data in the event of a catastrophe at the primary site.

- SnapVault archiving: *SnapVault* is archiving technology, designed for disk-to-disk Snapshot copy replication for standards compliance and other governance-related purposes. In contrast to a SnapMirror relationship, in which the destination volume contains only the Snapshot copies currently in the source volume, a SnapVault destination volume typically retains point-in-time Snapshot copies created over a much longer period.

- MetroCluster continuous availability: MetroCluster configurations protect data by implementing two physically separate, mirrored clusters. Each cluster synchronously replicates the data and SVM configuration of the other. In the event of a disaster at one site, an administrator can activate the mirrored SVM and begin serving data from the surviving site.

# OnCommand Workflow Automation Pack for ACI

NetApp's OnCommand Workflow Automation (WFA) tool has been extended, to include interacting with Cisco ACI APICs. WFA provides a management framework and command and workflow libraries (organized into modular "packs", that contain related functionality) to automate NetApp storage management tasks, such as provisioning, migration, decommissioning, data protection configurations, and cloning storage. The WFA Pack for ACI extends WFA to include APICs, bridging the automation gap between storage and network. The WFA Pack for ACI is discussed in detail in NetApp Technical Report TR-4588, which is available here: http://www.netapp.com/us/media/tr-4588.pdf. This Technical Report:

- Discusses how to obtain the pack, and import it into a WFA instance

- Shows how to connect the WFA instance to an ACI APIC

- Examines the new ACI-related commands, explaining how they work, and offers tips on how to use them:

  – Create or Remove Storage Contracts

  – Create or Remove EPG

  – Provide or Consume Storage Contract

  – Create VPC Bundle

  – Create Port Specific VPC Bundle

  – Add or Delete VLAN Bundle

- Examines the workflows that are included in the pack, explaining how they work:

  – Create or Remove Storage Contracts

  – Add or Remove VLAN tagged ifgrp to tn/app/epg

  – Provide or Consume Storage Contract

- Shows how to build two custom workflows using the WFA Designer, that configure both ONTAP and ACI components:

  – Add ifgrp/vPC

  – Add LUN with iSCSI Access

## VMware vSphere 6.5 Update1

VMware vSphere is a virtualization platform for holistically managing large collections of infrastructure (resources-CPUs, storage and networking) as a seamless, versatile, and dynamic operating environment. Unlike traditional operating systems that manage an individual machine, VMware vSphere aggregates the infrastructure of an entire data center to create a single powerhouse with resources that can be allocated quickly and dynamically to any application in need.

VMware vSphere 6.5 brings a number of improvements including, but not limited to:

- Added native features to the vCenter Server Appliance (e.g. Integrated Update Manager)

- vCenter Server Appliance High Availability

- vSphere Web Client and fully supported HTML-5 client

- VM Encryption and Encrypted vMotion

- Improvements to DR

For more information, please refer to the VMware vSphere documentation.

# Solution Design

## FlexPod Datacenter Design Overview

The FlexPod Datacenter solution is a flexible, scalable and resilient architecture for delivering business and application agility in enterprise and cloud deployments. The solution leverages an ACI fabric that takes a policy-based, application centric view of the network to provide end-to-end connectivity and services. This design guide assumes a pre-existing ACI fabric so the focus of this document, is on utilizing that fabric to deliver a shared infrastructure solution that can support the applications and services that a business needs.

The following aspects of the design are covered in this document:

- Addition of compute infrastructure – Cisco Unified Computing System
- Addition of storage infrastructure – NetApp All-Flash FAS Unified Storage
- Enabling shared services within the ACI fabric - application and network layer services
- Access to existing shared services outside the fabric – application and network services
- Enabling external connectivity to networks outside the ACI fabric – shared access
- Enabling firewall services for applications
- Multi-tenancy to support the organizational structure of the business
- Deploying virtualized applications that utilize all of the above – multi-tier applications

This solution also meets the following general design requirements:

- Resilient design across all layers of the infrastructure with no single point of failure
- Scalable design with the flexibility to add compute capacity, storage, or network bandwidth as needed
- Modular design that can be replicated to expand and grow as the needs of the business grows
- Flexible design choices in terms of solutions components, storage access, connectivity options etc.
- Ability to simplify and automate, including integrating with external automation and orchestration tools
- Interoperable, integrated and verified design ready for deployment

The infrastructure, connectivity, services and the overall design were then built and validated using the components outlined in the Solution Validation section of this document.

# FlexPod Datacenter Design Options

This FlexPod Datacenter solution provides the following two options for accessing storage.

- End-to-end IP-based Storage through ACI fabric using iSCSI

- Direct Connect Storage access using Fiber Channel (FC) or Fibre Channel over Ethernet (FCoE)

In both cases, Cisco UCS servers will use SAN boot while the virtual machine and bare-metal workloads can use SAN protocols or NFS or both for storage.

## Design 1: End-to-End IP-Based Storage using iSCSI

The FlexPod Datacenter solution supports an end-to-end IP-based storage solution based on iSCSI. In this design, the Cisco UCS servers (through Cisco 6300 Fabric Interconnects) and the NetApp AFF array connect into the ACI leaf switches using 40GbE links. Cisco UCS servers use iSCSI to SAN boot from boot volumes residing on the NetApp AFF array. This connectivity is enabled through the ACI fabric. Workloads running on Cisco UCS will also use the fabric to access iSCSI LUNs and/or NFS volumes on the same NetApp array.

An iSCSI based design as shown in Figure 9 was validated for this CVD.

Figure 9    End-to-End IP-Based Storage using iSCSI



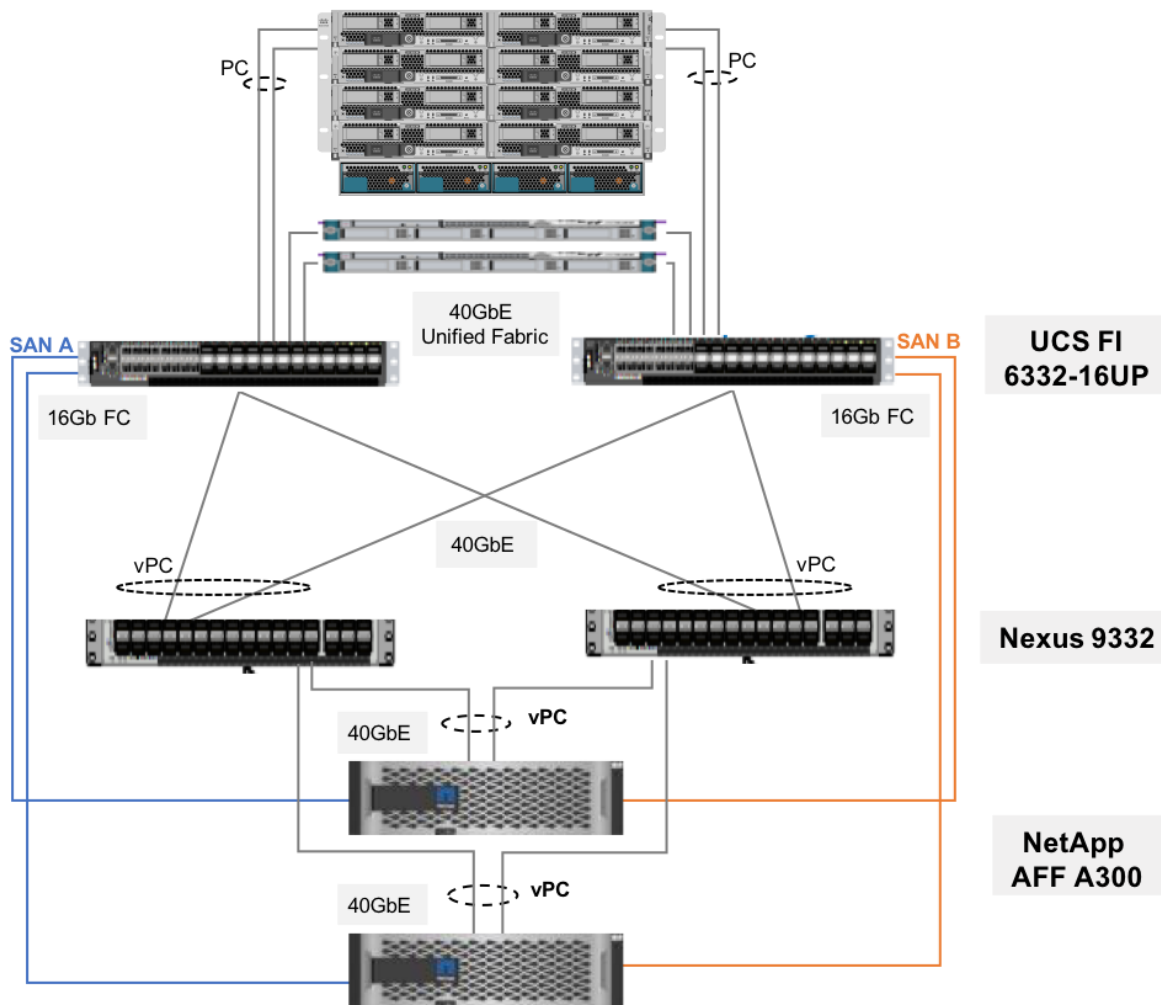## Design 2: Direct Connect Storage – SAN Access using FC or FCoE

The FlexPod Datacenter solution also supports FC or FCoE based SAN access by connecting the NetApp AFF arrays directly to Cisco UCS 6300 Fabric Interconnects using 16Gb/s FC links. This will require FC Zoning to be done in the Fabric Interconnects. The NetApp storage could also connect to the UCS FIs through a Cisco MDS-based FC network. In both these scenarios, the storage traffic for SAN boot and for accessing FC LUNs, will not traverse the ACI fabric when the UCS servers are in the same UCS domain (same FIs). The choice of FC vs. FCoE in direct connect mode is determined by the port and SFP type used on both the storage controller and Fabric Interconnects.

For this CVD, a FC-based design with NetApp storage directly connected to the Fabric Interconnects, as shown in Figure 10 was validated. To support NFS in this design, the NetApp arrays are also connected to the ACI fabric using 40GbE links. Designs using direct FCoE connectivity or FC connectivity through a MDS SAN network are supported but were not validated for this CVD.

Figure 10   Direct Connect Storage – SAN access using FC



## FlexPod Datacenter Design

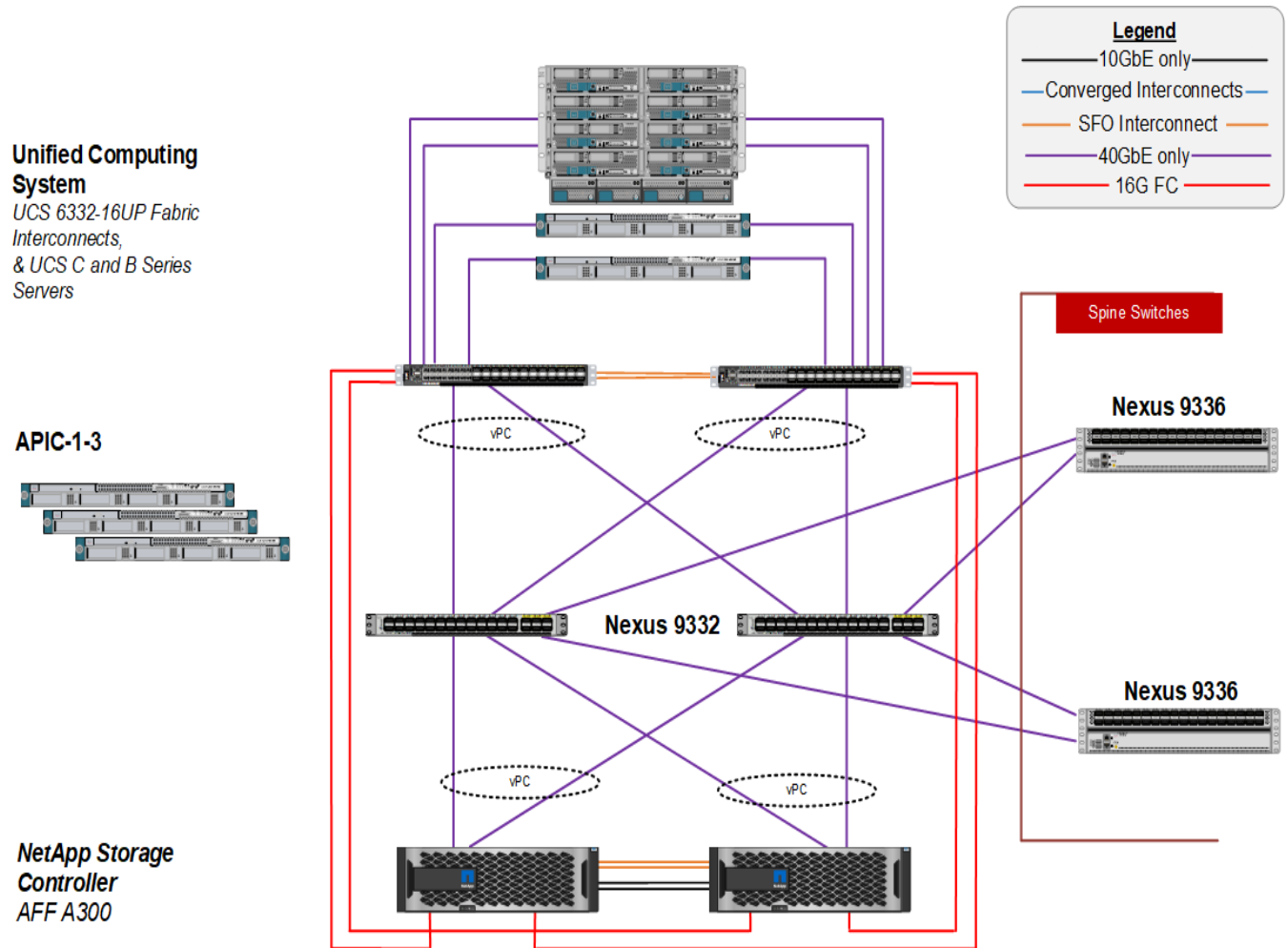In this section, we take a closer look at the FlexPod Datacenter design with ACI and IP-based storage. The ACI fabric is providing IP-based iSCSI access to NetApp storage. The NetApp array is providing both NFS volumes and iSCSI LUNs, including LUNs for SAN boot of Cisco UCS servers.

Figure 11 illustrates the end-to-end topology and the interconnections between the different components in the solution.

## Topology

Figure 11   FlexPod Datacenter with Cisco ACI and NetApp AFF Storage



## Connectivity Design – Compute Layer

The Cisco UCS platform provides the compute resources in the FlexPod Datacenter with Cisco ACI solution. The design supports both Cisco UCS B-series blade servers and Cisco UCS C-series rack-mount servers, connected and managed through a pair of Cisco UCS Fabric Interconnects running Cisco UCS manager.

### Blade Server Connectivity to Unified Fabric

Each Cisco UCS server is equipped with a Virtual Interface Cards (VIC) that aggregate all LAN and SAN traffic to and from the server across a single interface. Cisco VICs eliminate the need for separate physical interface cards on each server for LAN, SAN and management connectivity. Cisco VICs can be virtualized to create up to 256 virtual interfaces that can be dynamically configured as virtual network interface cards (vNICs) or virtual host bus adapters (vHBAs). These virtual interfaces will be presented and appear as standards-compliant PCIe endpoints to the OS. Blade and rack servers support different models of Cisco VIC. Cisco VICs are available in different form-factors and uplink speeds of up to 80Gbps.

## Blade Server Chassis Connectivity to Unified Fabric

The blade servers are housed in a Cisco UCS 5108 Blade Server Chassis that can support up to 8 half-width or 4 full-width blades. A blade server chassis have can have up to two fabric extenders (FEX) or I/O Modules (IOM) that connect the chassis to the Fabric Interconnects. Fabric Extenders operate in an active-active mode for data forwarding and active-passive mode for management functions

Fabric Extenders serve as a consolidation point for all blade server I/O traffic and provide 10/40GbE (FCoE) uplink connectivity from the chassis to the unified fabric, provided by a pair of Fabric Interconnects. Fabric Extenders extend the unified fabric to the chassis. Fabric Extenders also extend virtualization-aware networking through the unified fabric, from the Virtual Interface Cards (VIC) on the server to the Fabric Interconnects. FEX is managed as an extension of the fabric interconnects, simplifying diagnostics, cabling and operations with a single point of management and policy enforcement. This approach reduces the overall infrastructure complexity and enables the Cisco UCS to scale to multiple blade servers managed as a single, highly available management domain.

Table 3 provides a comparison of the 3 FEX models currently available for the blade server chassis. All FEX models are supported in this FlexPod design.

Table 3    Cisco UCS 2200 and 2300 Series Fabric Extenders

| Fabric Extender Model | External Links to FI (Blade Server Chassis Uplinks) | Internal Links to Servers *(via mid-plane)* | Total Uplink I/O to the chassis *(With redundancy)* |
|---|---|---|---|
| Cisco UCS 2204XP | 4 x 10GbE/FCoE (SFP+) | 16 x 10GbE | Up to 80Gbps |
| Cisco UCS 2208XP | 8 x 10GbE/FCoE (SFP+) | 32 x 10GbE | Up to 160Gbps |
| Cisco UCS 2304XP | 4 x 40GbE/FCoE (QSFP+) | 8 x 40GbE | Up to 320Gbps |

## Fabric Extender Connectivity to Unified Fabric

Each fabric extender connects to one Fabric Interconnect using multiple Ethernet (10GbE/40GbE) links – the number of links determines the uplink I/O bandwidth through that FEX. The number of links can be 1, 2, 4 or 8 depending on the model of FEX used. These links can be deployed as independent links (discrete Mode) or grouped together using link aggregation (port channel mode).

Figure 12   Fabric Extender to Fabric Interconnect Connectivity Options



In discrete mode, each server is pinned to a FEX link going to a port on the fabric interconnect and if the link **goes down, the server's connection also goes down through the FEX link.** In port channel mode, the flows from the server will be redistributed across the remaining port channel members. This is less disruptive overall and therefore port channel mode is recommended for this FlexPod design.

## Rack-Mount Server Connectivity to Unified Fabric

In FlexPod designs, the supported Cisco UCS C-Series servers can be either directly connected to the FIs using 10/40GbE links or through supported top-of-rack Cisco Nexus Fabric Extenders that connects to the FIs. FlexPod designs do require that these servers be managed by Cisco UCS Manager in order to ensure consistent policy-based provisioning, stateless computing and uniform management of the server resources, independent of the form-factor.

## Cisco UCS Fabric Interconnect Connectivity to Data Center Network

Fabric Interconnects are an integral part of the Cisco Unified Computing System, providing a unified fabric for integrated LAN, SAN and management connectivity for all servers that connect to the FIs (aka UCS domain). Fabric Interconnects provide a lossless and deterministic Fibre Channel over Ethernet (FCoE) fabric. Cisco UCS manager that manages the UCS domain, is embedded in the Fabric Interconnects.

This FlexPod design supports the following two models of Fabric Interconnects.

- Cisco UCS 6200 series fabric interconnects provide a 10GbE (FCoE) unified fabric with 10GbE uplinks for northbound connectivity to the datacenter network.

- Cisco UCS 6300 series fabric interconnects provide a 40GbE (FCoE) unified fabric with 40GbE uplinks for northbound connectivity to the datacenter network.

Fabric Interconnects support 802.3ad standards for aggregating links into a port-channel (PC) using Link Aggregation Protocol (LACP). Multiple links on each FI are bundled together in a port-channel and connected to upstream Nexus switches in the data center network. The port-channel provides link-level redundancy and higher aggregate bandwidth for LAN, SAN and Management traffic to/from the UCS domain.
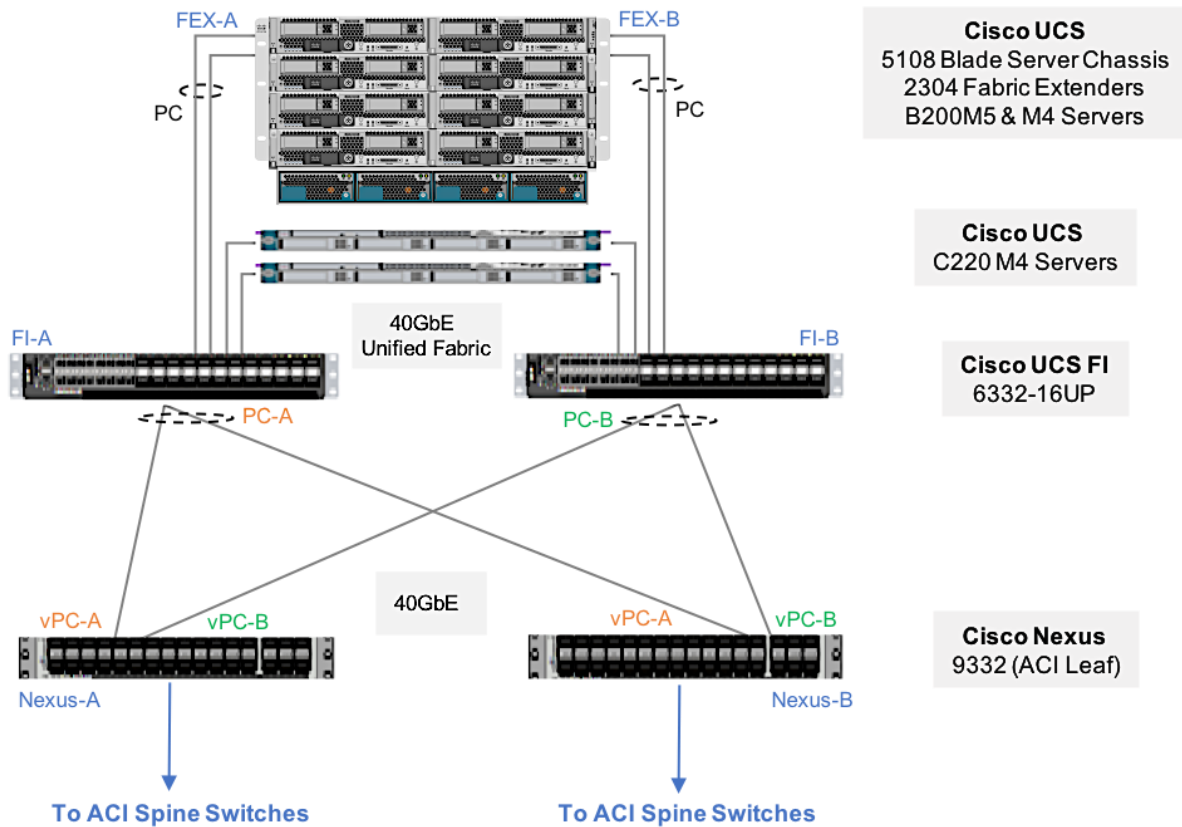
The Fabric Interconnect uplinks connect into a pair of upstream Nexus 9000 series switches (ACI Leaf switches) that are configured for virtual Port Channels (vPC). vPC allows links that are physically connected to two different Cisco Nexus 9000 Series devices to be bundled such that it appears as a "single logical" port channel to a third device (FI). vPC also provides higher aggregate bandwidth while providing both node and link-level fault tolerance.

In this design, each Fabric Interconnect connects into a pair of upstream Nexus 9000 ACI leaf switches. The links on each FI are bundled into a port-channel while links on Nexus leaf switches that connect to this FI are bundled into a vPC. This design provides link and node-level redundancy, higher aggregate bandwidth and the flexibility to increase the bandwidth as the uplink bandwidth needs grow.

## Validation – Compute Layer Connectivity

To validate the compute layer design, a Cisco UCS 5108 server chassis with Cisco UCS B200 M5 and B200 M4 half-width blade servers and two Cisco UCS C220 M4 rack mount servers are connected through a pair of Cisco UCS 6332-16UP Fabric Interconnects as shown in Figure 13.

Figure 13  Validated – Compute Layer Connectivity
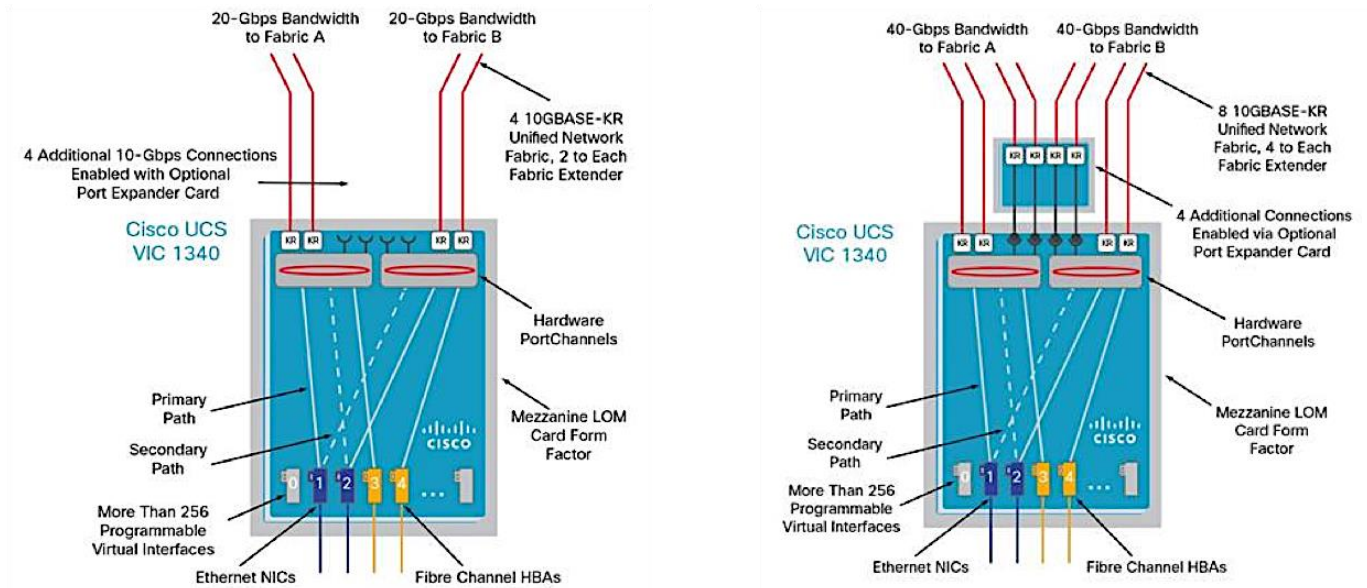


The blade server chassis is deployed using 2 x Cisco UCS 2304 FEX (IOMs), with each FEX connecting to one Fabric Interconnect, forming two distinct paths (Fabric-A, Fabric-B) through the unified fabric as follows:

- Fabric-A: 2 x 40GbE links from FEX-A to FI-A, links bundled into a port-channel

- Fabric-B: 2 x 40GbE links from FEX-B to FI-B, links bundled into a port-channel

This provides the blade server chassis with an aggregate uplink bandwidth of 160Gbps. Additional ports on each FEX can be used to further increase the bandwidth. For the 2304 FEX model, all 4 ports can be used for a total of 320Gbps of uplink bandwidth to a single blade server chassis.

The blade servers in the blade server chassis are each deployed with a single VIC 1340 adapter. The VIC 1340 adapter provides 40Gbps of uplink connectivity, 20Gbps through each Fabric (Fabric-A, Fabric-B) path. The uplink bandwidth can be increased to 80Gbps by using an optional port-expander, with 40Gbps through each Fabric path.

Figure 14    Cisco UCS Blade Server – VIC 1340 Uplink Connectivity



The rackmount servers are deployed with VIC 1385 adapters and directly connected to Fabric Interconnects, with one VIC (40GbE) port going to each FI, providing the rack servers with an aggregate uplink bandwidth of 80Gbps with high availability.

Figure 15    Cisco UCS Rack Server – VIC 1385 Uplink Connectivity

To connect to the upstream data center network, each FI is also connected to a pair of Nexus 9300 series leaf switches as follows:

- 2 x 40GbE links from FI-A to Nexus leaf switches (Nexus-A, Nexus-B), one to each Nexus switch

- 2 x 40GbE links from FI-B to Nexus leaf switches (Nexus-A, Nexus-B), one to each Nexus switch

The FI side ports are configured to be a port-channel, with vPC configuration on the Nexus leaf switches. This provides the UCS domain with redundant paths and 160Gbps of aggregate uplink bandwidth to/from the ACI fabric. The uplink bandwidth can be increased as needed by adding additional connections to the port-channel.

## Connectivity Design – Storage Layer

The FlexPod Datacenter with Cisco ACI solution is an end-to-end IP-based storage solution with iSCSI-based SAN access. The design can also support FC-based SAN access – see section on Design Options for more details. This design uses NetApp AFF A300 to provide the storage resources. NetApp Storage connects into the ACI fabric using dual 40GbE uplinks, configured for port-channeling to provide higher aggregate bandwidth and availability. Nexus Leaf switches that connect to the NetApp storage is configured for vPC to provide node-availability, in addition to link-availability and higher aggregate bandwidth.

The FlexPod Datacenter converged infrastructure supports a variety of NetApp controllers, including the AFF A-Series, AFF8000, FAS9000, FAS8000, FAS2600 and FAS2500 platforms. For a full list of supported controllers, see Interoperability Matrix for NetApp in Solution References.

FlexPod Datacenter architectures typically deploy FAS or AFF controllers running ONTAP data management software. Scalable and full-featured, ONTAP Data Management Software is ideal for converged infrastructure deployments. It enables the consolidation of all workloads onto a single system. All NetApp storage controllers are shipped with ONTAP installed and ready to begin using.

## Validation – Storage Layer Connectivity

To validate the storage layer design for IP-based storage access to application and boot volumes using iSCSI and NFS, two NetApp A300 arrays are deployed as a high availability (HA) pair and connected to a pair of Nexus leaf switches as shown in Figure 16. The NetApp A300 is deployed with 1.6TB SSDs and running clustered Data ONTAP 9.3 in a switchless cluster configuration.

Figure 16    Validated – Storage Layer Connectivity



This storage system can be easily scaled by adding disks and shelves to the existing HA pair, or by adding more HA pairs to the cluster. For two controllers to form an HA pair, several conditions must exist:

- Both controllers must have at least one (preferably two) SAS connections to each disk shelf, so that both controllers communicate with all disks.

- The NVRAM on both controllers must be in sync. (To improve write performance, write operations are cached to the battery-backed NVRAM in each controller, before being written to disk. In an HA pair, each controller caches its own write operations – and also that of its partner, so that in a takeover it knows where its partner left off. The two controllers copy their write operations to their partner over the high-speed HA interconnect.)

This ability to independently scale capacity and performance gives the storage admin a lot of options for efficiently handling growth and operations changes:

- If you are running out of disk space, but have plenty of controller performance headroom, you can add a disk shelf non-disruptively.

- If you are running out of controller performance, you can add nodes (controllers) to the cluster, and non-disruptively shift load and disk shelves over to the new controllers. That is a great way to migrate to a new model of controller, because only the partners in an HA pair have to be the same controller model: Different HA pairs in the same cluster can be different models.

For SAN environments, NetApp clustered Data ONTAP allows up to 6 HA pairs or 12 nodes. For NAS environments, it allows 12 HA pairs or 24 nodes to form a logical entity.

To connect to the upstream data center network, each NetApp array in the HA pair is connected to a pair of Nexus 9300 series leaf switches as follows:

- 2 x 40GbE links from each NetApp array to leaf switches, one link to each leaf (Nexus-A, Nexus-B),

- Port-channel configuration with 2 x 40GbE ports on each NetApp array

- vPC configuration on Nexus leaf switches, one vPC to each NetApp array. Each VPC has 2 links, one from each Nexus switch to NetApp.

The connectivity described above, provides each NetApp AFF A300 with redundant uplinks through separate leaf switches and 80Gbps (160Gbps for the HA pair) of bandwidth to the ACI fabric.

## Connectivity Design – Network Layer

The Cisco ACI fabric consists of discrete components that operate as routers and switches but are provisioned and monitored as a single entity. These components and the integrated management enable ACI to build programmable data centers that provide connectivity as well as advanced traffic optimization, security, and telemetry functions for both virtual and physical workloads.

Cisco ACI fabric is fundamentally different from other data center network designs. A basic understanding of ACI from a technology standpoint is therefore important for understanding an ACI-based design. To learn more about ACI, see ACI section of Solution Components and ACI resources listed in Solution References.

### Nexus 9000 Hardware

The ACI fabric is based on a spine-leaf architecture, built using Nexus 9000 series switches where each leaf switch connects to every spine switch, using high speed links and with no direct connectivity between leaf nodes or between spine nodes. Multiple models of Nexus 9000 series switches that support ACI spine and leaf functionality are available and supported in FlexPod.

In ACI, spine switches form the core of the ACI fabric and provide high-speed (40/100GbE) connectivity between leaf switches. A spine can be a:

- Modular Cisco Nexus 9500 series switch equipped with 40/100GbE capable line cards such as N9K-X9736PQ, N9K-X9736C-FX, N9K-X9732C-EX etc.

- Fixed form-factor Cisco Nexus 9300 series switch with 40/100GbE ports (e.g. N9K-C9336PQ, N9K-C9364C)

The edge of the ACI fabric are the leaf switches. Leaf switches are top-of-rack (ToR), fixed form factor Nexus switches such as N9K-C9332PQ, N9K-C93180LC-EX, N9K-C93180YC-EX/FX, Nexus 9300 PX series switches. These switches will typically have 40/100GbE uplink ports for high-speed connectivity to spine switches and access ports that support a range of speeds (1/10/25/40GbE) for connecting to servers, storage and other network devices.

Leaf switches provide access layer functions such as traffic classification, policy enforcement, traffic forwarding and serve as an attachment point to the ACI fabric. Nexus leaf switches also provide several advanced capabilities such as support for analytics in hardware, advanced traffic management, encryption, traffic redirection for L4-L7 services, etc.

### Edge Connectivity

Leaf switches at the edge of the ACI fabric, typically provide physical connectivity to the following types of devices and networks:

- Cisco APICs that manage the ACI Fabric

- Servers

- Storage

- Outside Networks **– These are networks within the customer's network** but outside the ACI fabric. These include management, campus, WAN and other data center and services networks.

- External Networks **–** These are networks that directly connect to the ACI fabric from outside the **customer's network**. This could be enabling cloud services and applications from the datacenter or for access to external cloud resources.

In this design, the following components are physically connected to the leaf switches:

- Cisco APICs that manage the ACI Fabric (3-node cluster)

- Cisco UCS Compute Domain (Pair of Cisco UCS Fabric Interconnects)

- NetApp Storage Cluster (Pair of NetApp AFF A300s)

- Management Network Infrastructure **in customer's** existing management network (Outside Network)

- Routers that serve as a gateway to campus, WAN and other parts of a **customer's existing network** (Outside Network)

Cisco ACI supports virtual Port-Channel (vPC) technology on leaf switches to increase throughput and resilience. Virtual Port channels play an important role on leaf switches by allowing the connecting devices to use 802.3ad LACP-based port-channeling to bundle links going to two separate leaf switches. Unlike traditional NX-OS vPC feature, the ACI vPC does not require a vPC peer-link to be connected nor configured between the peer leaf switches. Instead, the peer communication occurs through the spine switches, using the uplinks to the spines. The vPCs can therefore be created between any two leaf switches, in the ACI fabric through configuration, without having to do additional cabling between switches.

When creating a vPC between two leaf switches, the switches must be of the same hardware generation. Generation 2 models have -EX or -FX or -FX2 in the name while Generation 1 does not.

In this FlexPod design, vPCs are used for connecting the following access layer devices to the ACI fabric:

- Two vPCs, one to each Fabric Interconnect in the Cisco UCS Compute Domain

- Two vPCs, one to each NetApp storage array in the HA-pair

- Single vPC to management network i**nfrastructure in customer's existing management network**

The other access layer connections are individual connections - they are not part of a vPC bundle.

## ACI Fabric

The Cisco ACI fabric is a Layer 3, routed fabric with a VXLAN overlay network for enabling L2, L3 and multicast forwarding across the fabric. VXLAN overlays provide a high degree of scalability in the number of Layer 2 segments it can support as well as the ability to extend these Layer 2 segments across a Layer 3 network.  The ACI fabric provides connectivity to both physical and virtual workloads, and the compute, storage and network resources required to host these workloads in the data center.
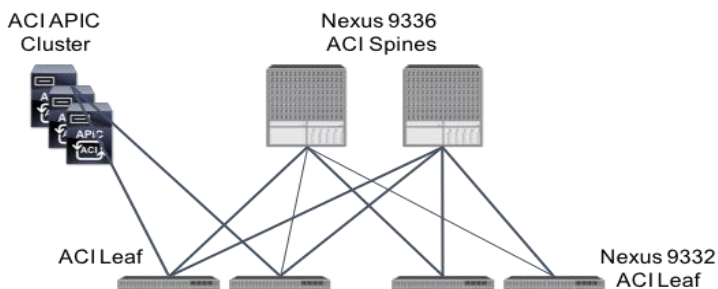
The ACI architecture is designed for multi-tenancy. Multi-tenancy allows the administrator to partition the fabric along organizational or functional lines into multiple tenants. The ACI fabric has three system-defined tenants (mgmt, infra, common) that gets created when a fabric first comes up. The administrator defines the user tenants as needed to meet the needs of the organization.  shared-services tenant for hosting infrastructure services such as Microsoft Active Directory (AD), Dynamic Name Services (DNS) etc.

This FlexPod design assumes that the customer already has an ACI fabric in place with spine switches and APICs deployed. The design guidance in this document therefore focusses on attaching different sub-systems (compute, storage) to the existing fabric to enable an end-to-end converged datacenter infrastructure based on Cisco UCS, NetApp Storage and Cisco ACI.

## ACI Fabric Connectivity Design

This design assumes an existing 40GbE ACI fabric, built as shown in Figure 17. It consists of a pair of Nexus 9336PQ spine switches, a 3-node APIC cluster and a pair of leaf switches that the Cisco APICs connect into using 10GbE – APICs only support 10GbE uplinks currently. The core provides 40GbE connectivity between leaf and spine switches. The fabric design can support other models of Nexus 9000 series switches, provided it has the required interface types, speeds and other capabilities. The design of the existing ACI fabric is outside the scope of this document.

Figure 17    FlexPod Datacenter – Existing ACI Fabric



In this FlexPod design, a pair of Nexus 9332PQ leaf switches are added to the existing ACI fabric to provide connectivity to downstream compute, storage and other network sub-systems as shown in Figure 18. The leaf switches use 40GbE links to connect to spine switches (Nexus 9336PQ) and for downstream connectivity to Cisco UCS servers and NetApp storage. The connectivity to outside networks and services are provided by the leaf switches in the existing ACI fabric for 10GbE access - the N9332 PQ switches only support 40GbE.

Figure 18    FlexPod Datacenter – ACI Fabric Sub-system Connectivity (High-Level)



The access layer connections on the ACI leaf switches to the different sub-systems are summarized below:

- A Cisco UCS Compute domain consisting of a pair of Cisco UCS 6300 Fabric Interconnects, connect into a pair of Nexus 9332PQ leaf switches using port-channels, one from each FI. Each FI connects to the leaf switch-pair using member links from one port-channel. On the leaf switch-pair, the links from each FI are bundled into a vPC. This design uses 2x40GbE links from each FI to leaf-switch pair to provide the UCS compute domain with an uplink bandwidth of 160Gbps. Additional links can be added to the bundle as needed, to increase the uplink bandwidth.

- A NetApp Storage cluster consisting of two NetApp AFF A300 arrays, connect into a pair of Nexus 9332PQ leaf switches using port-channels, one from each array. Each array connects to the leaf switch-pair using member links from one port-channel. On the leaf switch-pair, the links from each array are bundled into a vPC. This design uses 2x40GbE links from each array to leaf-switch pair to provide the NetApp storage domain with an uplink bandwidth of 160Gbps. Additional links can be added to the bundle as needed, to increase the uplink bandwidth. This design supports both NFS and iSCSI for storage access. This design can also support FC/FCoE SAN access by connecting through the UCS Fabric Interconnects. The NetApp storage arrays can be either directly attached using 16G FC or 10G FCoE or connected through a MDS-based SAN network using 16G FC to the UCS FIs. See Design Options section of this document for more details.

- To connect to an existing management network outside the ACI fabric, a single vPC is used to **connect to the customer's management network** infrastructure – in this case, a Nexus 5000 series switch. A 10GbE vPC from a pair of leaf switches connects to a port-channel on the Nexus 5k switch. This connectivity is enabled through a 10GbE-capable leaf switch pair in the existing ACI fabric since the Nexus 9332PQ leaf switches only supports 40GbE connections. From the ACI **fabric's perspective, this is a L2 bridged connection.**

- **To connect to customer's existing networks (e.g. Campus, WAN) outside the ACI fabric, 10GbE** links are used to directly connect two leaf switches to two Nexus 7000 series switches. Nexus 7k

switches provide reachability to other parts of **the customer's network. From the ACI fabric's** perspective, this is a L3 routed connection.

## ACI Fabric Design – Access Policies

Fabric Access Policies is an important aspect of the Cisco ACI architecture. Fabric Access Policies are defined by the Fabric Administrator and includes all the configuration and policies required to connect access layer devices to the ACI fabric. This must be in place before Tenant Administrators can deploy Application EPGs. These policies are designed to be reused as new leaf switches and access layer devices are connected to the fabric.

Fabric Access refers to access layers connections at the fabric edge to external devices such as:

- Physical Servers (Cisco UCS Rackmount servers, NetApp Storage Controllers)

- Layer 2 Bridged Devices(Switches, Cisco UCS FI)

- Layer 3 Gateway Devices(Routers)

- Hypervisors (ESXi) and Virtual Machine Managers (VMware vCenter)

Access Policies include any configuration that can applied to the above connections such as:

- Enabling PC, vPC on the links to the external devices

- Configuring Interface Policies (e.g. VLAN scope global, VLAN scope port-local, LACP, LLDP, CDP)

- Interface Profiles

Access Polices also include configuration and policies for leaf switches such as:

- VLAN pools

- Domain Profiles

- Switch Profiles

- Global policies such as AEP

The Fabric Access Policies used in the FlexPod design to connect to an outside Management network and UCS domains are shown in the figure below. Once the policies and profiles are in place, they can be re-used to add new leaf switches and connect new endpoints to the ACI fabric. Note that the Create an Interface, PC, and vPC wizard simplifies this task.

Figure 19   Fabric Access Policies to connect to UCS Domain and Management Network



## Fabric Access Design – Access Policies (VLAN Pools) and VLAN Design

Cisco ACI uses a VXLAN overlay network for communication across the fabric but use VLANs (or VXLAN) to communicate with devices outside the fabric. Unlike other data center architectures, VLANs in ACI are used to classify incoming traffic based on their VLAN tag but are not used to make forwarding decisions within the fabric. The traffic received on a VLAN are classified and mapped to an Endpoint group (EPG). EPG is a fundamental construct in ACI and forms the basis for policies and forwarding through the fabric. The EPG policies will determine the forwarding, and VXLAN tunnels will transport the traffic across the ACI fabric. For more details, see Cisco ACI Fundamentals listed in Solution References section.

Though VLANs are only used at the edge of the ACI fabric, they are still very tightly coupled to EPGs. To enable any forwarding within the fabric, a Tenant Administrator must first define an EPG and associate it with the physical infrastructure, which are typically access ports that connect to end devices that make up the EPG. The connectivity to an EPG endpoint on an access port is typically through a VLAN. An untagged port or VXLAN can also be used though they are not used in the FlexPod design. A port can provide connectivity to multiple EPG endpoints by defining multiple VLANs, one for each EPG. In the FlexPod design, the vPCs to compute, storage and outside networks all carry multiple VLANs, each mapping to different EPGs.

The Fabric Administrator defines and manages access configuration on a switch through fabric access policies. These policies include domains and VLAN pools, where the pools specify the allowed range of VLANs and the domains specify the scope of the associated vlan pool. In ACI, a domain can be physical (e.g. rackmount server, storage array), virtual (e.g. virtual machine manager or VMM) or external (e.g. L2 or L3 networks outside the ACI fabric). The domain associated with an EPG will determine the allowed range of VLANs on the EPG ports that connect to that domain. Multiple domains can be associated with a single EPG. Therefore, an EPG can have multiple VLANs in multiple domains to connect to all the endpoints in the EPG.

The endpoints could be spread across different leaf switches or on different ports of the same switch using same or different VLAN. Though not common, multiple endpoints on the same access port can use different VLAN IDs to map to the same EPG – this is the case for VMs deployed in the VMM domain and the hosts in the same UCS domain are part of the same EPG (e.g. Core-Services EPG).

When deploying an EPG, VLANs are allocated to provide connectivity to endpoints as outlined below. In the FlexPod design, EPG VLANs provide connectivity to endpoints on UCS compute, NetApp storage, networks outside the ACI fabric and virtual domains managed vCenter.

- Static Allocation is typically used on connections to a physical device in ACI. In the FlexPod design, static binding is used when connecting to any physical infrastructure. This includes connections to the FIs in the Cisco UCS domain, NetApp Storage Cluster and other networks outside the ACI fabric.

- Dynamic Allocation is typically used in ACI for virtual machines managed by a Virtual Machine Manager (VMM). In the FlexPod design, dynamic allocation is used when connecting to a virtual domain. APIC integrates with VMware vCenter managing the virtual domain to allocate VLANs as needed from a pre-defined VLAN pool. As new virtual machines come up in an EPG, vCenter will notify APIC and APIC will select a VLAN from the pool to allocate to the virtual endpoint.

The scope of an allocated VLAN on an interface can impact the overall VLAN design. If the VLAN scope is:

- Global, then the VLAN ID can only be associated with one EPG on a given leaf switch. The same VLAN ID can be mapped to a second VLAN ID but only on a different leaf switch.

If the EPGs on different leaf switches belong to the same bridge domain, they will be part of the same BPDU flooding domain and broadcast domain.

- Port-local, then the same VLAN ID can be mapped to multiple EPGs on the same leaf switch, if the EPGs are associated with different bridge domains. Also, the EPGs must be deployed on different ports. When the scope is port-local, the leaf switch will maintain translation entries (to/from fabric) for that VLAN on a per-port basis. Maintaining additional information does use up more hardware resources on that leaf switch. The FlexPod design uses port-local scope on all access ports since it provides more flexibility in VLAN allocation & usage.

An EPG is always associated with a bridge domain. See Cisco ACI fundamentals listed in Solution References section for more details on how EPGs related to Bridge Domains. Bridge domains in this design will be cover in a later section.

The tables below show the EPG VLANs used in the FlexPod design. These VLANs are enabled on the port-channels connecting leaf switches to access layer devices and provide compute, storage and management domains access to the ACI fabric.

**Table 4   Access Layer Connectivity on Leaf Switches – EPG VLANs to Management Switch**

| vPC to Management Switch | VLAN Name & ID | Purpose |
|---|---|---|
| Domain Name: `FP-Mgmt-Sw`<br>Domain Type: `External Bridged(L2)Domain`<br>VLAN Scope: `Port-Local`<br>Allocation Type: `Static`<br>VLAN Pool Name: `FP-Mgmt-Sw` | **IB-MGMT**<br>(118) | Provides connectivity to an existing In-Band management network, outside the ACI Fabric. This connection provides access to existing management and infrastructure services (same as Core-Services EPG in the ACI Fabric) |

**Table 5   Access Layer Connectivity on Leaf Switches – EPG VLANs to NetApp AFF Storage Cluster**

| vPC to NetApp AFF Cluster | VLAN Name & ID | Purpose |
|---|---|---|
| Domain Name: `NetApp-AFF`<br>Domain Type: `Bare Metal`(Physical)<br>VLAN Scope: `Port-Local`<br>Allocation Type: `Static`<br>VLAN Pool Name: `NetApp-AFF_vlans` | **SVM-MGMT**(219) | Provides access to Tenant SVMs on NetApp AFF Cluster |
| | **NFS**(3050) | Provides access to NFS volumes on NetApp AFF Cluster |
| | **iSCSI-A**(3010) | Provides access to boot, application data and datastore LUNs on NetApp AFF Cluster via iSCSI path-A |
| | **iSCSI-B**(3020) | Provides access to boot, application data and datastore LUNs on NetApp AFF Cluster via iSCSI path-B |

**Table 6   Access Layer Connectivity on Leaf Switches – EPG VLANs to Cisco UCS Compute Domain**

| vPC to Cisco UCS Fabric Interconnects | VLAN Name & ID | Purpose |
|---|---|---|
| Domain Name: `UCS`<br>Domain Type: `External Bridged(L2)Domain`<br>VLAN Scope: `Port-Local`<br>Allocation Type: `Static`<br>VLAN Pool Name: `UCS_vlans` | **Native**(2) | Security best practice to change Native Vlan from default of '1'; Used by control plane protocols |
| | **Core-Services**(318) | Provides access to services hosted on Cisco UCS - UCS hosts connect to ACI Fabric leaf switches |
| | **AV-IB-MGMT**(419) | To access In-Band Management Network for UCS hosts connected to ACI Fabric leaf switches and management VM hosted on the same UCS hosts |
| | **AV-vMotion**(3000) | To access vMotion network for VMs hosted on UCS hosts connected to ACI Fabric leaf switches |
| | **AV-Infra-NFS**(3150) | To access NFS volumes on NetApp AFF Cluster from UCS hosts and VMs through the ACI fabric |
| | **AV-Infra-iSCSI-A**(3110) | To access boot, application and datastore LUNs on NetApp AFF cluster via iSCSI Path-A from Cisco UCS hosts through the ACI fabric |
| | **AV-Infra-iSCSI-B**(3120) | To access boot, application and datastore LUNs on NetApp AFF cluster via iSCSI Path-B from Cisco UCS hosts through the ACI fabric |

Table 7   Access Layer Connectivity on Leaf Switches – EPG VLANs to Layer 3 Outside Network

| Redundant Connections from Leaf Switch pair to L3 Network | VLAN ID | Purpose |
|---|---|---|
| Domain Name: `Shared-L3-Out`<br>Domain Type: `External Routed(L3)Domain`<br>Allocation Type: `Static`<br>VLAN Pool Name: `VP-Shared-L3-Out` | 201 | Primary connection between 1st leaf switch and 1st L3 Gateway outside the fabric |
| | 202 | Primary connection between 2nd leaf switch and 1st L3 Gateway outside the fabric |
| | 203 | Secondary connection between 1st leaf switch and 2nd L3 Gateway outside the fabric |
| | 204 | Secondary connection between 2nd leaf switch and 1st L3 Gateway outside the fabric |

The access layer connectivity to a VMM domain is through the vPCs to the UCS domain hosting the virtual environment. APIC integration with VMware vCenter enables EPGs to deployed in the VMM domain. To communicate with a virtual endpoint in an EPG, VLANs are dynamically assigned based on VM events from VMware vCenter. As VMs come online, vCenter notifies the APIC and a VLAN is allocated from the pool. The EPG VLANs used in the FlexPod design for connectivity to virtual endpoints are listed in Table 8 .

Table 8   Access Layer Connectivity on Leaf Switches – EPG VLANs to VMM Domain

| APIC to VMM Domain Integration | VLAN ID | Purpose |
|---|---|---|
| Domain Name: `fpv-vc-vDS`<br>Domain Type: `VMM Domain`<br>Allocation Type: `Dynamic`<br>VLAN Pool Name: `VP-fpv-vc-vDS` | `[1100-1199]` | VLAN for Application EPGs hosted on Cisco UCS servers. The physical connectivity to the EPG virtual endpoints are through the vPCs to Fabric Interconnects in the Cisco UCS domain. APIC to VMM integration is used to dynamically assign VLANs as new virtual endpoints come online. |

## VLAN Design Guidelines and Best Practices

VLAN scalability (4096 VLANs) can be a limitation in traditional data center networks but since VLAN are only used at the edge to communicate with devices outside the fabric. The ACI guidelines and best practices that impact the VLAN design are summarized below. Some of the guidelines are dependent on other ACI design constructs that have not been covered yet but will be in the following subsections.
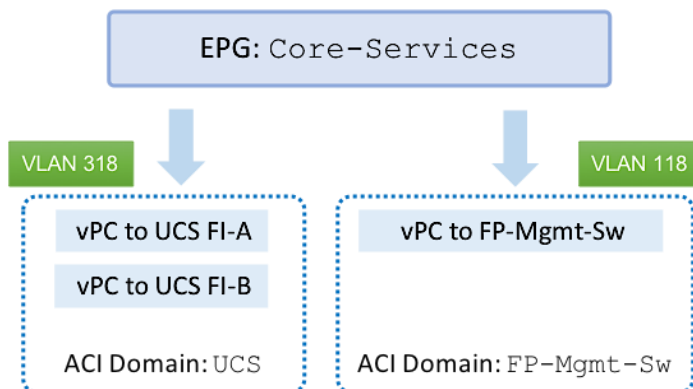
- Overlapping VLAN IDs can used on a single leaf switch or across leaf switches if the connecting devices are part of the same EPG.

- Overlapping VLAN IDs can be used on different EPGs of the same leaf switch if the different EPG devices are on separate ports, belong to different bridge domains and the vlan scope is port-local.

- Overlapping VLAN IDs can be used across different leaf switches and EPGs if the vlan scope is global. The EPGs can belong to the same bridge domain in this case but will be in the same broadcast domain.

- A VLAN pool can be shared by multiple ACI domains but a domain can only be mapped to one VLAN pool. VLAN pool and domains determine the range of VLANs allowed on the interfaces connecting to that domain and are part of the access policies established by the Fabric Administrator.

- When deploying an EPG with static binding to a physical interface, the VLAN ID specified must be from the allowed range of VLANs for that interface. The domain association for the EPG maps to a VLAN pool and this domain must also be associated with physical interface. The domain and the associated VLAN pool is mapped to the physical interface through the Access Entity Profile (AEP). AEP defines the scope of the VLAN (VLAN pool) on the physical infrastructure (port, PC or vPC). This ensures that the EPG VLAN deployed on the physical infrastructure is within the range of VLANs allowed on that infrastructure.

The VLAN guidelines used in the FlexPod design are as follows:
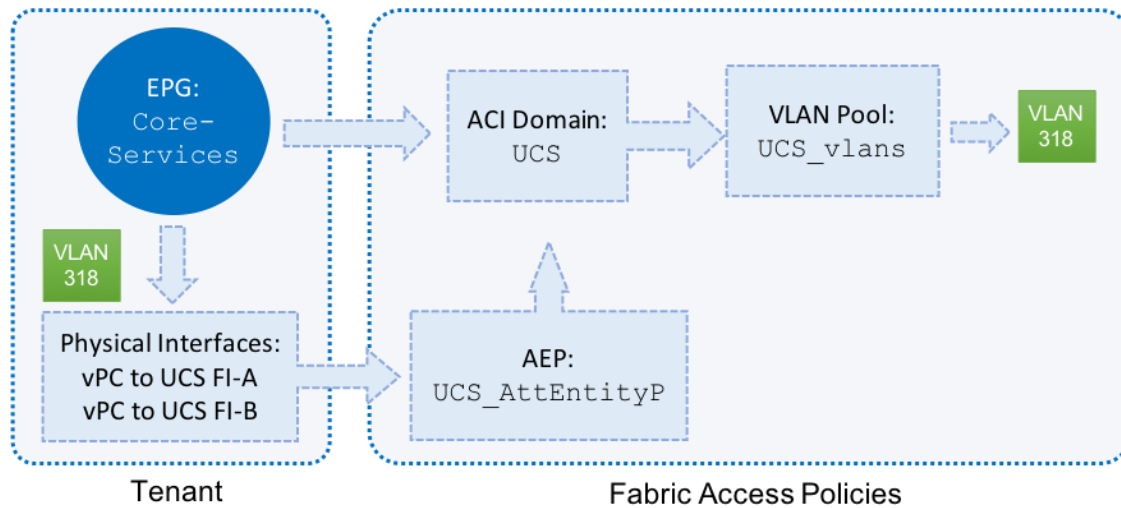
- EPGs will use the same VLAN IDs to connect to endpoints in the same ACI domain as shown in Figure 20. For example, the Core-Services EPG in the FlexPod design will use the same VLAN ID to connect to endpoints in the UCS domain though the connectivity is through different vPCs interfaces going to different FIs. Since both vPC interfaces connect to the same ACI domain (UCS) and same VLAN pool (UCS_vlans), the same VLAN ID is used. In this case, the EPG deployed on the two vPC interfaces connect to the same endpoint but the same is true for other endpoints in the same EPG and ACI domain. As new UCS domains are attached to the ACI fabric, this allows for endpoints in the new UCS domains to be added to the Core-Services EPG using the same VLAN ID, ACI Domain and VLAN pool. However, the Core-Services EPG deployed on the vPC to the Management Switch, will use a different VLAN ID as it is in a different ACI domain (FP-Mgmt-Sw) as shown in Figure 20.

Figure 20    FlexPod Design – EPG VLAN Allocation



- Static binding is used to deploy EPGs that connect to UCS, NetApp and outside networks. The EPG deployed by the Tenant Administrator and the access policies established by the Fabric Administrator must be aligned on the interface that the EPG is deployed on. For an EPG, the VLAN ID used to connect to an endpoint on a given interface must be in the VLAN pool specified by the Fabric Access Policy for that interface – see Figure 21.

Figure 21   EPG – Fabric Access Policy Relationship (Static Binding)



## Fabric Access Design – Access Policies (ACI Domains)

Domains in ACI are used to define how different entities (e.g. servers, network devices, storage) connect into the fabric and specify the scope of a defined VLAN pool. ACI defines four domain types based on the type of devices that connect to the leaf switch. They are:

- Physical domains are generally used for connections to bare metal servers or servers where hypervisor integration is not an option.

- External bridged domains are used to connect to Layer 2 devices outside the ACI fabric. These are referred to as Layer 2 Outside or L2OUT connections in ACI.

- External routed domains in ACI are used to connect to Layer 3 devices outside the ACI fabric. These are referred to as Layer 2 Outside or L2OUT connections in ACI.

The table below shows the ACI domains for the different access layer connections in the FlexPod design.

Table 9   Access Layer Connectivity on Leaf Switches – ACI Domains

| Access Connection | Domain Name | Domain Type |
|---|---|---|
| vPC to Management Switch | FP-Mgmt-Sw | External Bridged Devices(L2OUT) |
| vPC to NetApp AFF Cluster (AFF-1, AFF-2) | NetApp-AFF | Bare Metal(**Physical**) |
| vPC to Cisco UCS Fabric Interconnects (FI-A, FI-B) | UCS | External Bridged Devices(L2OUT) |
| Redundant Links to a pair of L3 Gateways | Shared-L3-Out | External Routed Devices(L3OUT) |
| VMM Domain | fpv-vc-vDS | VMM Domain |

Domains serve as a link between the fabric access policies defined by the Fabric Administrator and the ACI policy model and endpoint groups defined by the Tenant Administrator. Tenant Administrator will deploy EPGs and associate them with domains created by the Fabric Administrator. For static EPGs, the VLAN encapsulation used to classify the ingress traffic to a static EPG must be part of the vlan pool associated with that domain. For EPGs with dynamic members, the VLANs are dynamically assigned based on VM events from vCenter.

### Domain Design Guidelines and Best Practices

The design guidelines and best practices for ACI domains and their use in the FlexPod design are as follows.

- In ACI, a VLAN pool can be associated with multiple domains but a domain can only be associated with one VLAN pool. In the FlexPod design, the compute, storage and networks outside the ACI fabric are in different domains; a VLAN pool is created for each domain as shown in Table 10

Table 10   FlexPod Design – Domain to VLAN Pool Mapping

| Access Connection | Domain Name | Domain Type | VLAN Pool |
|---|---|---|---|
| vPC to Management Switch (Outside ACI Fabric) | `FP-Mgmt-Sw` | `External Bridged(L2) Domain` | `FP-Mgmt-Sw` `(118)` |
| vPC to NetApp AFF Cluster (AFF-1, AFF-2) | `NetApp-AFF` | `Bare Metal`(Physical) | `NetApp-AFF_vlans` `(219,3010,3020,3050)` |
| vPC to Cisco UCS Fabric Interconnects (FI-A, FI-B) | `UCS` | `External Bridged(L2) Domain` | `UCS_vlans` `(2,318,419,3000,3110,3120,3150)` |
| Redundant Links to a pair of L3 Gateways | `Shared-L3-Out` | `External Routed(L3) Domain` | `Shared-L3-Out_vlans` `(201,202,203,204)` |
| To VMM Domain | `fpv-vc-vDS` | `VMM Domain` | `VP-fpv-vc-vDS` `[1100-1199]` |

- When an EPG is deployed on a vPC interface, the domain (and the VLAN pool) associated with the EPG must be the same as the domain associated with the vPC interfaces on the leaf switch ports. When deploying an EPG on an interface, a VLAN in the domain that the interface connects to must be specified. The EPG must also be associated with a domain and VLAN pool and the VLAN specified on the interface must be part of that VLAN pool for the EPG to be operational.

- In addition to VLAN pools, ACI domains are also associated to Attachable Entity Profiles (AEP) – AEP will be covered later in this section. Multiple domains can be associated with a AEP. Domains link VLAN pools to AEP.  Without defining a VLAN pool in an AEP, a VLAN is not enabled on the leaf port even if an EPG is provisioned.
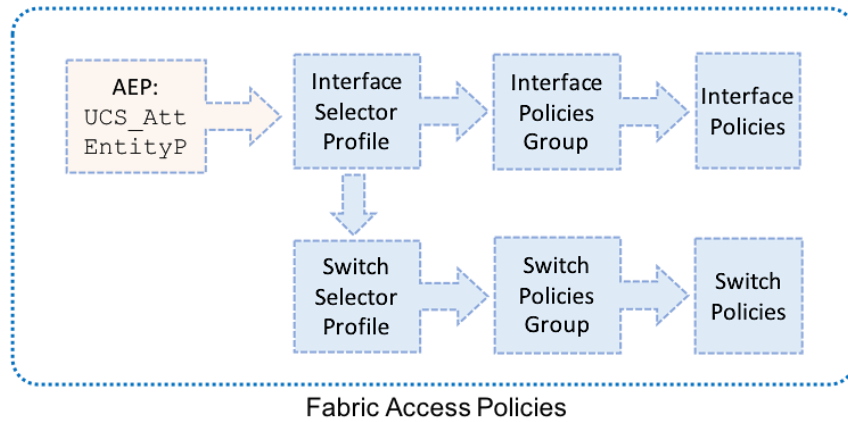
## Fabric Access Design – Access Policies (Attachable Entity Profile)

Attachable Entity Profile (AEP) is an ACI construct for grouping external devices with common interface policies. AEP is also known as a Attachable Access Entity Profile (AAEP).

ACI provides multiple attachment points for connecting access layer devices to the ACI fabric. Interface Selector Profiles represents the configuration of those attachment points. Interface Selector Profiles are the
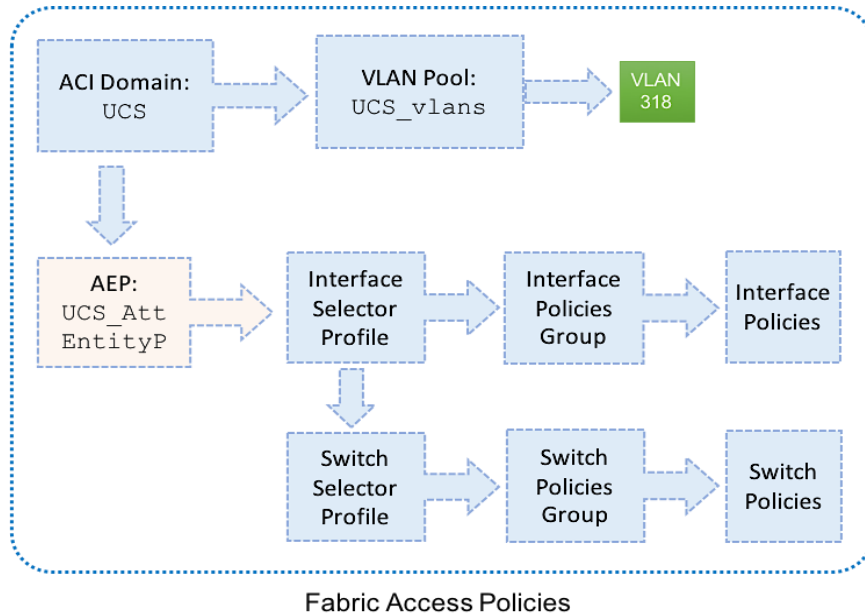
consolidation of a group of interface policies (e.g. LACP, LLDP, CDP) and the interfaces they apply to. The Interface Profiles and AEPs can be reused across multiple switches if the policies and ports are the same.

Figure 22   AEP Mapping to Interface and Switch Policies for UCS Domain



Fabric Access Policies

AEP also link ACI domains (and VLAN pools) to the physical infrastructure through Interface and Switch Selector Profiles, thereby defining the scope of the VLAN pool on the physical infrastructure. VLAN pools and Domains are access policies that can be reused across multiple switches and by linking an AEP to Interface and Switch Profiles as shown below, AEP specify the VLAN range allowed on the switches and interfaces. When deploying an EPG on a leaf port, the VLAN specified must be part of a pool (through the domain) linked to an AEP for the VLAN to be operational.
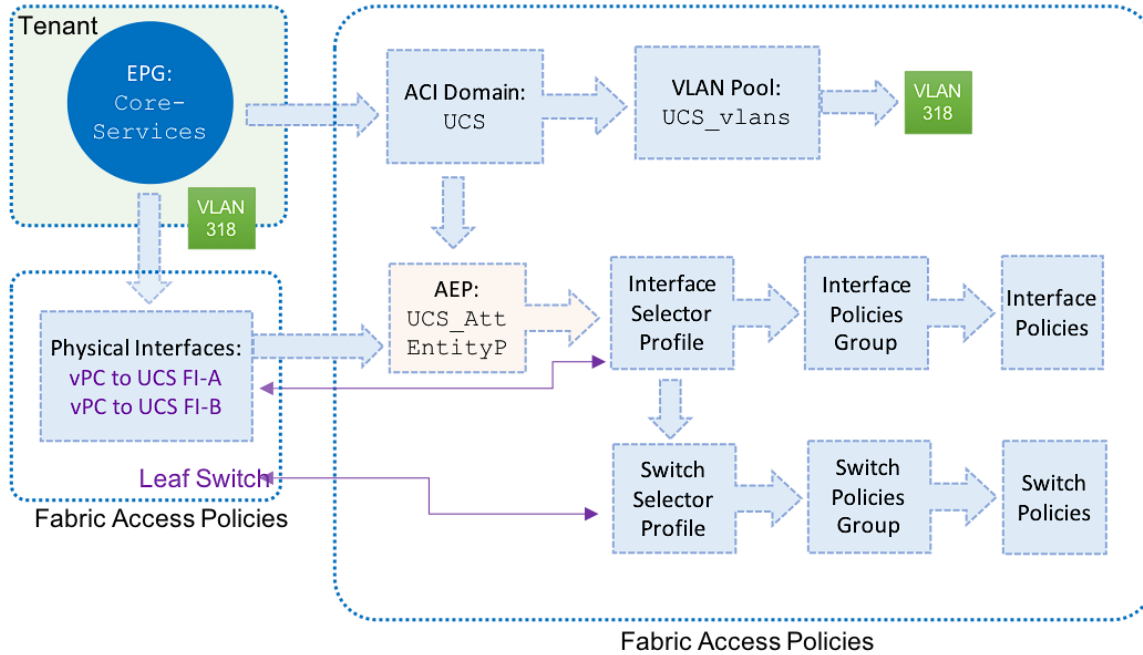
Figure 23   AEP Mapping to VLAN Pools and Domains for UCS Domain



Fabric Access Policies

In VMM domains, the associated AEP provides the interface profiles and policies (CDP, LACP) for the virtual environment.

In summary, an EPG must be deployed on a physical interface (port, PC, vPC) to send/receive traffic and an AEP must be mapped to a VLAN pool (through domain) in order to enable the EPG VLAN on that interface. The figure below shows the relationships and mappings for the Core-Services EPG in the FlexPod design.

**Figure 24    EPG to AEP Mapping for Core-Services EPG**



## AEP Design Guidelines and Best Practices

The design guidelines and best practices for AEPs in ACI and their use in the FlexPod design are as follows.

- Multiple domains can be associated with a single AEP. In the FlexPod design, the AEP associated with interfaces going to the UCS domain are linked to two domains as shown in the Table below.

- Multiple AEPs are required to support overlapping VLAN pools, enable/disable infrastructure VLAN, to limit the scope of VLANs or to support different virtual switch policies associated with a VMM domain. In the FlexPod design, multiple AEPs are to provide flexibility in the design. There is a one-to-one mapping between AEPs and Domains and VLAN Pools in this design - see Table below.

- Note that AEPs can be easily provisioned using the access port configuration wizard (Create an interface, PC, and vPC). However, in order to fully realize the benefits of policy-driven architecture that ACI provides and to design for fabric-wide policy-reuse, requires an understanding of ACI constructs and the inter-dependencies.

The table below shows the AEP to Domain Mapping in the FlexPod design.

Table 11    AEP to Domain Mapping

| Access Connection | VLAN Pool | Domain Type | AEP |
|---|---|---|---|
| vPC to Management Switch | `FP-Mgmt-Sw` `(118)` | `External Bridged` `Devices(`**L2OUT**`)` | `__ui_I105_106)FP-MGMT-Sw-` `Ports-21-1` |
| vPC to NetApp AFF Cluster | `NetApp-AFF_vlans` `(219,3010,3020,3050)` | `Bare Metal(`**Physical**`)` | `NetApp-AFF_AttEntityP` |
| vPC to Cisco UCS Fabric Interconnects | `UCS_vlans` `(2,318,419,3000,3110,` `3120,3150)` | `External Bridged` `Devices` | `UCS_AttEntityP` |
| Redundant Links to a pair of L3 Gateways | `Shared-L3-Out_vlans` `(201,202,203,204)` | `External Routed` `Domain(`**L3OUT**`)` | `FP-Shared-L3-` `Out_AttEntityP` |
| VMM Domain | `VP-fpv-vc-vDS` `(1100-1199)` | `VMM Domain` | `UCS_AttEntityP` |

## ACI Fabric Design – EPGs, Contracts and Application Profiles

### End Point Groups

An End Point Group (EPG) in ACI is a logical entity that contains a group of endpoints. Endpoints can be physical or virtual and connect directly or indirectly to the ACI fabric. Endpoints have an identity (address) and a location. The endpoint grouping is not tied to the physical or logical topology – they can be distributed across the fabric. The grouping can be based on requiring a common set of policies or services or providing a common set of services or other functions. By grouping them, the endpoints can be managed as a group rather than individually.

For example, the Core-Services EPG in the FlexPod design is formed on the basis of providing common infrastructure services such as AD, DNS, DHCP etc. Another EPG, IB-MGMT is based on being part of a common management network for ESXi hosts .

In the FlexPod design, various application tiers, host interface ports for iSCSI, NetApp LIFs for SVM-Management,  NFS volumes and iSCSI datastores are all placed in separate EPGs. The table below shows the EPGs defined in the FlexPod design. Application EPGs are not listed but can be added as needed.

Table 12    FlexPod Design – EPGs

| EPG | Purpose |
|---|---|
| Core-Services | For Virtual Endpoints that provide core services such as AD, DNS, DHCP |
| IB-MGMT | For ESXi hosts – interface to management network. Also for VMs that need access to it. |
| SVM-MGMT | For NetApp's Infrastructure SVM management interfaces. Also for VMs that need access to it. |
| vMotion | For VMware ESXi hosts – interface to vMotion network |
| iSCSI-A | For ESXi hosts and NetApp Controllers – iSCSI Path A interface (to provide and access iSCSI LUNs) |
| iSCSI-B | For ESXi hosts and NetApp Controllers – iSCSI Path B interface (to provide and access iSCSI LUNs) |
| NFS | For ESXi hosts and NetApp Controllers – interface to NFS network (to provide and access NFS volumes) |

EPGs can be static or dynamic depending whether the endpoints are added to the EPG using a static binding or dynamically. Static EPG is where an endpoint is added to the EPG by an Administrator. Addition of virtual machines to an EPG is an example of a dynamic EPG. As new VMs come online, vCenter will triggers the addition of VMs to EPGs. See Fabric Access Design – Access Policies (VLAN Pools) and VLAN Design section for more details.

## Application Profiles

An application profile models application requirements and contains one or more EPGs as necessary to enable multi-tier applications and services. The FlexPod design uses multiple application profiles to define multi-tier applications as well as to establish storage connectivity. Table 13  shows the Application Profiles and EPGs used in this FlexPod design.

## Contracts

Contracts define inbound and outbound traffic filter, QoS rules and Layer 4 to Layer 7 redirect policies. Contracts define the way an EPG can communicate with another EPG(s) depending on the application requirements. Contracts are defined using provider-consumer relationships; one EPG provides a contract and another EPG(s) consumes that contract. Contracts utilize filters to limit the traffic between the applications to certain ports and protocols. Table 13  shows the EPG contracts used in this FlexPod design.

Table 13    FlexPod Design – EPGs, Application Profiles and Contracts

| Application Profile | EPG | Contracts | Notes |
|---|---|---|---|
| **Core-Services** | **Core-Services** | `common-Allow-Core-Services` (Provider) | Provides access to Core-Services VMs |
| **IB-MGMT** | **IB-MGMT** | `common-Allow-Core-Services` (Consumer) | Enables ESXi hosts and management VMs access to Core-Services VMs |
| | **SVM-MGMT** | `common-Allow-Core-Services` (Consumer) | Enables NetApp SVM Management Interfaces to access Core-Services VMs |
| **Host-Conn** | **vMotion** | --- | No access to endpoints outside the EPG |
| | **iSCSI-A** | --- | No access to endpoints outside the EPG |
| | **iSCSI-B** | --- | No access to endpoints outside the EPG |
| | **NFS** | --- | No access to endpoints outside the EPG |

## Fabric Access Design – Multi-tenancy, VRFs and Bridge Domains

### Tenants

Tenant is a logical container for grouping applications and their networking and security policies. This container can represent an actual tenant, an organization, an application or a group based on some other criteria. Tenants also enable domain-specific management through a Tenant Administrator. A tenant represents a unit of isolation from a policy perspective. All application configurations in Cisco ACI are part of a tenant. Tenants can include multiple Layer 3 contexts or VRFs, multiple bridge domains per VRF, and EPGs to provide additional segmentation within the bridge domains.

ACI provides two categories of tenants: User Tenants and System Tenants. System Tenants include Common, Infra and Mgmt Tenants. For more information on system tenants, see ACI Fundamentals listed in the Solution References section.

The FlexPod design uses the Common tenant to host core services such as AD, DNS, etc. ACI provided Common Tenant is designed for shared services that all tenants need.

This design also includes a user-tenant called FPV-Foundation to provide:

- compute to storage connectivity for SAN boot, to access iSCSI LUNs

- compute to storage connectivity for accessing Infrastructure datastores using NFS

- access to vMotion network

- access to ESXi and SVM management networks

To deploy applications, it is expected that additional tenants will be created.

### VRF

A virtual routing and forwarding (VRF) instance in ACI is a tenant network. VRF is a unique Layer 3 forwarding domain. A tenant can have multiple VRFs and a VRF can have multiple bridge domains. In the FlexPod design, a single VRF is created in each tenant created.

### Bridge Domains

A bridge domain represents a Layer 2 forwarding construct within the fabric. The bridge domain can be further segmented into EPGs. Like EPGs, bridge domains can also span multiple switches. A bridge domain can contain multiple subnets, but a subnet is contained within a single bridge domain. One or more bridge domains together form a tenant network. A bridge domain represents the broadcast domain and has global scope. EPGs must be associated with bridge domain.

Bridge Domains are an important consideration in the FlexPod design. When a bridge domain contains endpoints belonging to different VLANs (outside the ACI fabric), a unique MAC address is required for every unique endpoint. NetApp storage controllers, however, use the same MAC address for an interface group and all the VLAN interface ports defined for that interface group on that storage node. As a result, all the LIFs on a NetApp interface group end up sharing a single MAC address even though these LIFs belong to different VLANs.

To overcome potential issues caused by overlapping MAC addresses, multiple bridge domains need to be deployed for correct storage connectivity.

### FlexPod Design – Tenants, VRFs and Bridge Domains

The table below shows the Tenants, VRFs and Bridge Domains in the FlexPod design.

Table 14    FlexPod Design – Tenants, VRFs and Bridge Domains

| Tenant | VRF | Bridge Domain |
|---|---|---|
| Common | vrf-FP-Common-IB-MGMT | BD-FP-common-Core-Services |
| FPV-Foundation | FPV-Foundation | BD-FPV-Foundation-Internal |
| | | BD-FPV-Foundation-iSCSI-A |
| | | BD-FPV-Foundation-iSCSI-B |
| | | BD-FPV-Foundation-NFS |

## Fabric Access Design – Virtual Machine Manager Integration and Networking

Integration with the Virtual Machine Manager (VMM) or VMware vCenter managing the virtual environment allows ACI to control the virtual switches running on the ESXi hosts and extend the fabric access policies to these switches. The integration also automates the deployment tasks associated with connecting a virtual switch to the ACI fabric.

### Virtual Machine Networking

VMM integration allows Cisco APIC to control the creation and configuration of the virtual switches running on the hypervisor. The virtual switch can be a VMware vSphere Distributed Switch (VDS) or a Cisco ACI

Virtual Edge (AVE). Once the virtual distributed switches are deployed, APIC communicates with the switches to create port-groups that correspond to EPGs in the ACI fabric and send network policies to be applied to the virtual machines.

### VMM Domain Profiles

A VMM Domain profile defines a VMM domain and specifies the connectivity policies for connecting VMware vCenter to the ACI fabric. VMMs with common policies can be mapped to a single domain. The access credentials included in the domain enable APIC to VMM communication. This communication channel is used for querying vCenter for info such as VM names and vNIC details, to create port-groups, to push VM policies and to listen for VM events such as VM creation, deletion etc. A VMM domain provides VM mobility within the domain and support for multi-tenancy within the fabric.

### VLAN Pools

VLAN Pool is a shared resource that can be shared by multiple domains including a VMM domain. However a VMM domain can only be associated with one VLAN pool. VLAN pool specifies the allowed range of VLANs in the VMM domain. Typically, VLANs in a VMM domain pool are dynamically assigned to an EPG by the APIC. APIC also provisions the VMM domain VLAN on the leaf ports based on EPG events, either statically or based on VM events from vCenter.

### Endpoint Groups

A VMM Domain can be associated with multiple EPGs and an EPG can span multiple VMM domains. APIC creates a port-group for every EPG associated with a VMM domain in the ACI fabric. A VMM domain with multiple EPGs will have multiple port groups defined on the virtual distributed switch. To position an application, the application administrator deploys the VMs on VMware vCenter and places the VM NIC(s) into the appropriate port group for the application tier.

### Attachable Entity Profile

AEPs associated the domain (and VLAN pool) to the physical infrastructure and can enable VMM policies on leaf switch ports across the fabric.

### VMM Integration with Cisco UCS Blade Server

The blade server architecture factors into the VMM integration, where VMs are deployed on Cisco UCS Blade servers connected through a pair of Fabric Interconnects. VMM integration with Cisco UCS B-series servers requires the following:
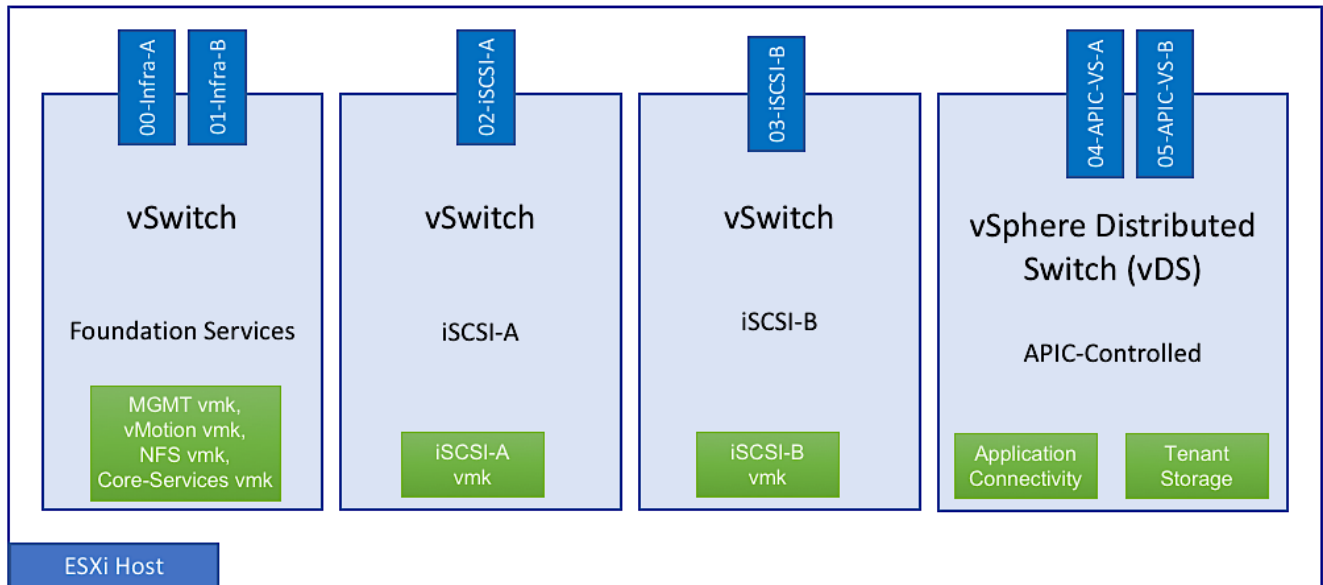
- Cisco UCS vNICs must be configured to use CDP or LLDP - both cannot be enabled.

- The VLAN pool or the allowed range of VLANs for the VMM domain must be allowed on the UCS vNICs and on the FI uplinks connecting to the leaf switches

### Virtual Switching Architecture

In the FlexPod design, tenant applications are deployed on port-groups in the APIC-controlled distributed switch (VMware vDS) but for other infrastructure connections such as vCenter access using in-band management, vSphere vMotion and iSCSI storage access, a vSphere vSwitch per connection is used. To support this multi-vSwitch environment, multiple vNIC interfaces are setup the services profiles for the UCS hosts and the VLANs for storage, management and vMotion are then enabled on the appropriate vNIC interfaces. Figure below shows the distribution of VMkernel ports and VM port-groups on an iSCSI

connected ESXi server. For an ESXi server, supporting iSCSI based storage access, In-band management and vMotion traffic is handled by a Foundation Services vSwitch and iSCSI-A and iSCSI-B traffic is handled by two dedicated iSCSI vSwitches. The resulting ESXi host configuration is therefore a combination of 3 vSwitches and a single APIC-Controlled distributed switch which handles application (tenant) specific traffic.

**Figure 25  FlexPod – Virtual Switching Design**



## Enabling Access to Core Services and Storage Management

### Access to Core Services

To provide ESXi hosts and VMs access to an existing management segment where core-services such as Active Directory (AD), Domain Name Services (DNS), etc. reside, inter-tenant contracts are utilized. Cisco ACI fabric provides a pre-defined tenant named common to host the common services that can be easily shared by other tenants in the system. The policies defined in the common tenant are usable by all the tenants without any special configurations. By default, in addition to the locally defined contracts, all tenants **in the ACI fabric can "consume" the contracts "provided" in the** common tenant.

In the FlexPod design, access to core services is provided as shown in the figure below. To provide this access:

- A common services segment is defined where core services VMs connect. The Core-Services EPG is associated with the APIC-controlled virtual switch VMM domain.  A separate services segment ensures that the access from the tenant VMs is limited to only core services' VMs.

- A static port mapping is used to link a separate, isolated In-Band management network to Core-Services.

- The EPG for the core services segment Core-Services is defined in the common tenant.

- The tenant VMs access the core services segment by consuming contracts from the common tenant.

- The contract filters can be configured to only allow specific services related ports.

- The tenant VMs access the core services segment using their EPG subnet gateway.

- Since the tenant VMs reside in separate subnets than the Core-Services VMs, routes must be configured in the Core-Services VMs and hosts to reach the Application tenant VMs.  For this lab implementation a supernet route with destination 172.18.0.0/16 was put into each Core-Services VM.  Routes needed to be shared across VRFs since the two EPGs were in different tenants.

- Unique IP subnets have to be used for each EPG connecting to Core-Services
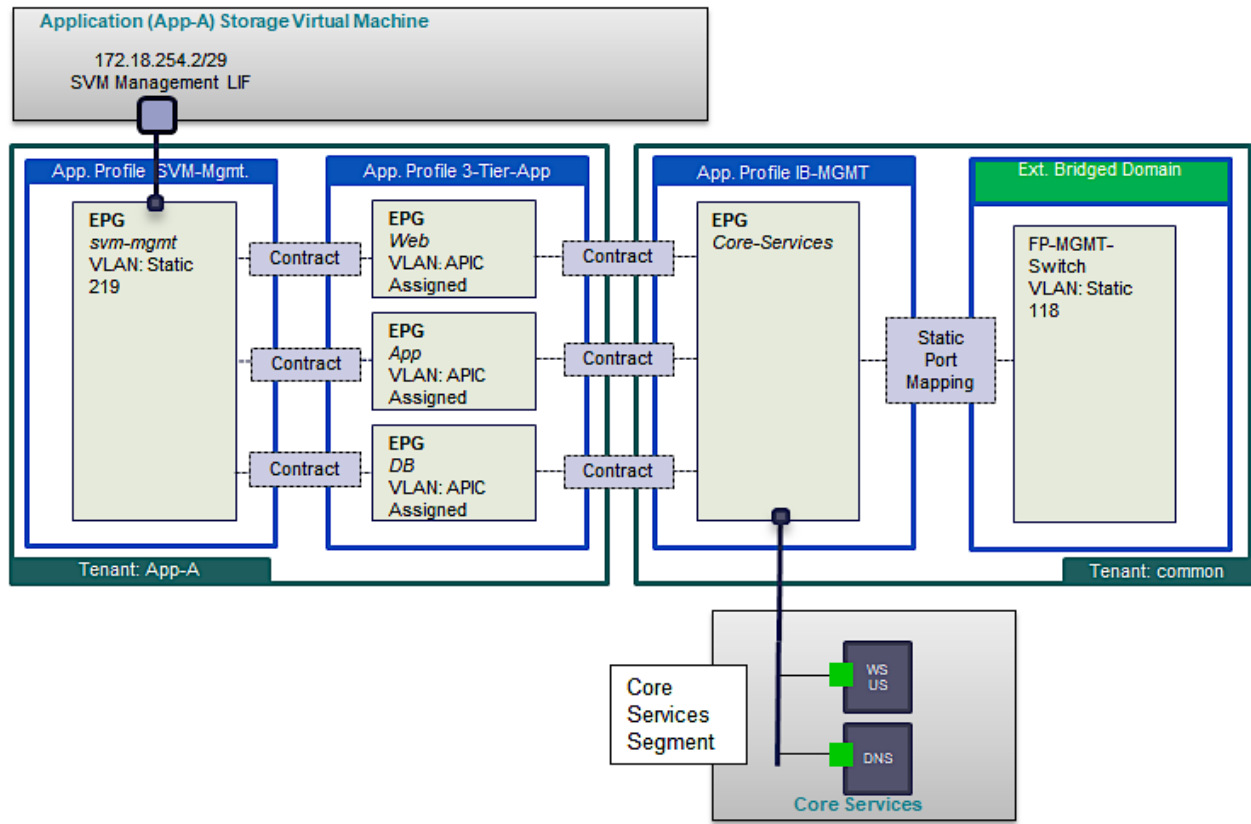
## Accessing SVM Management

Some applications such as NetApp Snap Drive require direct connectivity from the application (SharePoint, Exchange, SQL, etc.) VMs to the management LIF on the tenant SVM. To provide this connectivity securely, a separate VLAN is dedicated for each tenant to define the management LIF. This VLAN is then statically mapped to an EPG in the application tenant as shown in Error! Reference source not found.. Application Ms can access this LIF by defining and utilizing contracts. The SVM management LIF can also be connected to Core-Services to allow tenant ESXi hosts to use Snapdrive to configure tenant CSVs.

> Note:  When an application tenant contains mappings for NetApp LIFs for storage access (iSCSI, NFS etc.), a separate bridge domain is required for the SVM management LIF because of the overlapping MAC addresses. Ensure that only one type of storage VLAN interface is connected to a given bridge domain.

The figure below shows both "provider" Core-Services EPG in the tenant common and the consumer EPGs "Web", "App" and "DB" in tenant "App-A".

Figure 26   Access Core-Services and Storage Management



## Onboarding Infrastructure Services

In an ACI fabric, all the applications, services and connectivity between various elements are defined using ACI tenants, application profiles, bridge domains and EPGs. In the FlexPod design, the tenant configured to provide infrastructure services is called FPV-Foundation. The FPV-Foundation tenant provides the following services and connectivity:

- Provides ESXi hosts and VMs access to existing management network outside the ACI fabric for core services such as AD, DHCP, DNS etc.

- Enables core services VMs to reach the NetApp Infrastructure SVM management interface

- Enables connectivity between hosts for vMotion

- Compute to Storage Connectivity to access ISCSI datastores

- Compute to Storage Connectivity to access NFS datastores

### ACI Design for Foundation Tenant

The application profiles, EPGs and Bridge Domains in the FPV-Foundation tenant are shown in Table 15  .

Table 15    FPV-Foundation Tenant – Application Profiles, EPGs and Bridge Domain

| Tenant/VRF | Application Profile | EPG | Bridge Domain Associated with EPG |
|---|---|---|---|
| FPV-Foundation/ FPV-Foundation | IB-MGMT | IB-MGMT | BD-FP-common-Core-Services(Common **Tenant**) |
| | | SVM-MGMT | BD-FPV-Foundation-Internal |
| | Host-Conn | vMotion | BD-FPV-Foundation-Internal |
| | | iSCSI-A | BD-FPV-Foundation-iSCSI-A |
| | | iSCSI-B | BD-FPV-Foundation-iSCSI-B |
| | | NFS | BD-FPV-Foundation-NFS |

The  FPV-Foundation tenant consists of single VRF since there is no overlapping IP address space requirements in this design.

The bridge domains (BD) in the FPV-Foundation tenant are:

- BD-FPV-Foundation-iSCSI-A: BD associated with EPGs for iSCSI-A traffic

- BD-FPV-Foundation-iSCSI-B: BD associated with EPGs for iSCSI-B traffic

- BD-FPV-Foundation-NFS: BD associated with EPGs for NFS traffic

- BD-FPV-Foundation-Internal: This BD hosts EPGs for all other FPV-Foundation tenant traffic

While all the EPGs in a tenant can theoretically share the same bridge domain, overlapping MAC address usage by NetApp storage controllers on the interface groups across multiple VLANs determines the actual number of bridge domains required. As shown in Figure 27, the FPV-Foundation tenant connects to two iSCSI LIFs and one NFS LIF to provide storage connectivity to the infrastructure SVM. Since these three LIFs on each storage controller share the same MAC address, a separate BD is required for each LIF.

The two application profiles in this tenant are:

- IB-MGMT: This application profile provides connectivity to existing management network through the Common tenant (details covered later in this section) and reachability to NetApp SVM Management interface.

- Host-Conn: This application profile supports compute to storage connectivity and vMotion.

The ACI constructs for infrastructure services including an overview of the connectivity and relationship between various ACI elements is covered in the figures below.

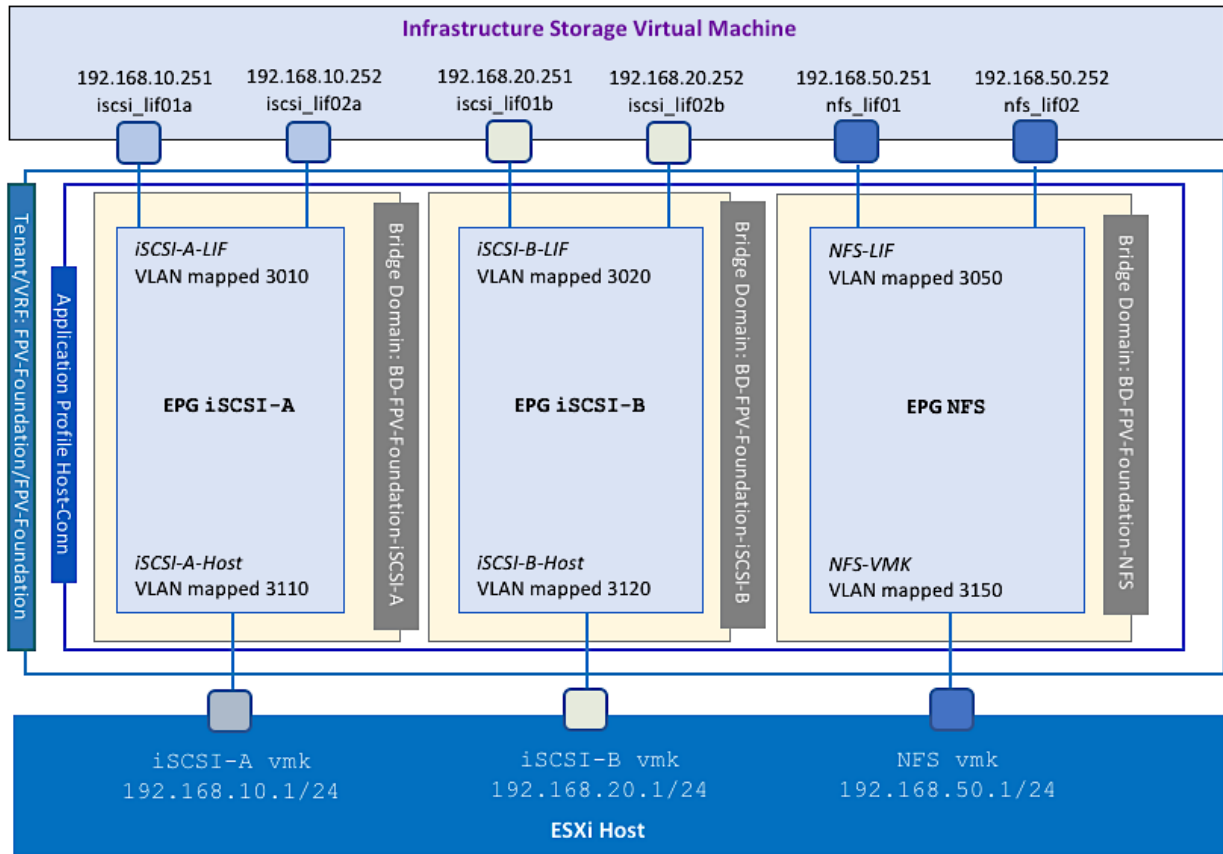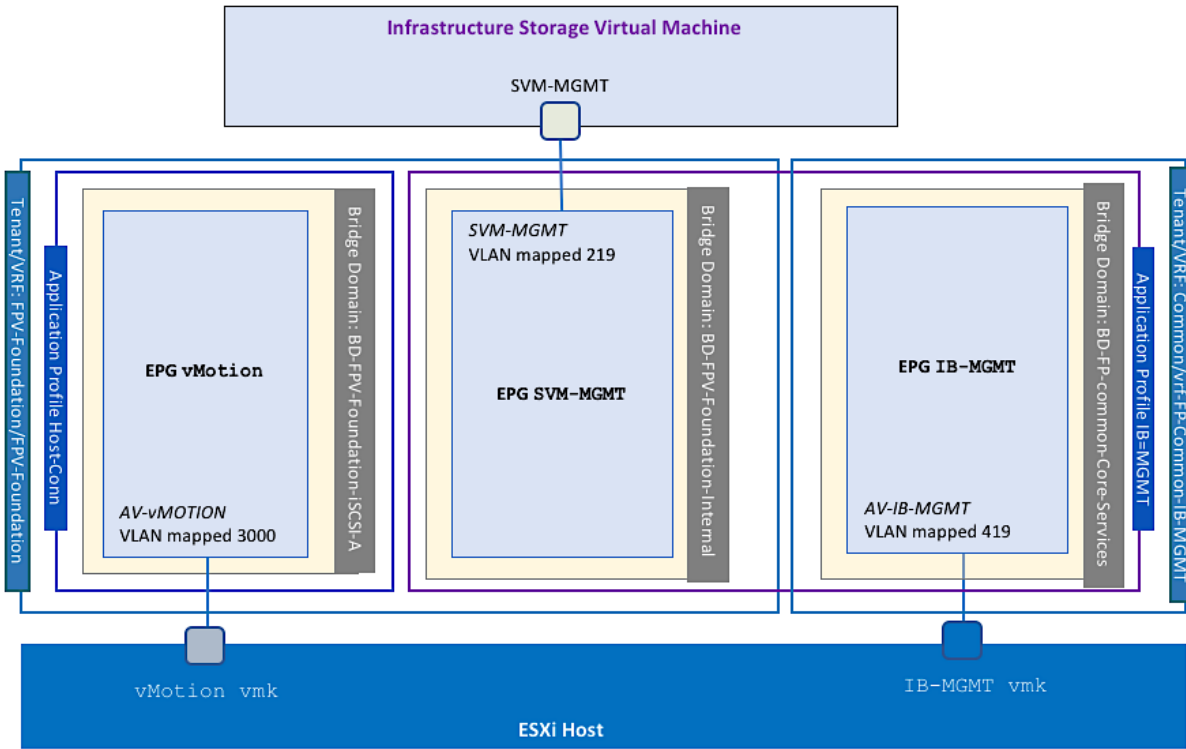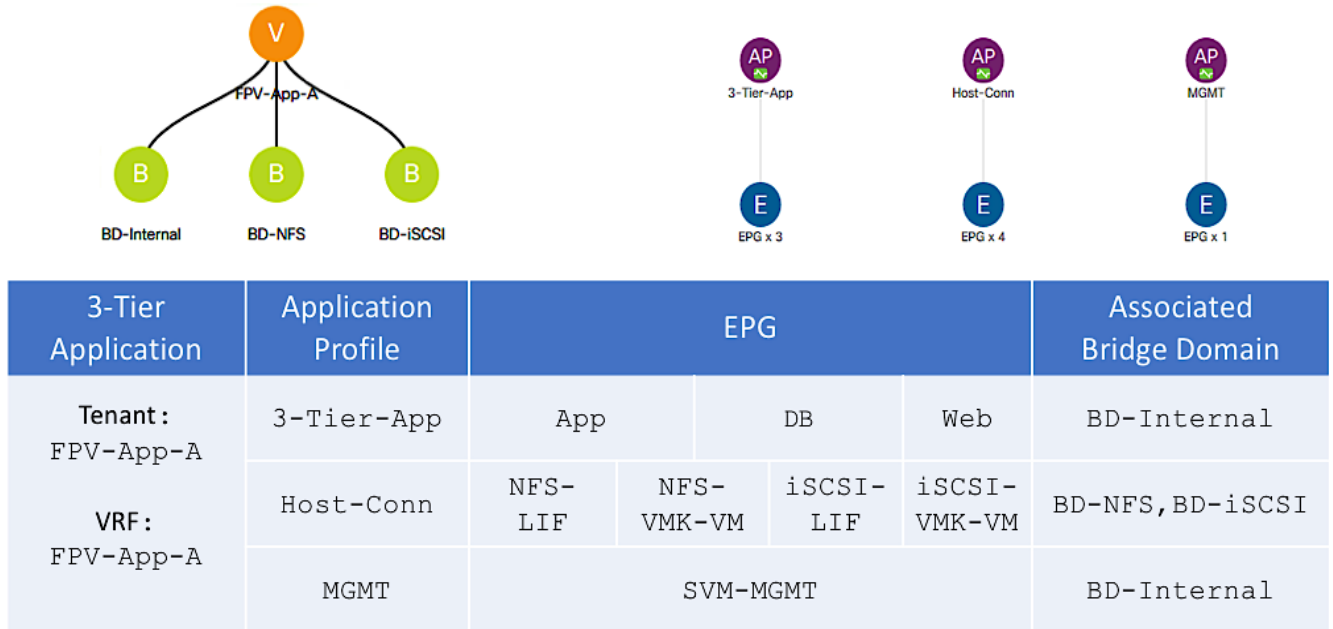Figure 27   Foundation Tenant Design – Compute to Storage Connectivity

Figure 28   Foundation Tenant Design – vMotion and Management



## Onboarding Multi-Tier Applications

The ACI constructs for a multi-tier application deployment include defining a new tenant, VRF(s), bridge domain(s), application profile(s), end point group(s), and the contract(s) to allow communication between various tiers of the application. Deploying a sample 3-tier application requires a new tenant, VRF, Bridge Domains, Application Profiles and EPGs as shown below.

Figure 29    FlexPod Design – ACI Constructs for Deploying a 3-Tier Application



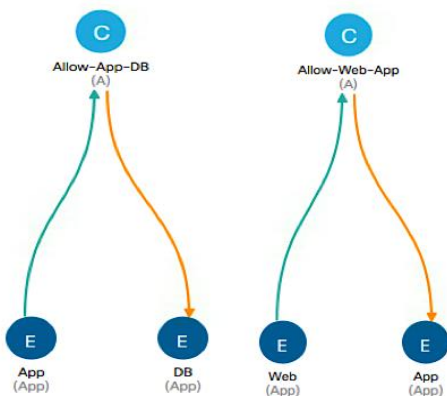| 3-Tier Application | Application Profile | EPG | | | | Associated Bridge Domain |
|---|---|---|---|---|---|---|
| Tenant: FPV-App-A | 3-Tier-App | App | | DB | Web | BD-Internal |
| VRF: FPV-App-A | Host-Conn | NFS-LIF | NFS-VMK-VM | iSCSI-LIF | iSCSI-VMK-VM | BD-NFS,BD-iSCSI |
| | MGMT | SVM-MGMT | | | | BD-Internal |

Multiple Bridge Domains are necessary under the Tenant and VRFs. As stated earlier, though all EPGs could be under a single Bridge Domain, the use of overlapping MAC addresses by NetApp storage controllers across multiple interface groups and VLANs, requires separate bridge domains for iSCSI, NFS and other traffic (SVM-MGMT and Application EPGs). Application VMs do not use overlapping mac-addresses so they can share the BD for SVM-MGMT.

The ACI constructs for a 3-tier application deployment is a little more involved than deploying the infrastructure tenant (FPV-Foundation) covered in the last section. In addition to providing ESXi host to storage connectivity, various tiers of the application also need to communicate amongst themselves as well as with the storage and common or core services (DNS, AD etc.). Appropriate contracts are provided to enable hosts access hosts to storage and to allow traffic between various application tiers as outlined below.
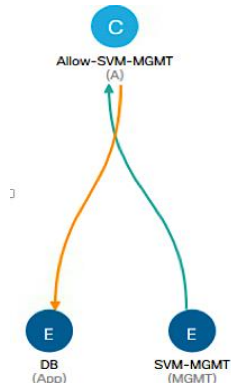
- Web, DB, and App EPGs are deployed in the VMM domain to provide a VM network for Web, Application and Database servers. Contracts are defined to allow traffic between the different tiers.

Figure 30    3-Tier Application: Contracts between Application Tiers

61

- Contracts are provided by the SVM-MGMT EPG to the application to tiers to enable access to SVM for storage provisioning and backup if necessary.

**Figure 31   3-Tier Application: Contracts to Allows Access to SVM for Storage Provisioning**



- NFS and iSCSI Contracts are also provided to allow the different application tiers access to NFS and iSCSI datastores

**Figure 32   3-Tier Application: Contract to allow access NFS and iSCSI datastores**
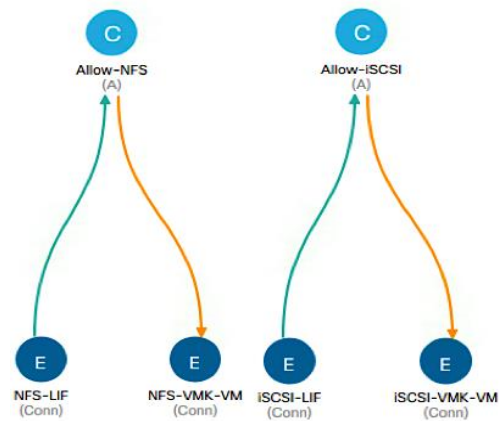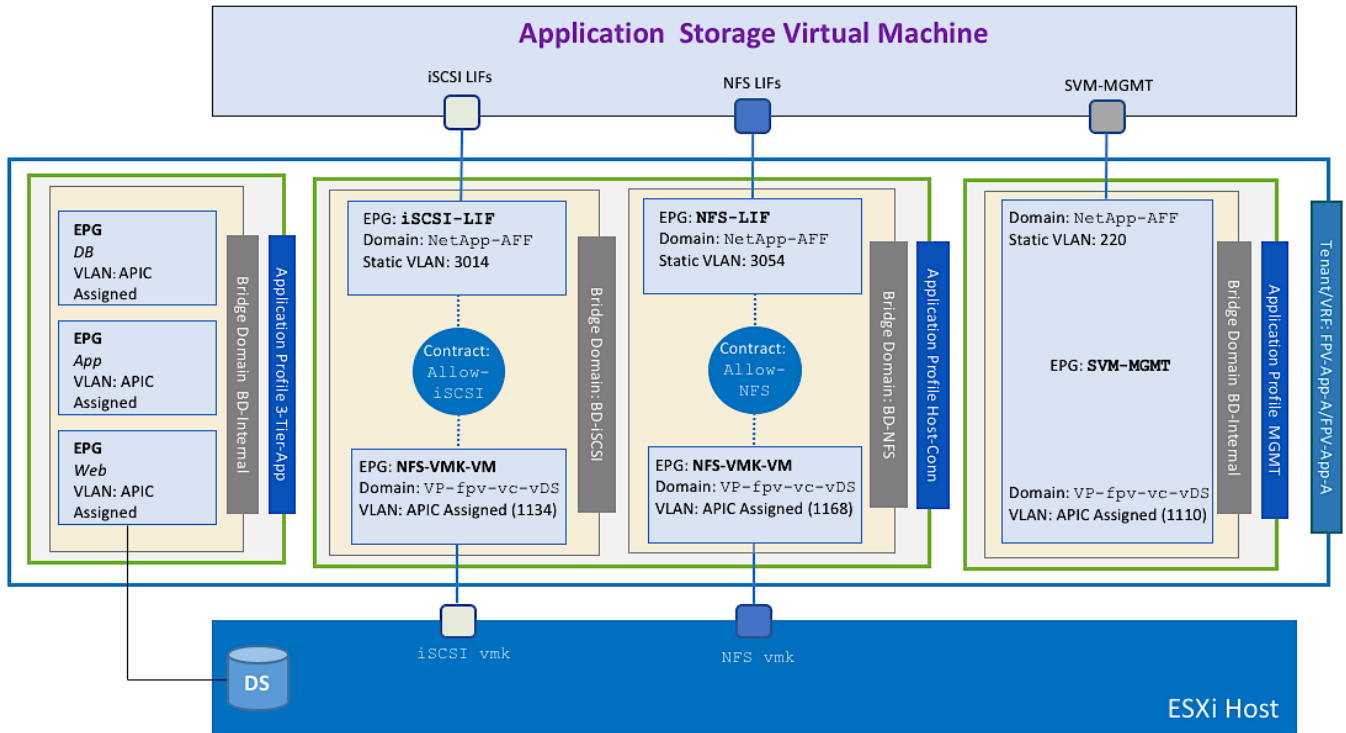


Figure 33 provides an overview of the constructs, the connectivity and the relationship between the various ACI elements for a 3-tier Application.

Figure 33    FlexPod Design – Sample 3-Tier Application



## External Network Connectivity - Shared Layer 3 Out

In order to connect the ACI fabric to existing infrastructure, the leaf nodes are connected to a pair of core infrastructure routers/switches. In this design, a Cisco Nexus 7000 was configured as the core router. Figure 34 shows the connectivity details from the Shared_L3_Out External Routed Domain in the "common" tenant and the "common/default" VRF.  Figure 35 shows how tenants with other VRFs are connected to the Shared_L3_Out via contracts. Tenant network routes can be shared with the Nexus 7000s using OSPF and external routes from the Nexus 7000s can be shared with the tenant VRFs. Routes can also be shared across VRFs.

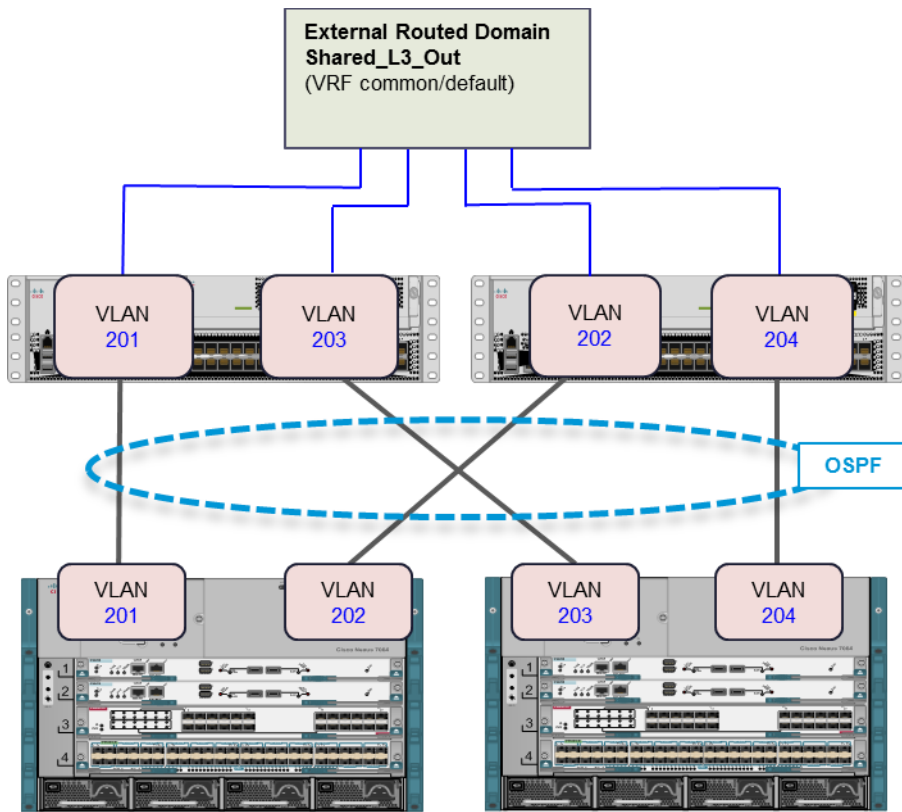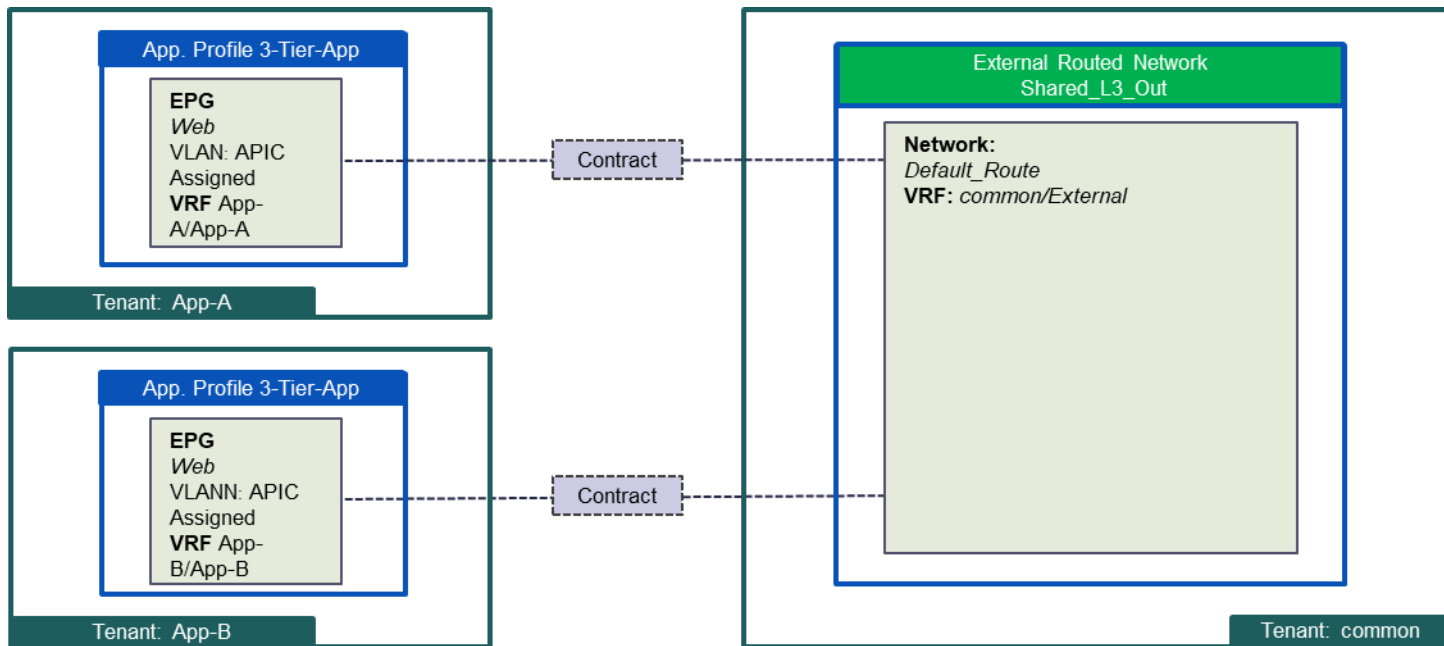Figure 34   ACI Connectivity to Existing Infrastructure



Figure 35   Connecting Tenant Networks to Shared Layer 3 Out



Some of the design principles for external connectivity are as follows:

- Each Leaf switch is connected to both Cisco Nexus 7000 switches for redundancy.

- A unique VRF is defined for every tenant. App-A/App-A and App-B/App-B are two such VRFs shown in Figure 35.

- Unique VLANs are configured to provide multi-path connectivity between the ACI fabric and the core infrastructure. VLANs 201-204 are configured as shown.

- On the ACI fabric, OSPF is configured to share routes. The ACI fabric learns a default route from the core router and each tenant advertises one or more "public" routable subnets to the core infrastructure across VRFs.

## Enabling Firewall Services for Tenant Applications

To provide L4-L7 services for multi-tier applications, ACI provides a L4-L7 VLAN stitching feature. In this design, two Cisco ASAv virtual firewalls are deployed in the common Tenant. The ASAv firewall devices connect into the ACI Fabric using the APIC-integrated VMware vDS. Inside and Outside EPGs and corresponding vDS port-groups are created for the inside and outside segments of a firewall.

When a tenant application requires firewall services, a L4-L7 service graph with the ASAv instance is deployed, resulting in new port-groups in the tenant and the ASAv network interfaces are moved to the new port-groups. In this design, firewalls are used between the Tenant Web EPG in the multi-tier application and the ACI fabric's shared L3Out interface.

To provide management connectivity to Cisco ASAv, the VM's management interfaces are added to the SVM MGMT EPG and managed from Core Services EPG. A previously established contract is used between SVM-MGMT EPG and Core-Services EPG to enable management access to ASAv.

The Inside and Outside Interfaces on ASAv are initially placed in a single bridge domain in the common Tenant but moved to separate Bridge Domains when a tenant application deploys a L4-L7 service graph to enable firewall services. Note that the Inside and Outside interfaces must connect to different Bridge Domains when in use.  In this case, since the outside interface of the firewall connects to shared L3Out, the Outside bridge domain is in the same common Tenant as shared L3Out. The inside interface will be in the Web EPG and therefore the same bridge domain.

# Other Design Considerations

## High Availability

FlexPod Datacenter with ACI is designed, at every layer – compute, network and storage, to be fully redundant. There is no single point of failure from a device or traffic flow perspective.  Link aggregation technologies continue to play an important role in an ACI-based design as it did in previous FlexPod designs based on NX-OS standalone mode.  Link aggregation or port channeling (PC) based on 802.3ad standards is utilized across all layers of the FlexPod design to provide higher bandwidth and resiliency on the uplinks. Cisco Unified System, NetApp storage arrays and Cisco Nexus 9000 series switches all support port channeling using Link Aggregation Control Protocol (LACP) to aggregate and load-balance traffic across multiple links. Cisco Nexus 9000 series extend this technology and support virtual Port Channels (vPC) to provide resiliency at the node-level, in addition to the link level resiliency that a PC provides. Typically, vPC requires more configuration than a PC but in ACI, vPC configuration is simpler.
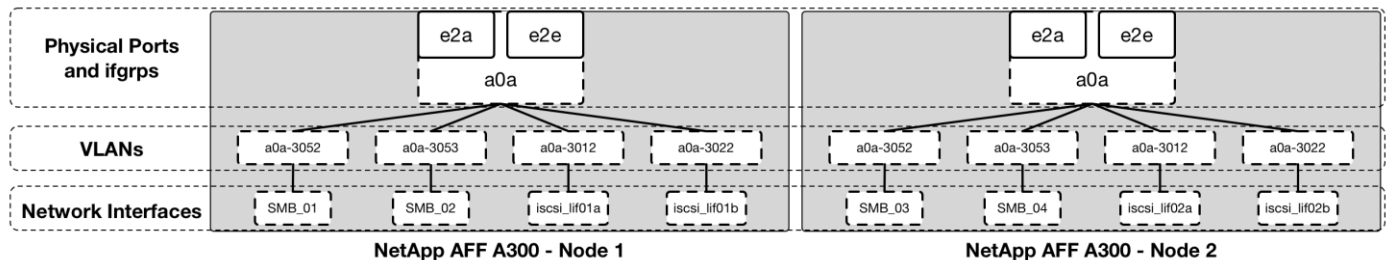
## SAN Boot

SAN boot for the Cisco UCS servers is considered a best practice in the FlexPod Datacenter solution. SAN boot takes advantage of the stateless compute capabilities of Cisco UCS, and enables the operating system to be safely secured by the NetApp All Flash FAS storage system, providing better performance and security. In this design, SAN boot was validated with iSCSI and Fibre Channel (FC). FCoE is also supported.

In the iSCSI validation:

- The physical network data ports on each NetApp controller were combined into LACP ifgroups (one ifgroup per controller), for Layer 1 connectivity

- Two VLANs (one for each iSCSI fabric) was created on the ifgroup on both NetApp controllers, for Layer 2 connectivity

- Four logical interfaces (LIFs; one for each iSCSI fabric, on each controller) was created for Layer 3 connectivity

- The iSCSI traffic passes through the Cisco ACI fabric, and ACI policy can be applied to iSCSI traffic. Each ESXi host ends up with 2 active optimized paths, and 2 active un-optimized paths to its boot LUN. This network interface and port layout can be seen in Figure 36.

Figure 36  iSCSI SAN Boot



In the FC validation, there are no ifgroups and no VLANs; there are two physical FC fabrics, and each controller has a logical FC interface on each fabric, for a total of 4 FC network interfaces. Each ESXi host ends up with 2 active optimized paths, and 2 active un-optimized paths to its boot LUN. The separate FC ports were directly connected to the Cisco UCS Fabric Interconnects. FC traffic does not pass through the ACI fabric and does not have ACI policy applied.

After boot, ESXi will recognize that its LUNs are on an ALUA-aware array, and will use its native ALUA multipathing (MPIO) software to manage its paths to its LUNs on the NetApp controllers. MPIO software is required in any configuration which provides more than a single path to the LUN. Because there are multiple paths to the LUN, SAN network interfaces are not configured to fail over to an available port. Instead, the network interface will be disabled and the host chooses a new optimized path on a different network interface. ALUA is an industry standard protocol for identifying optimized paths between a storage system and a host. ALUA enables the initiator (ESXi) to query the target (the NetApp controller) about path attributes, such as primary and secondary paths. ALUA can be enabled on an ONTAP interface group (igroup) and is automatically enabled for all configurations described in this guide.

## Cisco Unified Computing System – QoS and Jumbo Frames

FlexPod accommodates a myriad of traffic types (Cluster, SMB, iSCSI, control traffic, etc.) and is capable of absorbing traffic spikes and protect against traffic loss. Cisco UCS QoS system classes and policies can be configured to deliver this functionality. In this validation effort the FlexPod was configured to support jumbo frames with an MTU size of 9000. Enabling jumbo frames allows the FlexPod environment to optimize throughput between devices while simultaneously reducing the consumption of CPU resources.

**Note:** When setting up Jumbo frames, it is important to make sure MTU settings are applied uniformly across the stack to prevent packet drops and negative performance.

# Solution Validation

The FlexPod Datacenter solution was validated for data forwarding by deploying IOMeter VMs. The system was validated for resiliency by failing various pieces of the system under load. The tests executed include:

- Failure and recovery of network paths to AFF nodes, ACI switches, and fabric interconnects.

- SSD removal to trigger an aggregate rebuild.

- Storage link failure between one of the AFF nodes and the ACI fabric.

- Storage controller failure and takeover by surviving controller.

## Validated Hardware and Software

Table 16 describes the hardware and software versions used in validating this FlexPod Datacenter solution.

> Cisco, NetApp, and VMware have interoperability matrixes that should be referenced to determine support for any specific implementation of FlexPod. The matrixes can be accessed here.

Table 16   Validated Hardware and Software Versions

| Layer | Device | Image | Comments |
|---|---|---|---|
| Compute | Cisco UCS FI 6200 and 6300 Series, UCS B-200 M5, B-200 M4, UCS C-220 M4 with Cisco UCS VIC 1340 and 1385, 2304 IOMs | 3.2(3a) * (Infra  and Server Bundle) | Initial Validation on 3.2(2d) |
| Network | Cisco APIC | 3.1(1i) | |
| | Cisco Nexus 9000 ACI | n9000-13.1(1i) | |
| Storage | NetApp AFF A300 | ONTAP 9.3 | |
| Software | Cisco UCS Manager | 3.2(3a) * | Initial Validation on 3.2(2d) |
| | VMware vSphere | 6.5 Update 1 * | VMware vSphere Patches – Custom ESXi ISO with patches ESXi650-201803401-BG and ESXi650-201803402-BG applied after initial validation |
| | Cisco VIC nenic Driver | 1.0.13.0 | |
| | Cisco VIC fnic Driver | 1.6.0.36 | |

* Initial validation was completed using earlier versions of Cisco UCS and VMware releases. However, Cisco and VMware released Speculative Execution vulnerability (Spectre & Meltdown) patches in updated software releases (shown below) after validation was complete.  These patches and releases were installed and limited runs of validation tests were performed to check for continued behavior. Cisco recommends and supports the updated releases and patches for this CVD.

# Summary

FlexPod Datacenter with Cisco ACI and VMware vSphere 6.5U1 provides an optimal shared infrastructure foundation to deploy a variety of IT workloads that is future proofed with 40Gb/s iSCSI, NFS or 16 Gb/s FC, with either delivering 40Gb Ethernet connectivity. Cisco and NetApp have created a platform that is both flexible and scalable for multiple use cases and applications. From virtual desktop infrastructure to SAP®, FlexPod can efficiently and effectively support business-critical applications running simultaneously from the same shared infrastructure. The flexibility and scalability of FlexPod also enable customers to start out with a right-sized infrastructure that can ultimately grow with and adapt to their evolving business requirements.

# Solution References

## Compute

Cisco Unified Computing System:

http://www.cisco.com/en/US/products/ps10265/index.html

Cisco UCS 6300 Series Fabric Interconnects:

http://www.cisco.com/c/en/us/products/servers-unified-computing/ucs-6300-series-fabric-interconnects/index.html

Cisco UCS 5100 Series Blade Server Chassis:

http://www.cisco.com/en/US/products/ps10279/index.html

Cisco UCS 2300 Series Fabric Extenders:

https://www.cisco.com/c/en/us/products/collateral/servers-unified-computing/ucs-6300-series-fabric-interconnects/datasheet-c78-675243.html

Cisco UCS 2200 Series Fabric Extenders:
https://www.cisco.com/c/en/us/products/collateral/servers-unified-computing/ucs-6300-series-fabric-interconnects/data_sheet_c78-675243.html

Cisco UCS B-Series Blade Servers:

http://www.cisco.com/en/US/partner/products/ps10280/index.html

Cisco UCS C-Series Rack Mount Servers:

http://www.cisco.com/c/en/us/products/servers-unified-computing/ucs-c-series-rack-servers/index.html

Cisco UCS VIC Adapters:

http://www.cisco.com/en/US/products/ps10277/prod_module_series_home.html

Cisco UCS Manager:

http://www.cisco.com/en/US/products/ps10281/index.html

Cisco UCS Manager Plug-in for VMware vSphere Web Client:

http://www.cisco.com/c/en/us/td/docs/unified_computing/ucs/sw/vmware_tools/vCenter/vCenter_Plugin_Release_Notes/2_0/b_vCenter_RN_for_2x.html

## Datacenter Network Fabric

Cisco Nexus 9000 Series Switches

http://www.cisco.com/c/en/us/products/switches/nexus-9000-series-switches/index.html

Cisco Application Centric Infrastructure – Datacenter and Virtualization

https://www.cisco.com/c/en_au/solutions/data-center-virtualization/aci.html

Cisco Application Centric Infrastructure – Cisco Datacenter

https://www.cisco.com/go/aci

Cisco ACI Fundamentals
https://www.cisco.com/c/en/us/td/docs/switches/datacenter/aci/apic/sw/1-x/aci-fundamentals/b_ACI-Fundamentals.html

Cisco ACI Infrastructure Release 2.3 Design Guide

https://www.cisco.com/c/en/us/solutions/collateral/data-center-virtualization/application-centric-infrastructure/white-paper-c11-737909.pdf

Cisco ACI Infrastructure Best Practices Guide

https://www.cisco.com/c/en/us/td/docs/switches/datacenter/aci/apic/sw/1-x/ACI_Best_Practices/b_ACI_Best_Practices.html

## Storage

NetApp ONTAP 9:

http://www.netapp.com/us/products/platform-os/ontap/index.aspx

NetApp AFF A300:

http://www.netapp.com/us/products/storage-systems/all-flash-array/aff-a-series.aspx

NetApp AFF A-series controllers:

https://hwu.netapp.com/Controller/Index

NetApp OnCommand:

http://www.netapp.com/us/products/management-software/

NetApp SnapCenter:

http://www.netapp.com/us/products/management-software/snapcenter/

## Virtualization Layer

VMware vCenter Server:

http://www.vmware.com/products/vcenter-server/overview.html

VMware vSphere:

https://www.vmware.com/products/vsphere

## Interoperability Matrixes

Cisco UCS Hardware Compatibility Matrix:

https://ucshcltool.cloudapps.cisco.com/public/

VMware and Cisco Unified Computing System:

http://www.vmware.com/resources/compatibility

NetApp Interoperability Matrix Tool:

http://support.netapp.com/matrix/

FlexPod Technical Specification:

https://www.netapp.com/us/media/tr-4036.pdf

# About the Authors

John George, Technical Marketing Engineer, Data Center Solutions Engineering, Cisco Systems Inc.

John has been designing, developing, validating, and supporting the FlexPod Converged Infrastructure for over seven years. Before his roles with FlexPod, he supported and administered a large worldwide training network and VPN infrastructure. John holds a Master's degree in Computer Engineering from Clemson University.

Dave Derry, Technical Marketing Engineer, Converged Infrastructure Engineering, NetApp Inc.

Dave Derry is a Technical Marketing Engineer in the NetApp Converged Infrastructure Engineering team. He focuses on producing validated reference architectures that promote the benefits of end-to-end data center solutions and cloud environments.

## Acknowledgements

For their support and contribution to the design, validation, and creation of this Cisco Validated Design, the authors would like to thank:

- Archana Sharma, Cisco Systems, Inc.
- Haseeb Niazi, Cisco Systems, Inc.
- Ramesh Isaac, Cisco Systems, Inc.
- Lindsey Street, NetApp