CAI2 - Simplified Compute Operations

For the AI ready data center

Chris Merkel - Solutions Engineer - DC Architecture



Agenda

- 1. Intro
- 2. Platform Overview
- 3. Intersight Operations
- 4. Closing

Platform Overview

Cisco UCS Compute Portfolio - 2009

MAINSTREAM ENTERPRISE SERVERS

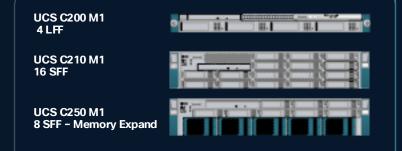
AI SERVERS

UCS 5108 Chassis

UCS B200 M1



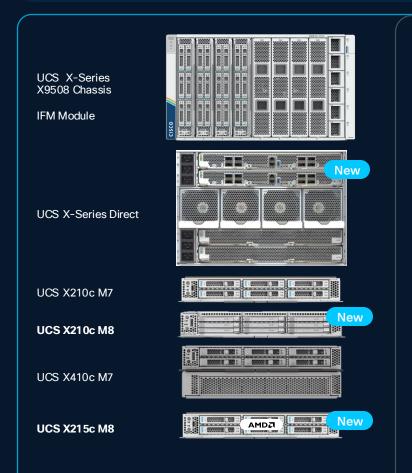


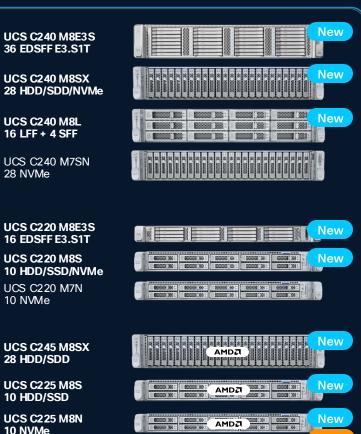




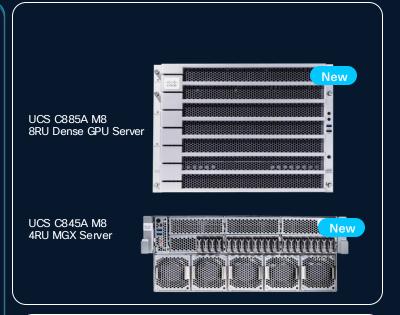
Cisco UCS Compute Portfolio

MAINSTREAM ENTERPRISE SERVERS





AI SERVERS



EDGE COMPUTING

Coming Late 2025!

Agenda

- 1. Intro
- 2. Platform Overview
- 3. Intersight Al operations
- 4. Closing

X-Series Portfolio

COMPUTE

X210c Compute Node

- 2-Socket, single slot servers
- · Three Generations: M6*, M7 and M8
- Intel 3rd Gen (Ice Lake) 4th Gen (Sapphire Rapids) 5th Gen (Emerald Rapids) and 6th gen (Granite Rapids) Xeon CPUs





X410c Compute Node

· Intel 4th Gen Xeon CPU

• Up to 64 DDR5 DIMMs

4-Socket, dual slot servers

X215c Compute Node

- 2-Socket, single slot servers
- M8 with AMD 4th/5th gen EPYC **CPU**



* End of sale

4th,5th and 6th Gen FI

FABRIC

- 25/100G ports
- Unified ports: 32G FC 64XX, 6536). 64G FC (66xx)
- · Supports VIC 1400, 14000 and 15000 series



UCS X-Series Direct

· Scale at the edge with X-series advantage for 1-16 servers



25/100G IFM

• 8 x 25/100G connectivity



4th and 5th Gen VIC

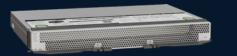
• 25/100G connectivity for both blades and racks



X-FABRIC AND PCIE NODE

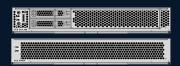
X-Fabric

- Based on native PCle Gen 4
- Provides GPU acceleration to enterprise application
- No backplane or cables = Easy upgrades



GPU Node and Front Mezz GPUs

- · Nvidia A16, Nvidia L40S, Nvidia L4, AMD MI210
- and Nvidia H100-NVL GPUs today in various configurations

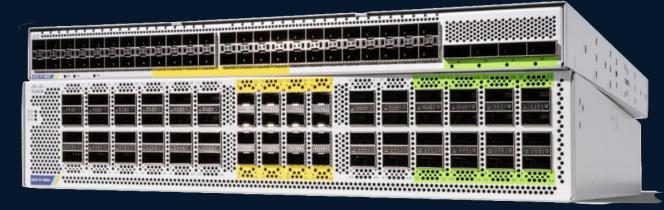


6th Generation UCS Fabric Interconnects

Redefining Compute Connectivity

- Simplified rack and modular server transition to 64G Fibre Channel
- Densest 25/50G or 100G converged fabric POD for x86 compute
- Comprehensive storage protocol support for compute
- Cloud scale management and visibility





Cisco UCS X210c M8

2S Compute Node for a wide range of workloads, including virtualization, web, collaboration, cloud, and bare-metal applications

Up to 172 Cores

2 x Intel® Xeon®6 Scalable Processors 6500P/6700P CPUs

Up to 8TB Memory

32x 6400 MT/s DDR5 Up to 256G Per DIMM Support for MRDIMMs (up to 8000MT/s)



Up to 9 Drives

Up to 9 E3.S drives, or 6 SAS/SATA/NVMe (tri-mode)

25G -100G

mLOM/Mezz and M.2 support

GPU Support with X440p & Front Mezz

Nvidia A16, Nvidia L40S, Nvidia L4, Nvidia H100 NVL

Cisco UCS X215c M8

2S Compute Node for a wide range of workloads, including virtualization, web, collaboration, cloud, and bare-metal applications (DB and analytics)

Up to 320 Cores

2 x AMD 4th/5th Gen EPYC CPUs

Up to 6TB Memory

24x 6000 MT/s DDR5 Up to 256G Per DIMM



Up to 6 Drives

Up to 6 SAS/SATA/NVMe (trimode) 25G -100G

mLOM/Mezz and M.2 support

GPU Support with X440p & Front Mezz*

Nvidia A16, Nvidia L40S, Nvidia L4*, Nvidia H100 NVL

UCS X-Series Direct: Expansion!

FCS RELEASE (Q3CY24)

Up to 8 compute nodes

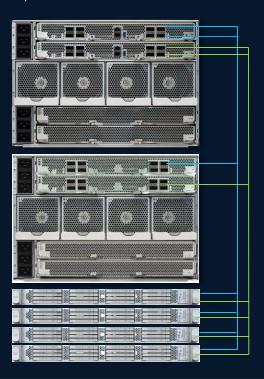
- × X210c M6, X210c M7, X215c M8
- X410c M7
- Support of X-Fabric + X440p
- 4th & 5th Gen VIC
- Full config support of X-series compute Node



Now Available

Up to 16 X-Series compute nodes

- IFM-100G for 2nd Chassis
- Plus: Up to 4 rack server



Al Specific Compute

Cisco Compute Al Portfolio

Address Al workloads with visibility, consistency, and control

Validated solutions for AI with compute, network, storage, and software

Build the model Training

Optimize the model Fine-tuning and RAG

Use the model Inferencing



UCS Dense GPU Servers

Dense compute for demanding Al



UCS Blade (w/GPU Expansion) and Rack

Full stack AI with compute and networking

Cisco UCS C845A Versatile Al Server

Optimized for GenAl

Address a variety of Al use cases

Modular NVIDIA MGX™ architecture
enables configuration flexibility

Built on air-cooled enterprise
rack design that easily fits in
your data center

Scalable Al performance

Choose the number of GPUs that match your use case

Increase the number of GPUs when your workloads demand it _____

Easily create clusters of Al servers to scale out and meet increasing needs

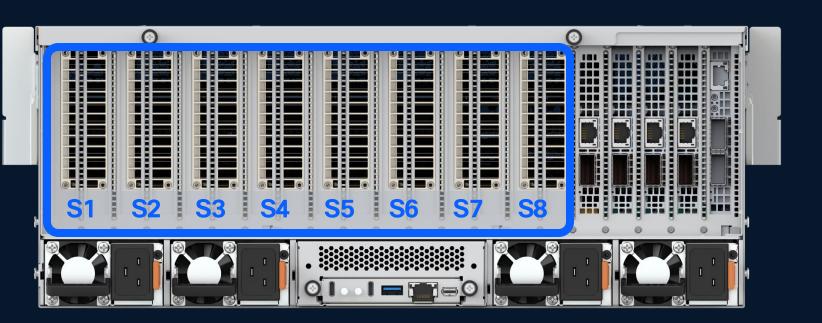
Consistent management

Cisco Intersight delivers a unified management paradigm across all your UCS servers

Manage your Al server with the same tool as your traditional servers



UCS C845A



2/4/6/8



H100 NVL, H200 NVL, RTX 6000 Pro, L40S GPUs



2/4/6/8

AMD

MI210 GPUs

300W/GPU

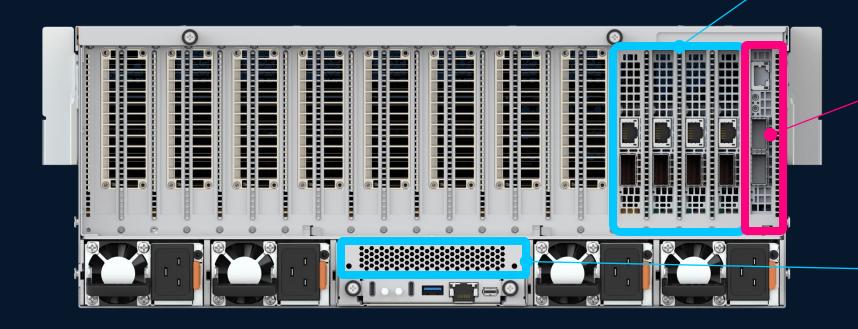


UCS C845A

4xNvidia CX-7 or

BlueField-3 3140H

for east-west GPU traffic



1x

Nvidia CX-7 or BlueField-3 3220

for north-south front-end traffic

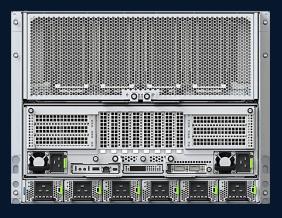
1x OCP 3.0

Intel X710

for north-south host management traffic

UCS C885A M8 Massive performance for Al at scale





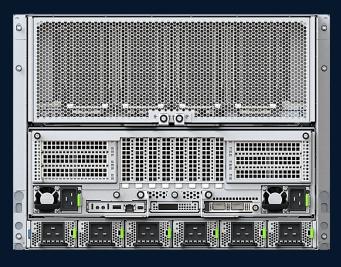
HGX 8RU 19" EIA Rack
8x Nvidia H100/H200/, AMD MI300X/MI350X GPUs
AMD Genoa 400W/Turin 500W TDP
8x PCle5 x16 HHHL + 5x PCle5 x16 FHHL NICs

TRAINING - LARGE / MEDIUM / SMALL MODELS FINE TUNING LARGE MODEL INFERENCING RAG

- Service Providers
- Financial Services
- Manufacturing

- Healthcare and Life Sciences

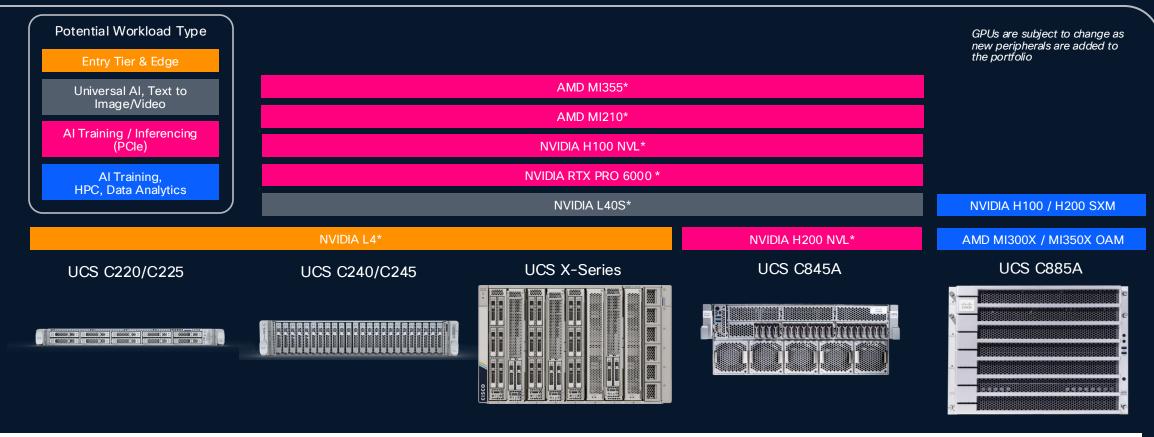
C885A M8 Specifications





Product Specifications						
Form Factor	HGX 8U 19" EIA Rack					
Compute + Memory	 2 AMD SP5 Genoa/Turin, 96 Core, TDP 400W, Up to 3.7 GHz 24 DDR5 RDIMMs Up to 6,000 MT/S 					
Storage	 1 PCle3 x4 M.2 NVMe (Boot Device) 16 PCle5 x4 2.5" U.2 NVMe SSD (Data Cache) 					
GPU	8x H100 or 8x H2008x MI300X or 8x MI350X					
Network Cards	 8 PCle5 x16 HHHL for E-W NIC ConnectX-7, BF3 B3140H 5 PCle5 x16 FHHL for N-S NIC BF3 B3220, B3240 2 OCP 3.0 SFF 					
Front IO	• 2 USB 2.0, 1 ID BTN, 1 Power Button					
Rear IO	 1 USB 3.0 A, 1 USB 3.0 C, mDP, 1 ID BTN, 1 Power Button, 1 USB 2.0 C (for debugging), 1 RJ45 (mgmt.) 					
Power Supply	 Up to 6 54V 3kW and 2 12V 2.4kW, N+1 redundancy 					

Al Platform Considerations: UCS GPU Options



Max GPUs 3 2-8* 2-8* 2-8 8

^{*} NOTE: GPU Form Factor and GPU model support may vary between AMD and Intel Platforms (i.e. c220/c225, c240/c245, and x210c/x215c). Check the spec sheet for each platform to determine maximum GPU support based on GPU selection

Al Infrastructure requirements by workload type

Workload type	Small* model inferencing / RAG	Large* model inferencing / RAG	Fine Tuning	Small* model training	Large* model training
Al Accelerator (GPU)	Maybe 0-2 GPU per server	Yes 2-4 GPU per server	Yes 4-8 GPU per server	Yes 2-4 GPU per server	Yes 8 GPU per server
Cluster size	0 - 100's Scales horizontally with request rate	0 - 100's Scales horizontally with request rate	20-80 GPUs ~10's of servers	20-80 GPUs ~10's of servers	>1000 GPUs 100s of servers
Intra-host GPU interconnect (e.g. NVLink)	No	Maybe**	Maybe	Yes	Yes
Inter-host GPU interconnect ("Backend" network)	No	Probably Not**	Maybe 200-400G RoCE	Yes 200-400G RoCE	Yes 200-400G RoCE/Infiniband
"Frontend" network	"Any"	100-400G	100-400G	100-400G	100-400G



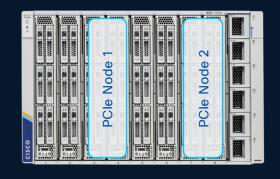
^{* &}quot;Small" and "Large" here are relative to the size of the GPU. i.e. Small can fit in memory of a single GPU. ** Depends on parallelism strategy, e.g. Tensor Parallel, Pipeline Parallel, etc.

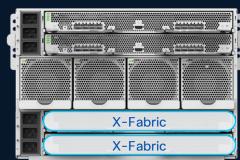
UCS-X GPU Expansion

2nd Gen X580p PCle Node and X9516 X-Fabric









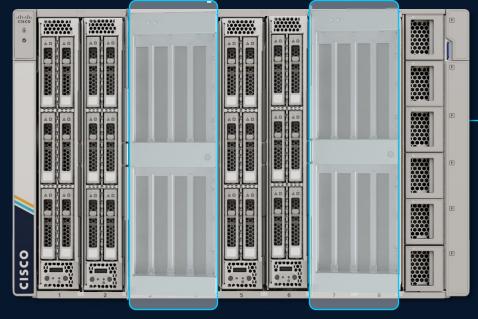
UCS X-Fabric Technology with PCle Node

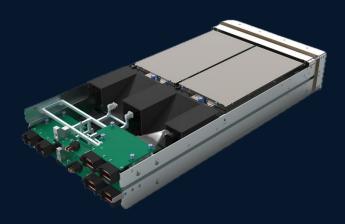
- \bigcirc
- PCle Switching with PCle Gen 5 connectivity
- 4x FHFL or HHHL GPUs per PCle node
- Intra-host GPU interconnect with NVLink
- Intersight policy-based Management



Inter-host scaling with RDMA over Al Fabric

UCS X580p PCle Node

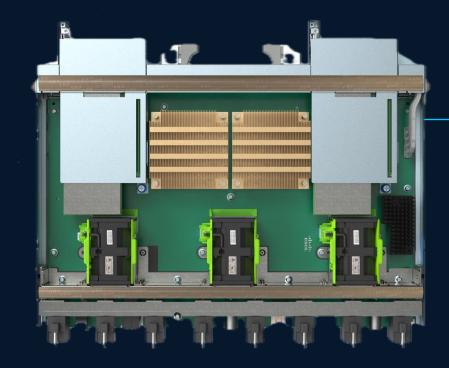




- Double wide PCle node for 4x FHFL GPU and PCle G5 GPU support
 - Nvidia H200-NVL, RTX PRO 6000 & L40S
- Support multiple vendors: Nvidia, AMD*/Intel*
- NVLink bridge support
- Support up to 600W FHFL GPU
- Managed PCle node with BMC support
- Policy based GPU management
- Ability to share GPUs across two Compute nodes

* AMD & Intel GPUs support will be post FCS

UCS X9516 X-Fabric





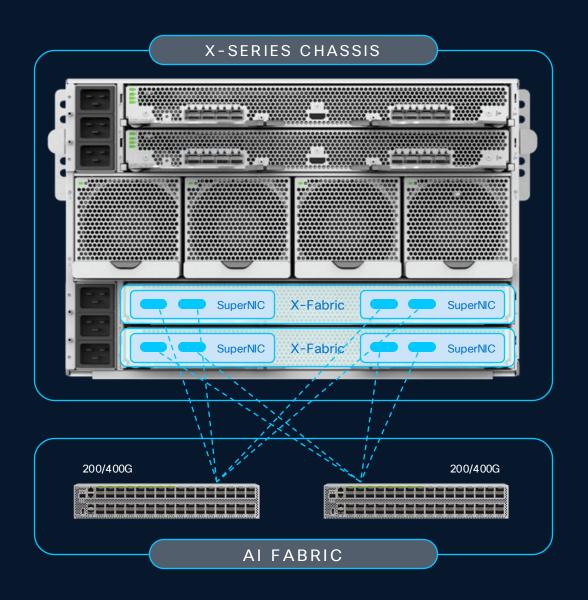
- PCle Gen5 Switching
- 2x CEM Slots to support HHHL NIC cards
 - ConnectX 7 (2x 200GB & 1x 400G)
- Managed XFM Modules with BMC support
- GPU Direct Support over RDMA
- GPU Backend (East-West Traffic) network support

Al Cluster Expansion

GPU-to-GPU connectivity

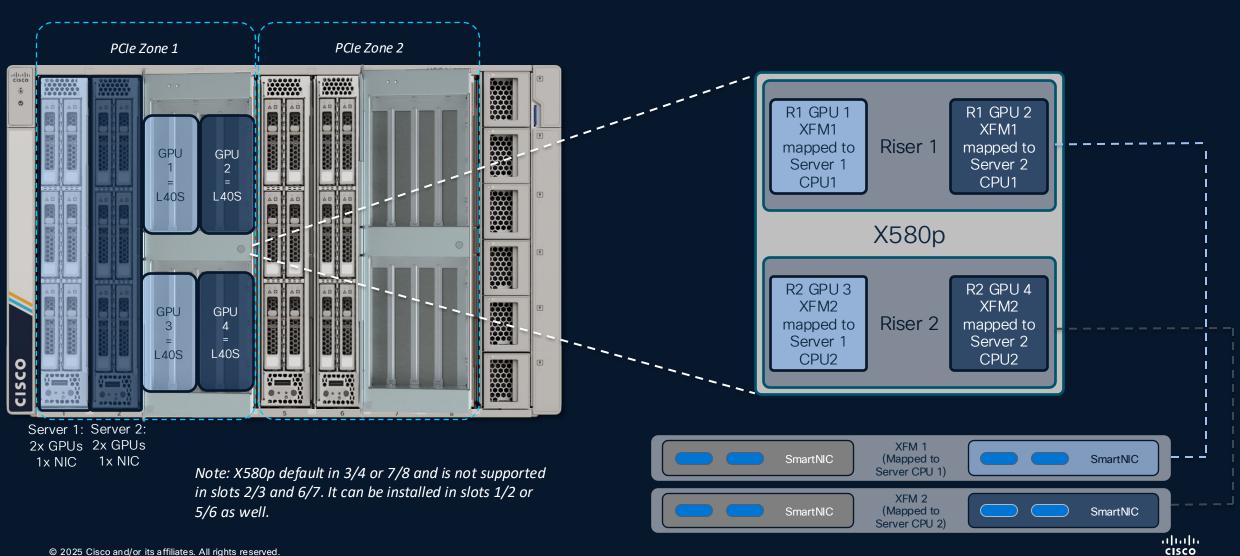
with XFM external ports





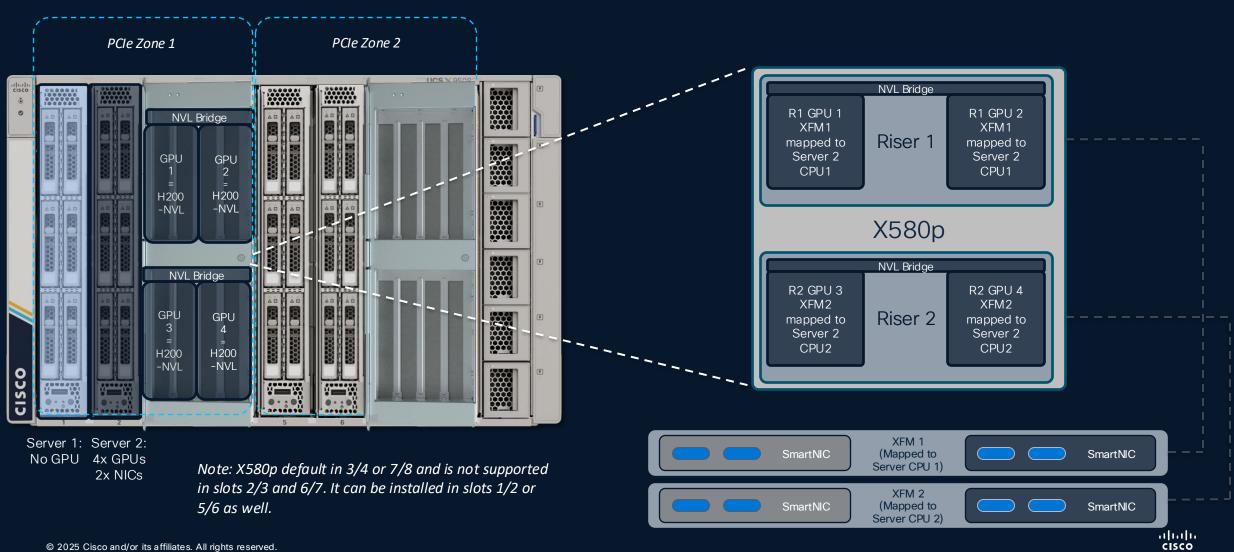
X580p - GPUs 1/3 Mapped to Server 1 and GPUs 2/4 Mapped to Server 2

(1x NIC mapped to each server)



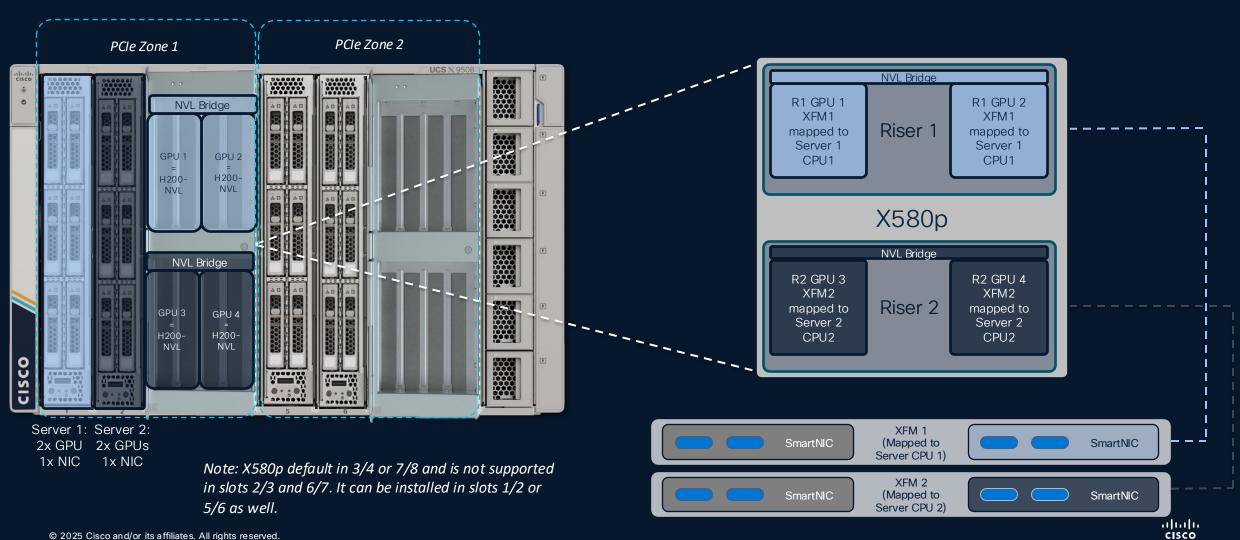
X580p - All GPUs Allocated to Server 2 w/NVL Bridge

(2x NICs mapped to one server)



X580p - 2x GPUs Allocated to Server 1 and 2x 2 w/NVL Bridge

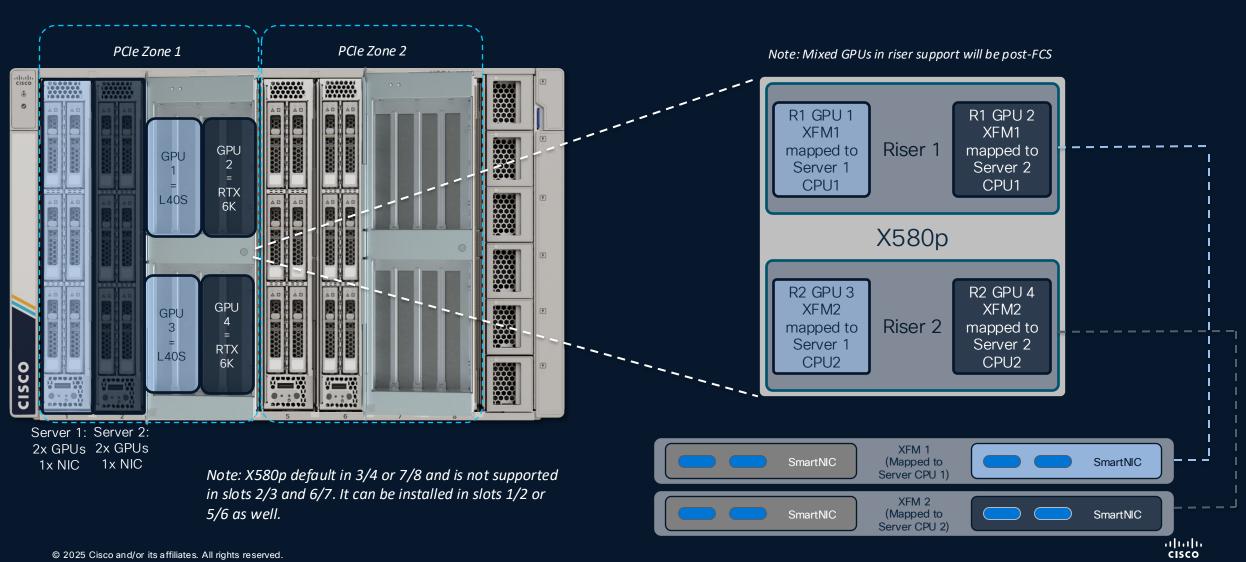
(1x NIC mapped to each server)



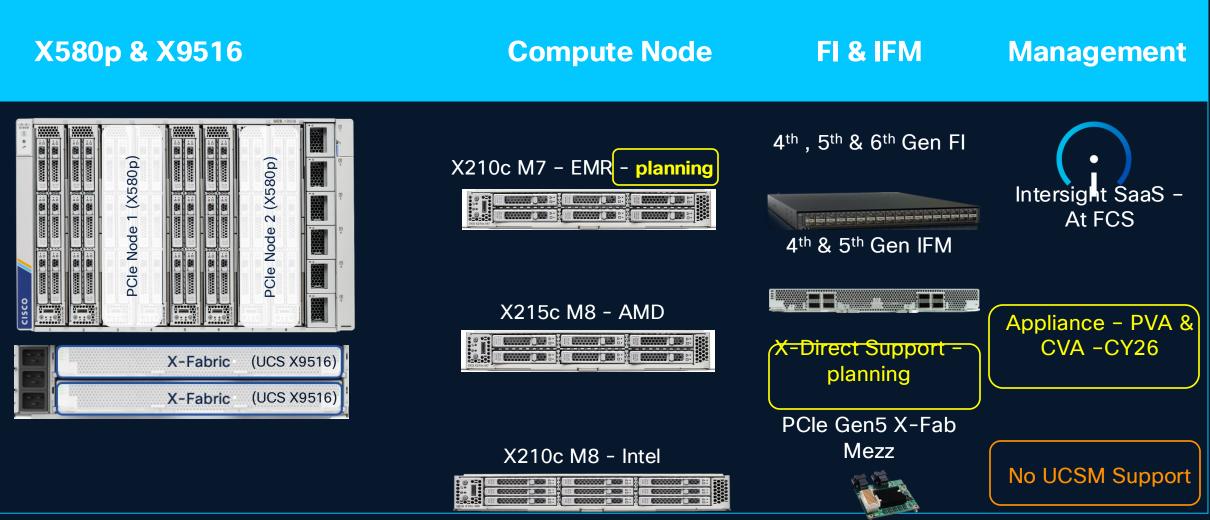
© 2025 Cisco and/or its affiliates. All rights reserved Cleveland Engage 10/2/2025

X580p - GPUs 1/3 Mapped to Server 1 and GPUs 2/4 Mapped to Server 2

(1x NIC mapped to each server)



Platform Support



Scalable Management

for the Al ready data center

Cisco Intersight

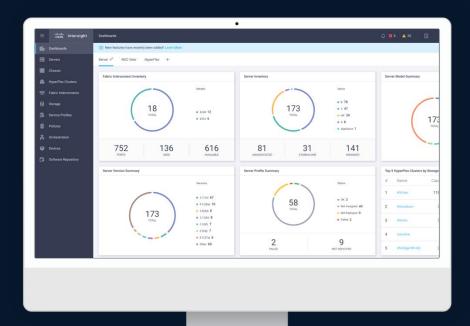
Unified Operating Model

Consistent operational model globally, from DC to edge, at cloud scale

Secure operations with built-in advisories and continuous risk mitigation

Simplified operation with Aldriven capabilities including Connected TAC, Proactive RMA, and Predictive Insight

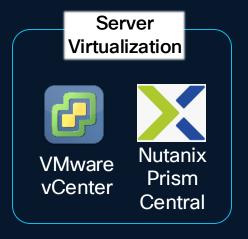
Automate deployments, configuration, workflows, and day-0 to day-N tasks



Supported Solutions













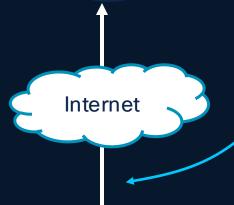
SaaS <u>US</u> or <u>EMEA</u> Supported Systems Appliance <u>US</u> or <u>EMEA</u> Supported Systems

Connecting Cisco solutions to Intersight (SaaS example)

Device Connector is built-in:

- Embedded in target management plane
- Delivered with management plane firmware





Device Connector

Native API

Cisco Target

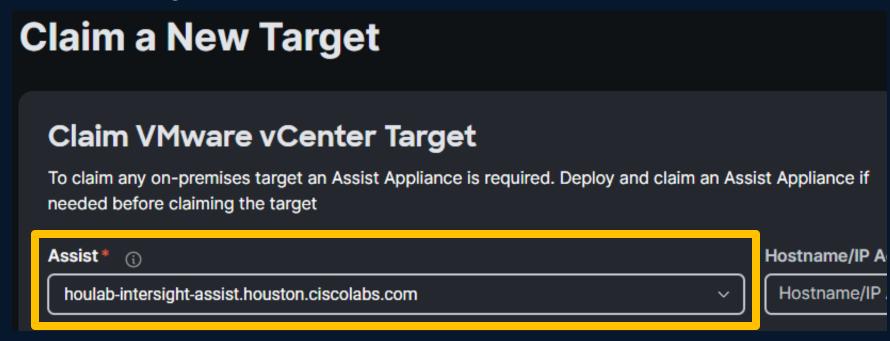
Device Connector connection:

- Highly secure
- Outbound single-destination
 HTTPS 443
- Durable websocket enabling bidirectional communication

*Optional proxy config available

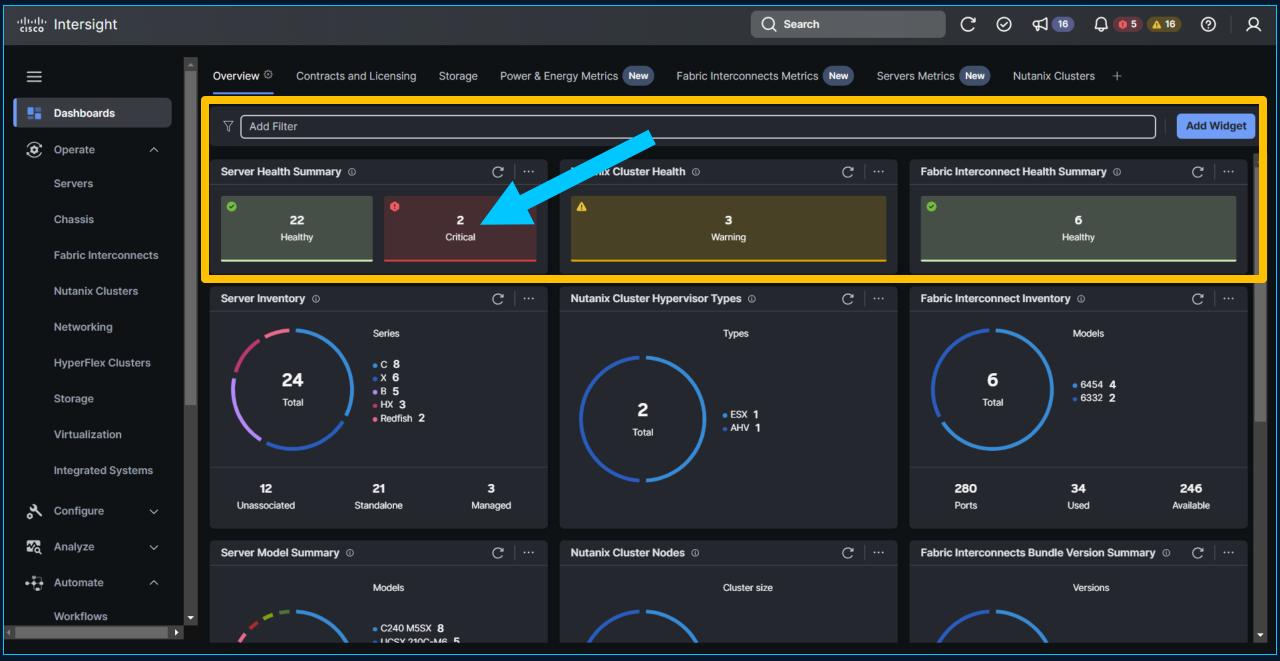
Connecting 3rd party solutions to Intersight

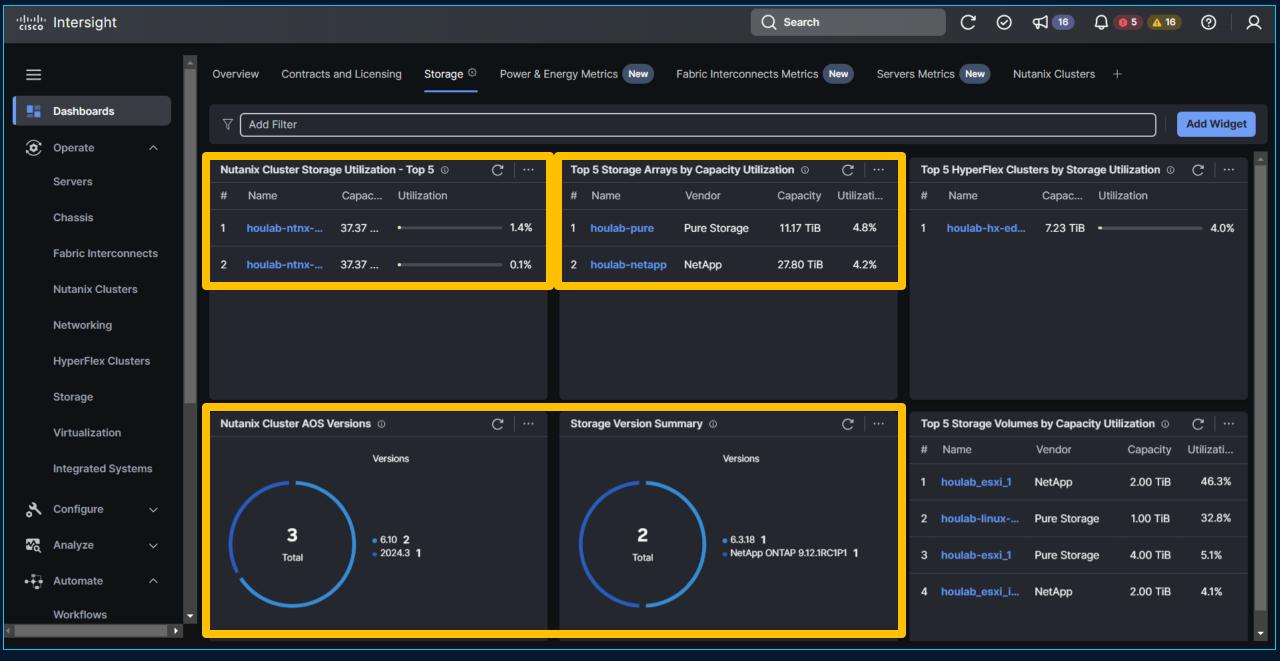
- Assist via mode selection of Intersight Appliance OVA
- Manages connections to 3rd party solutions via APIs through this device



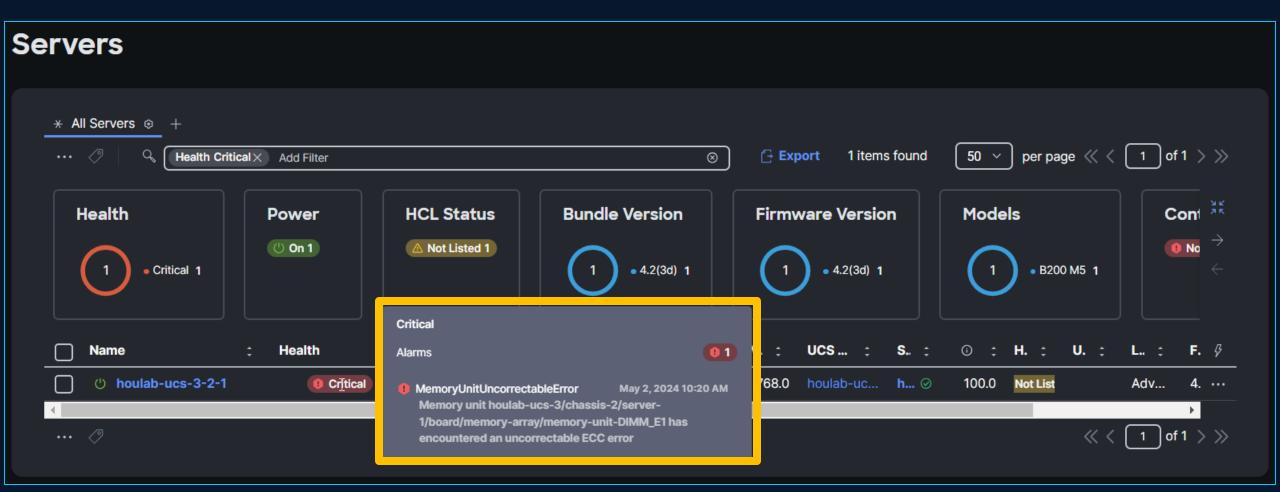
*Assist functionality is also part of the Connected and Private appliance model

Visibility



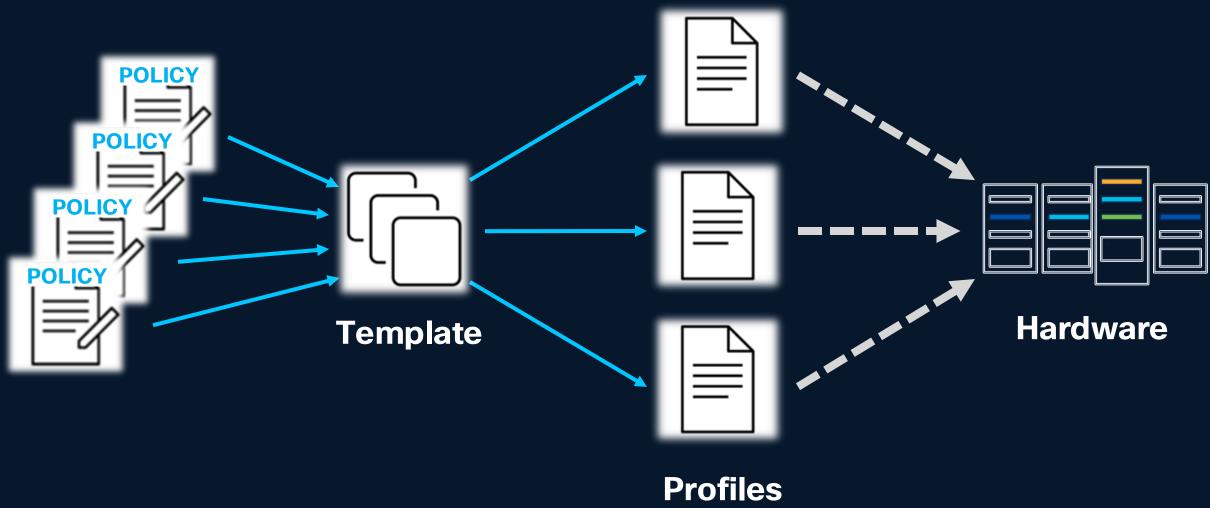


Proactive RMA - example

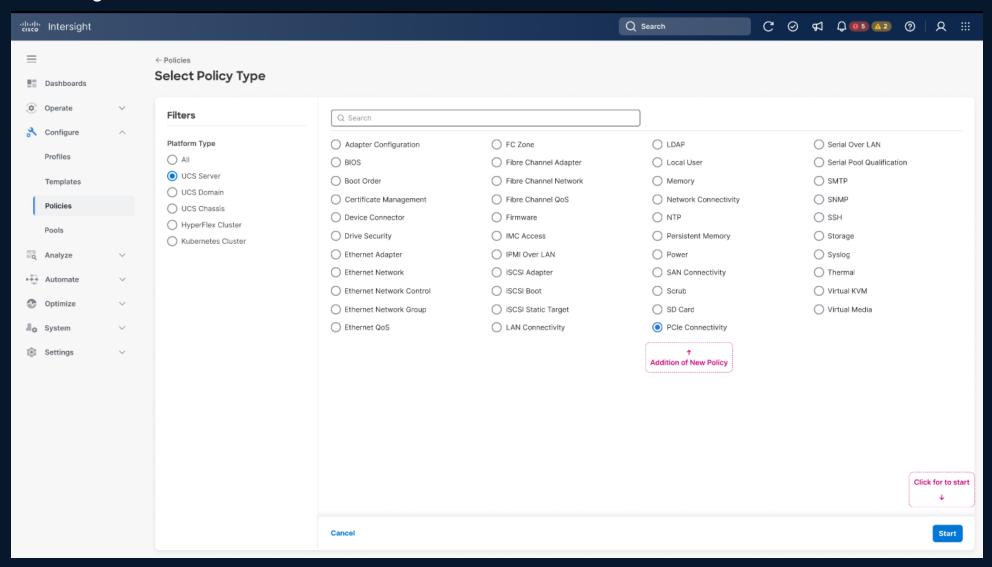


Policy Based Configuration

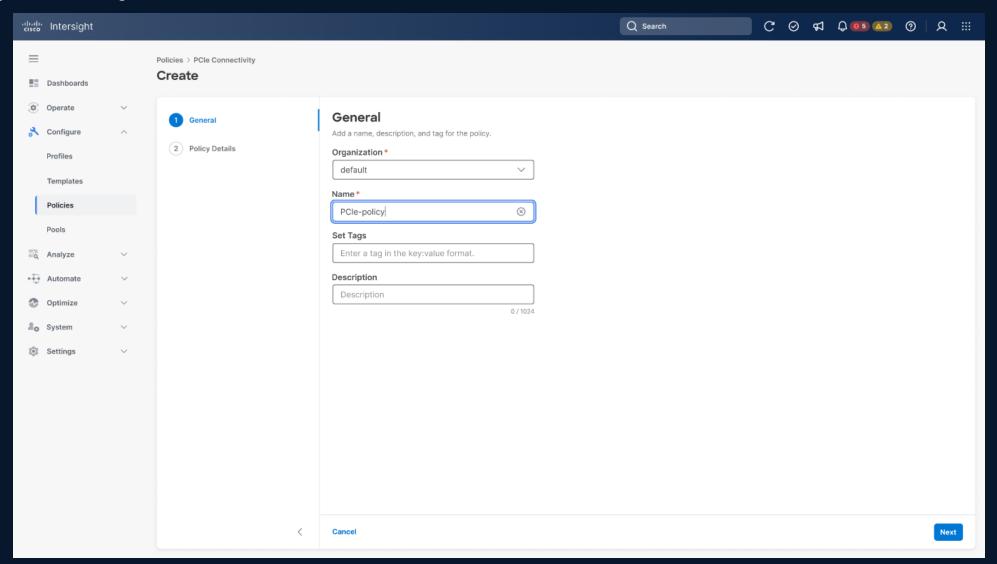
Policy for consistent configuration



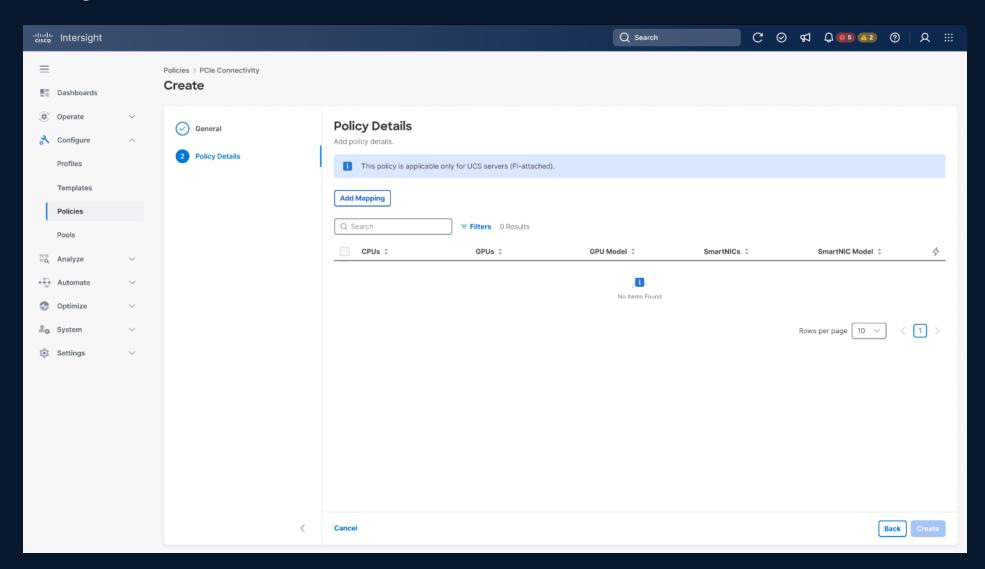
Intersight PCIe Connectivity Policy Creation

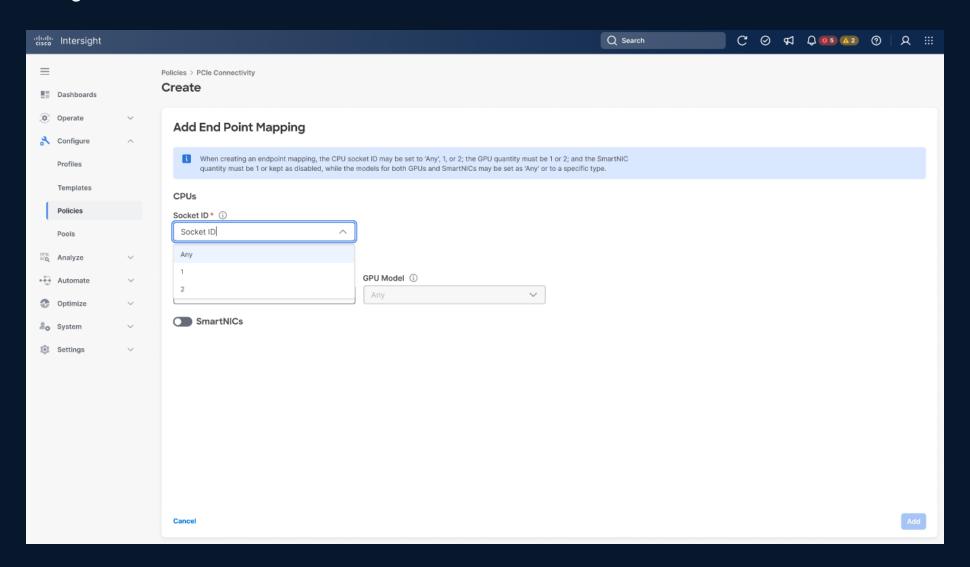


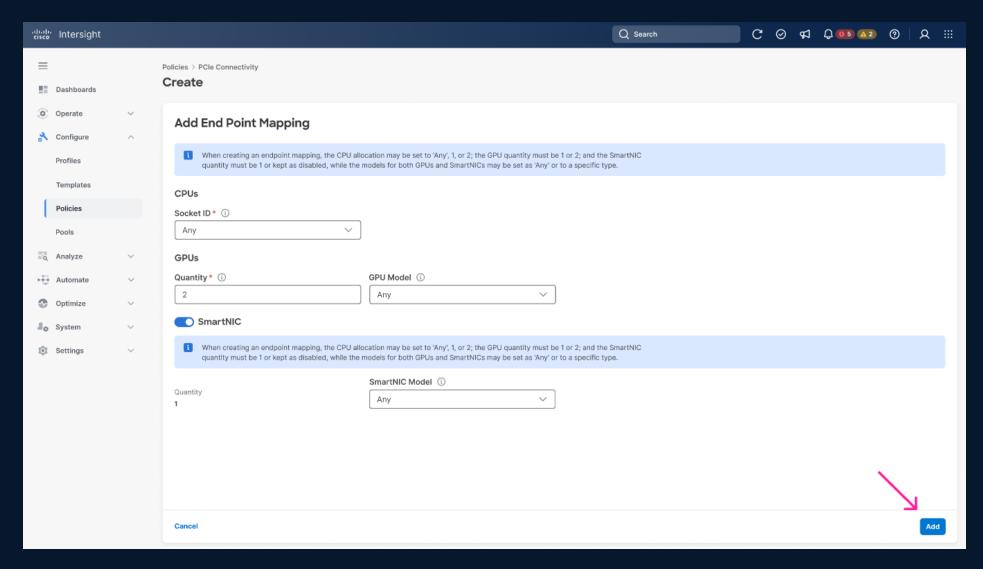
Intersight PCIe Policy Details

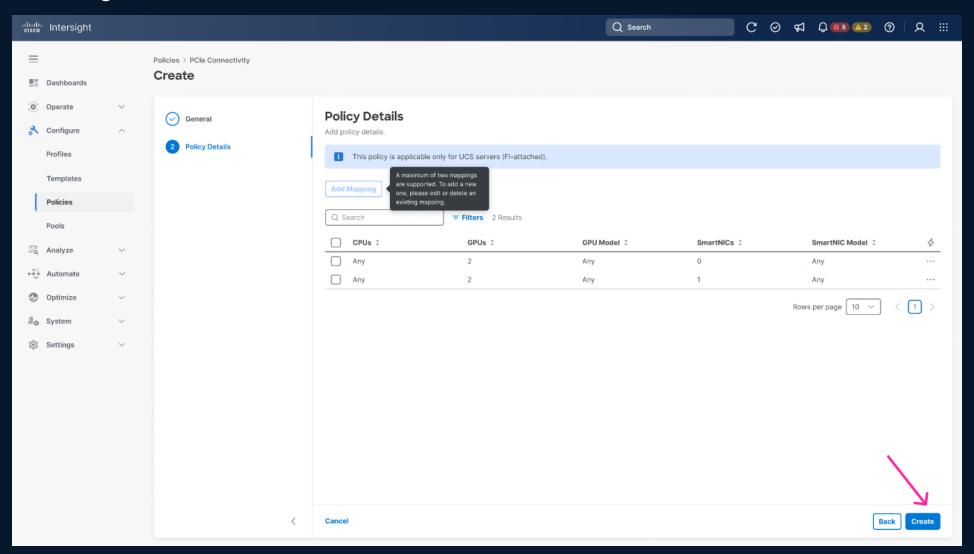


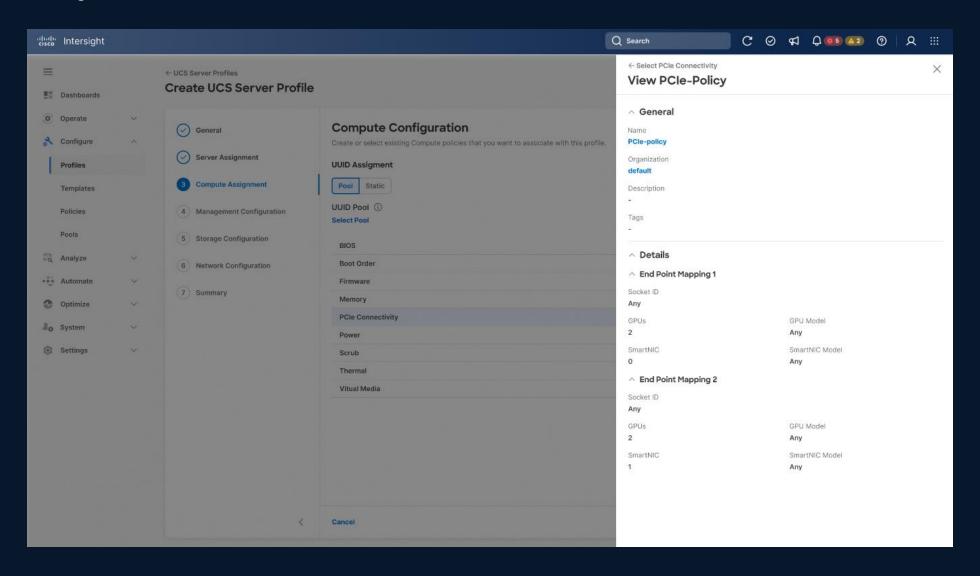
Intersight PCIe Policy Details





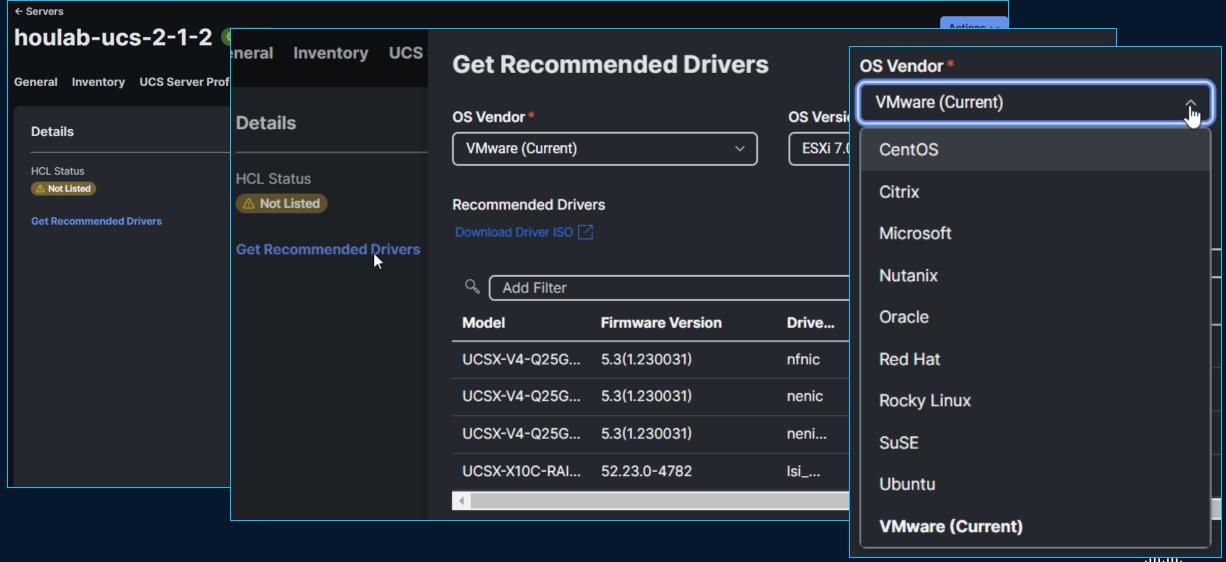


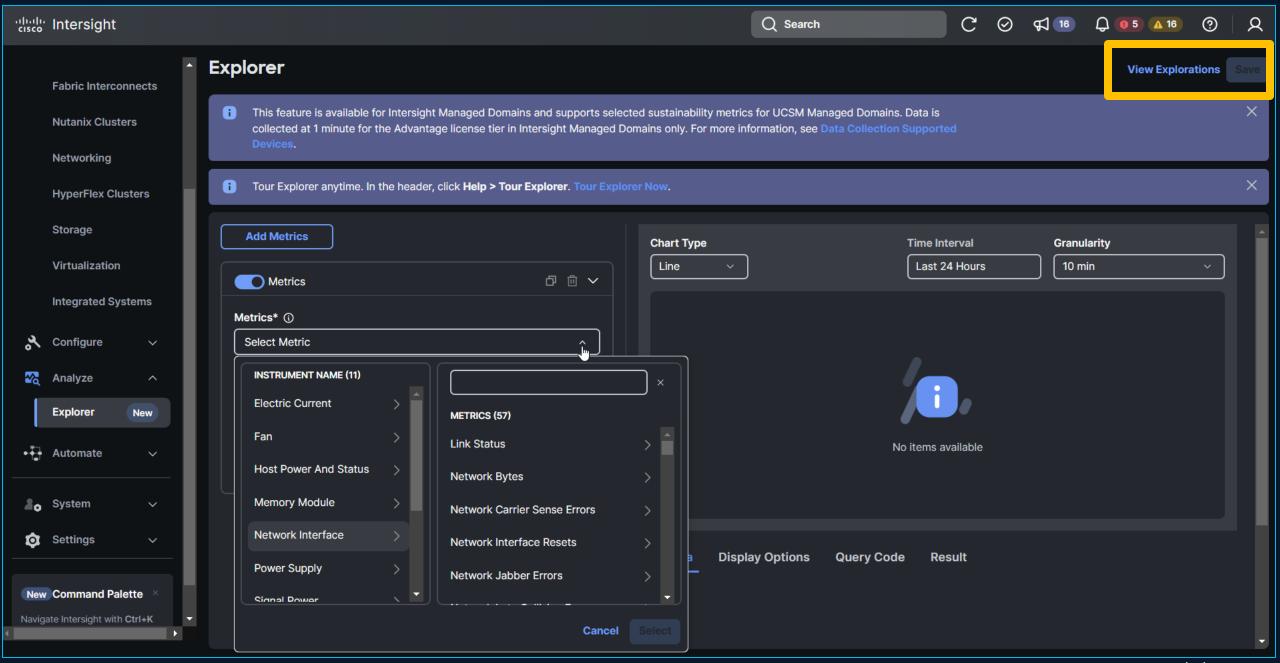


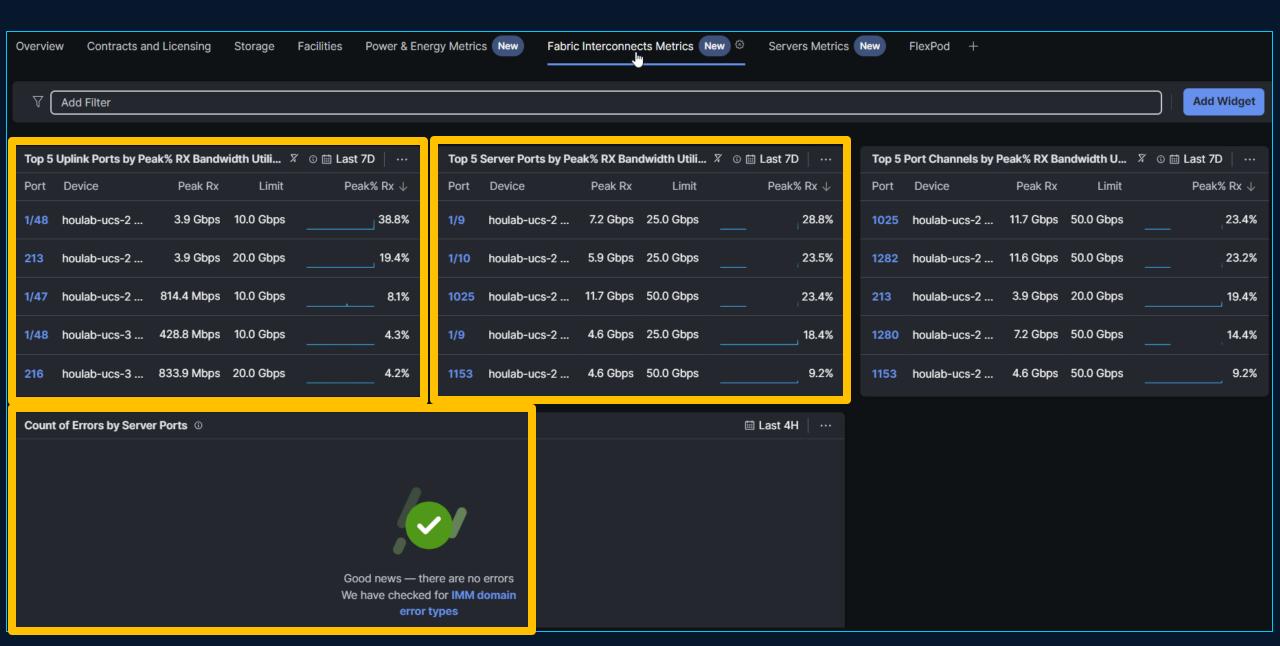


Day 2 Monitoring

HCL (Hardware Compatibility List)







Closing

- UCS offers the flexible, policy-based compute for any Al deployment
- Intersight provides a single place to configure and mange all Cisco Compute
- 3. Intersight integrates with 3rd party platforms to give a complete view of the DC
- 4. Intersight offers expanded configuration capabilities through its automation capabilities

Thank you



.1|1.1|1. CISCO