

Simplifying Compute Infrastructure for the AI Ready Data Center

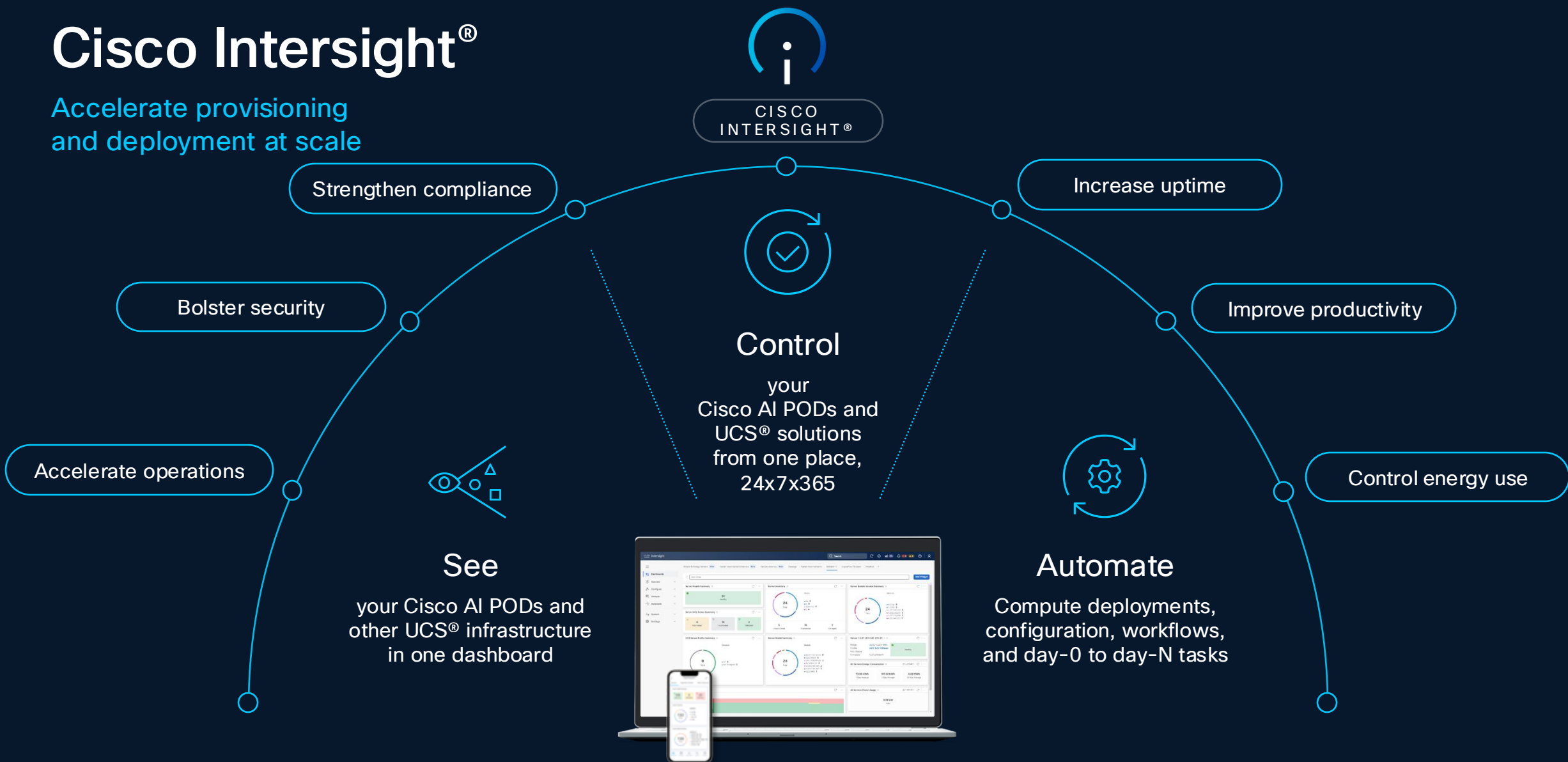
Terence Gibson
C&AI Solutions Engineer – US Commercial

John Rice
C&AI Solutions Engineer – US Public Sector



Cisco Intersight®

Accelerate provisioning
and deployment at scale



Intersight fleet management



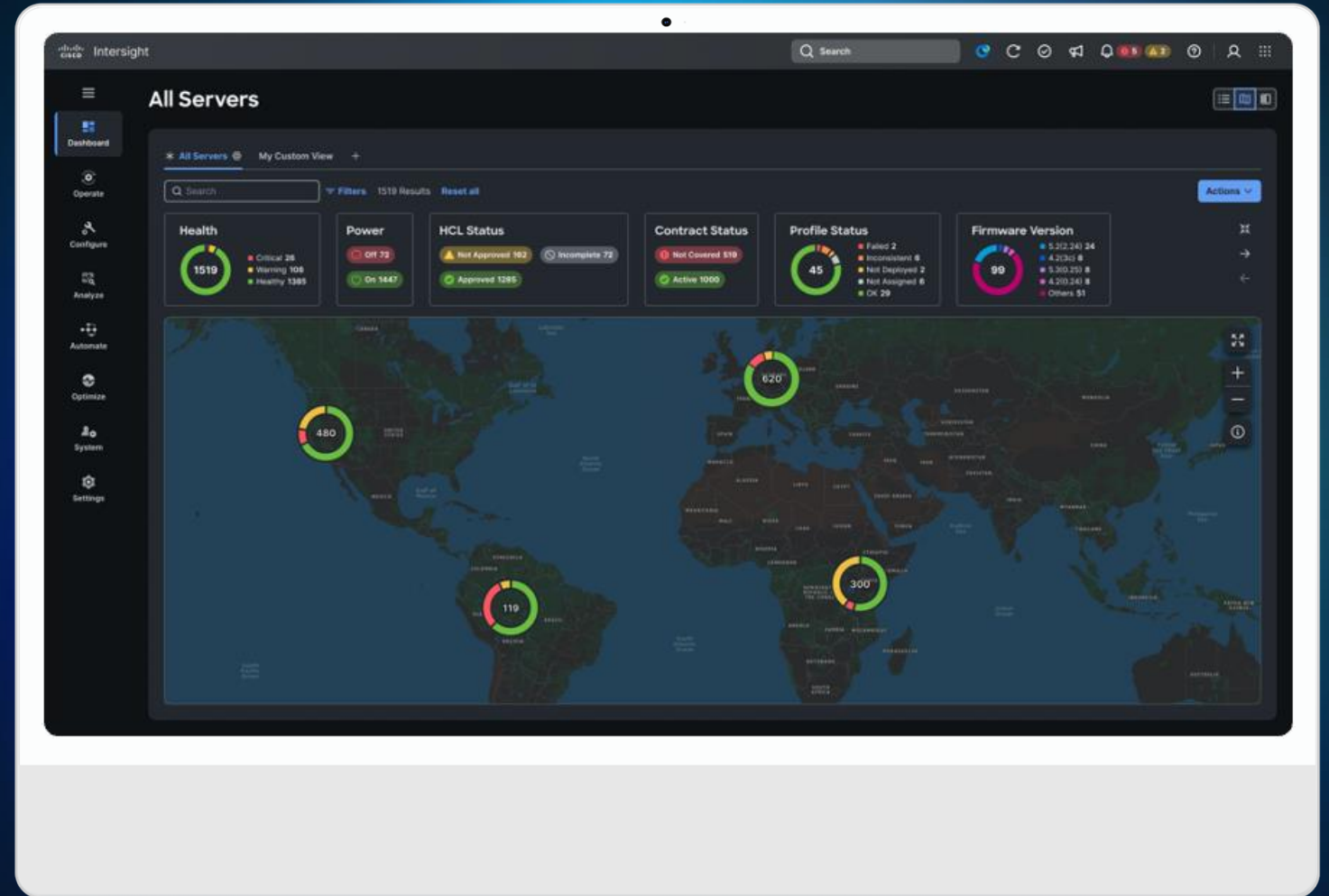
Simplified onboarding & Zero-touch provisioning



Automated lifecycle management



Global visualization



Compute AI portfolio

Address AI 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

RTX PRO SERVER

Supporting RTX PRO 6000 Blackwell Server Edition GPUs



Cisco UCS®
GPU-dense servers
PCIe and NVLink Servers



Cisco UCS blade (with GPU extensions) and
rack servers



Enterprise AI edge

Dense compute for demanding AI

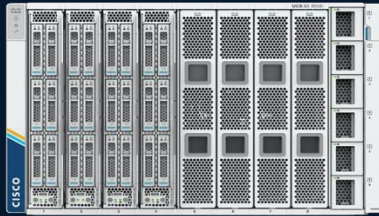
Full-stack AI with compute and networking

Cisco UCS Compute Portfolio

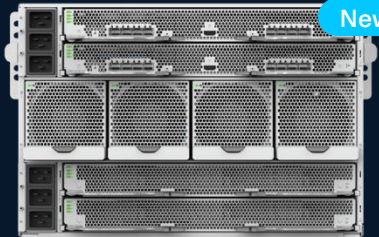
MAINSTREAM ENTERPRISE SERVERS

UCS X-Series
X9508 Chassis

IFM Module



UCS X-Series Direct



UCS x580p M8



UCS X210c M7



UCS X210c M8



UCS X215c M8



UCS X410c M7



UCS C240 M8E3S
36 EDSFF E3.S1T



New

UCS C240 M8SX
28 HDD/SDD/NVMe



New

RTXPRO

UCS C240 M8L
16 LFF + 4 SFF



New

UCS C240 M7SN
28 NVMe



UCS C240 M6S
14 SSD/HDD Media drive



UCS C240 M6N
14 NVMe Media Drive



UCS C220 M8E3S
16 EDSFF E3.S1T



New

UCS C220 M8S
10 HDD/SSD/NVMe



New

UCS C220 M7N
10 NVMe



UCS C245 M8SX
28 HDD/SDD



New

UCS C225 M8S
10 HDD/SSD



New

UCS C225 M8N
10 NVMe



New

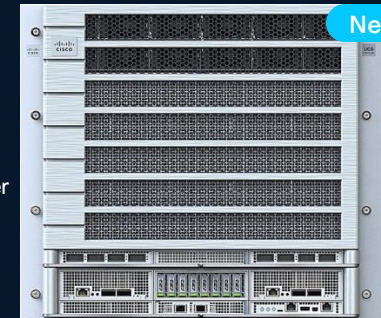
AI SERVERS

UCS C885A M8
8RU Dense GPU Server



New

UCS C880A M8
10RU Dense GPU Server



New

UCS C845A M8
4RU MGX Server



New

RTXPRO



Consolidate rack workloads



AI/ML

Accelerated VDI

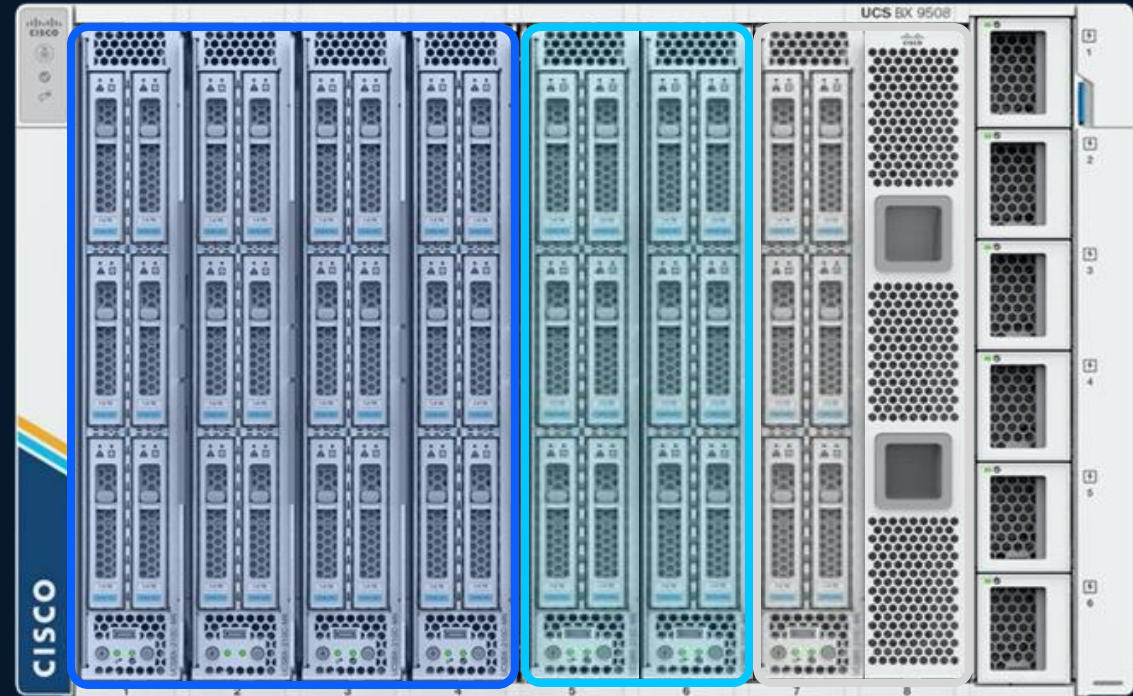


Big Data, SDS, Containers



Traditional blade workloads

UCS® X-Series with X-Fabric



Up to 2,048

cores
per chassis

24

GPUs
per chassis



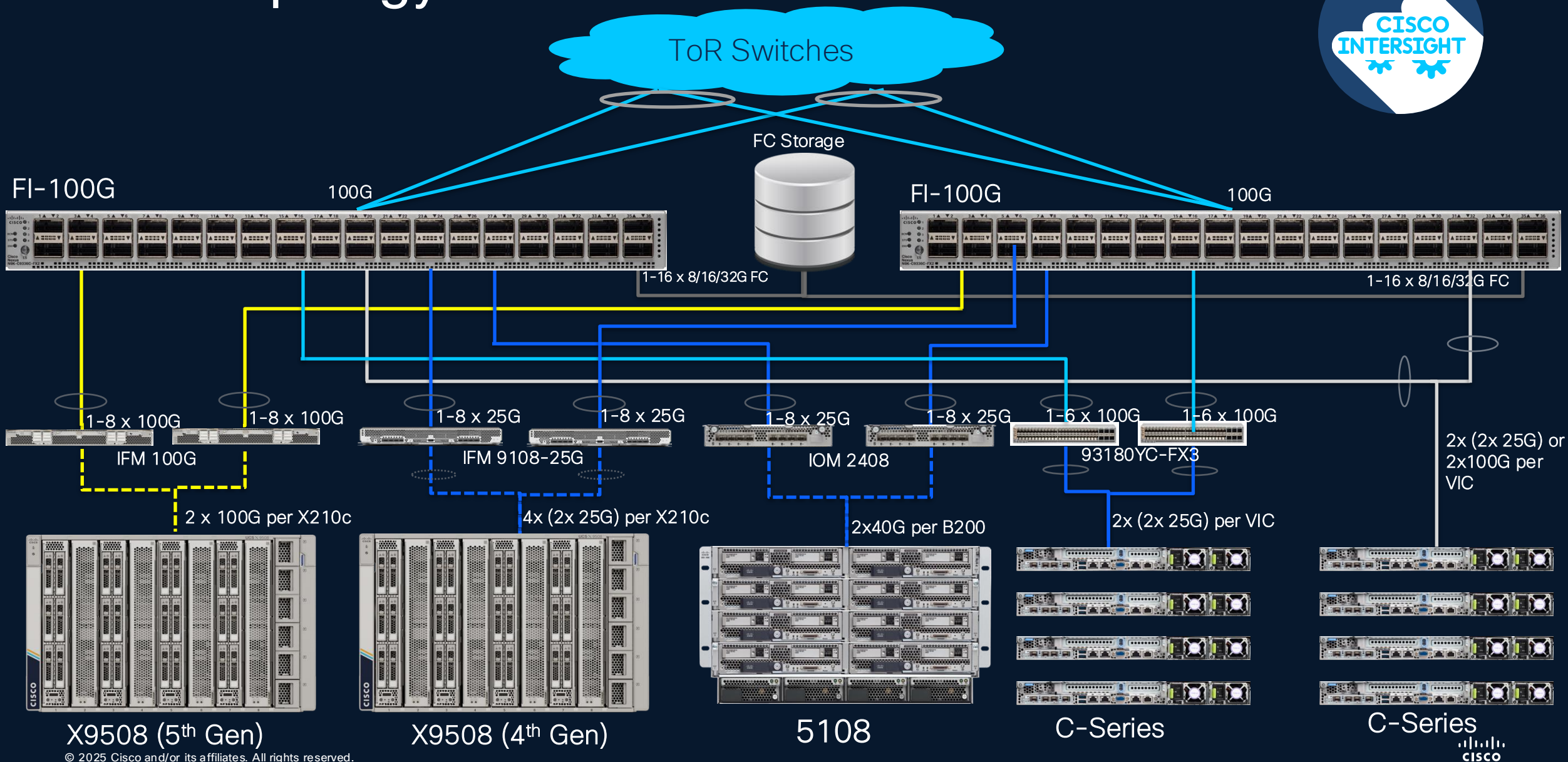
200G

bandwidth to
compute node

736 TB

of storage

Fabric Topology



X9508 (5th Gen)

X9508 (4th Gen)

5108

C-Series

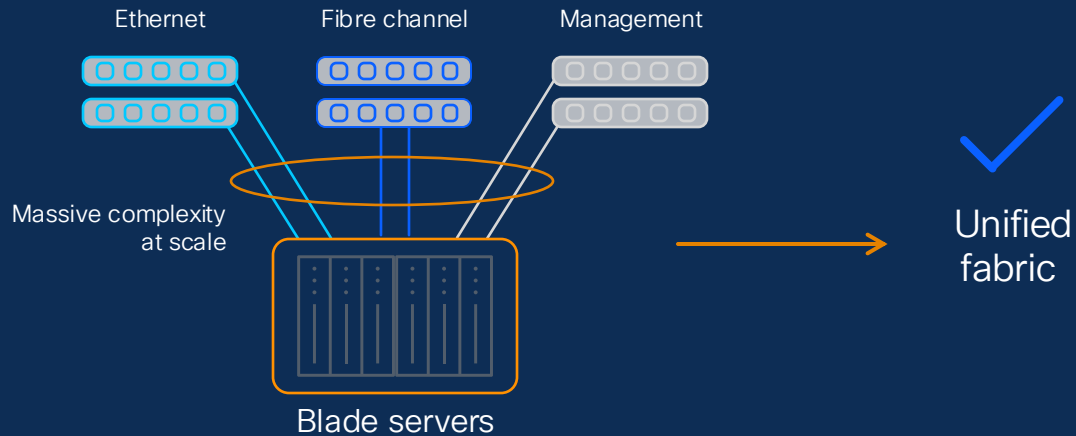
C-Series



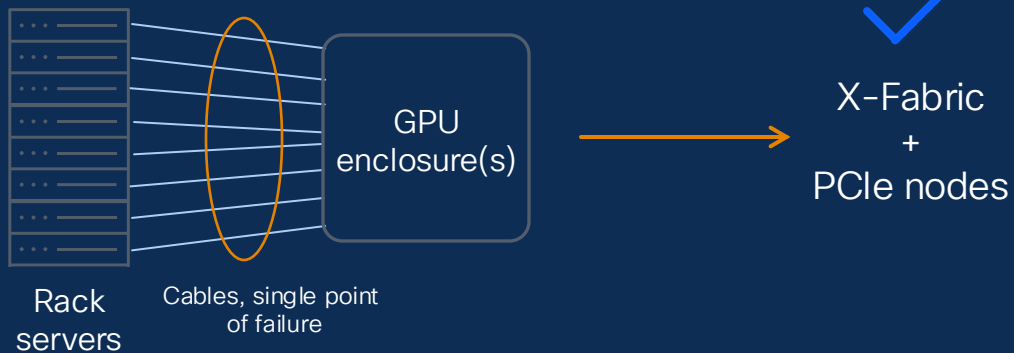
Industry-leading simplicity

Conventional approaches

1 | Silos of multiple Ethernet and SAN fabrics and adapters



2 | Complex PCIe connectivity to external accelerators

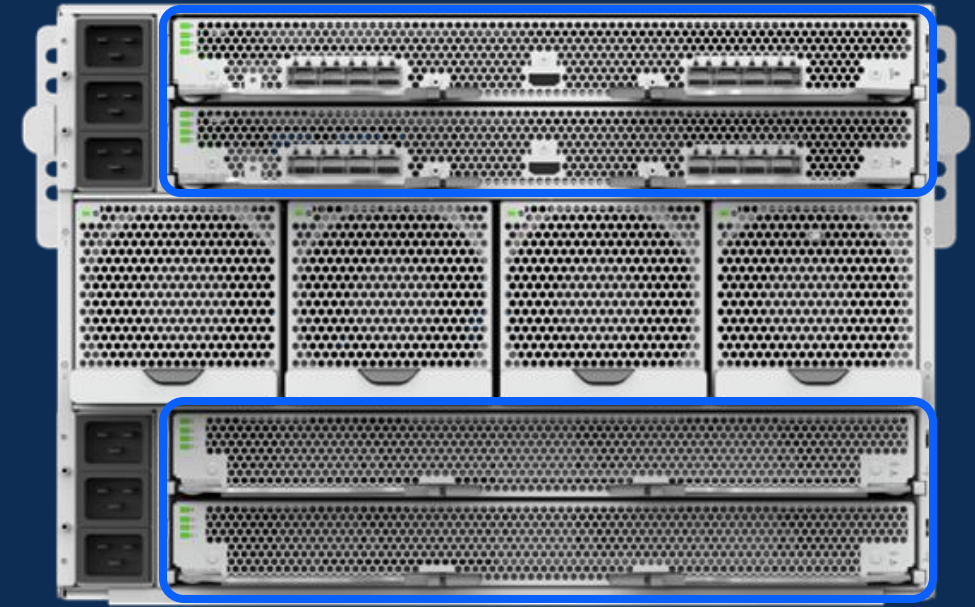


Cisco solution

UCS® X-Series



Cisco Intersight®



UCS X-Fabric Technology

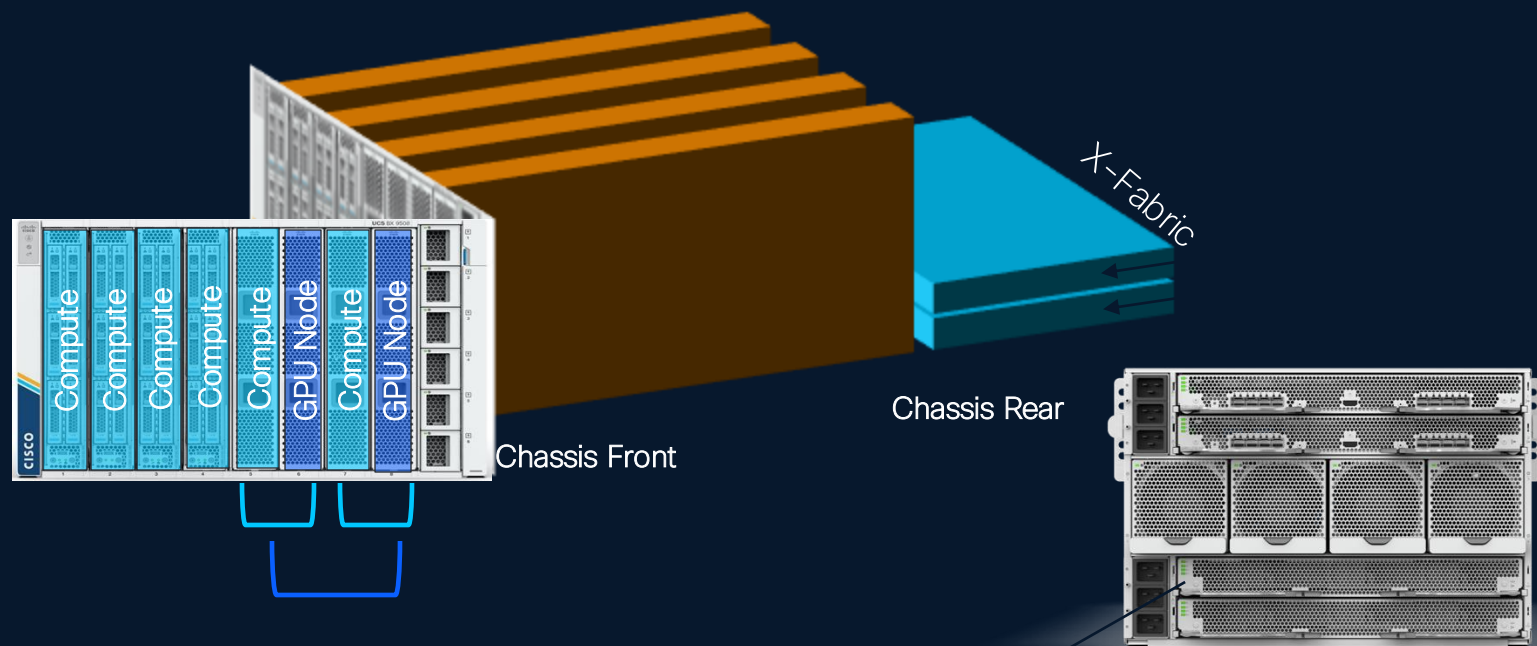
Open, modular design enables compute and accelerator node connectivity

Open standards: PCIe
4/5/6, CXL*

No midplane nor cables =
easy upgrades

Expandability to address
new use cases in future
(memory & storage nodes)

*CXL fabrics are dependent on future processors



UCS X-Fabric Technology

- Internal Fabric interconnects nodes
- Industry standard PCIe, CXL Traffic
- Upgrade to future generations

Common AI Challenges



Unclear business objectives & priorities

Unclear direction hinders cross team collaboration, creates confusion, and hampers acquisition of necessary skills



Complex AI infrastructure deployment

Lack of high-performance infrastructure with integrated compute, network, storage, and AI software can stall AI projects



Security vulnerabilities

AI models, frameworks, apps, and supporting infrastructure represent a new cyberattack surface



Network performance & Security challenges

Model training and inferencing generates a lot of traffic, slowing networks and also results in new attack surface

Cisco Differentiation



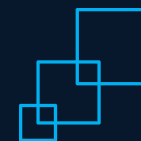
The Security

Security-first architecture enables safe enterprise AI



The Network

High-performance integrated AI networking enables efficient model training and inferencing



The Assurance

Pre-validated AI infrastructure stack with flexible deployment options improves data scientists and developer productivity

Introducing: Cisco AI PODs

A scalable architecture, built to support any AI workload simply & efficiently



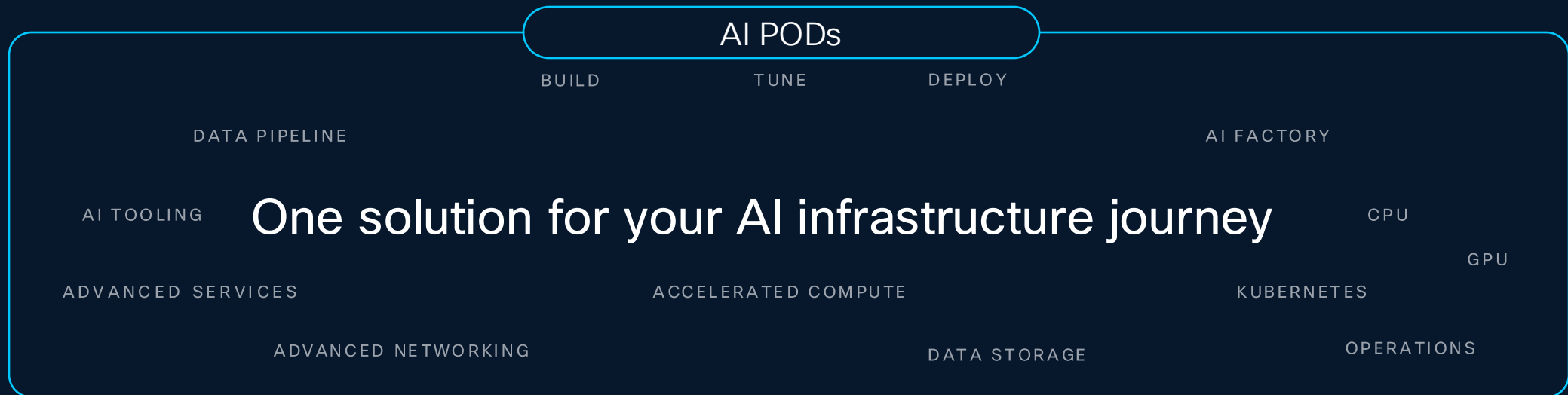
Training



Optimization



Inference



Cisco AI PODs

A scalable architecture, built to support any AI workload simply & efficiently

Deploy AI with confidence
Cisco CVD, NVIDIA ERA

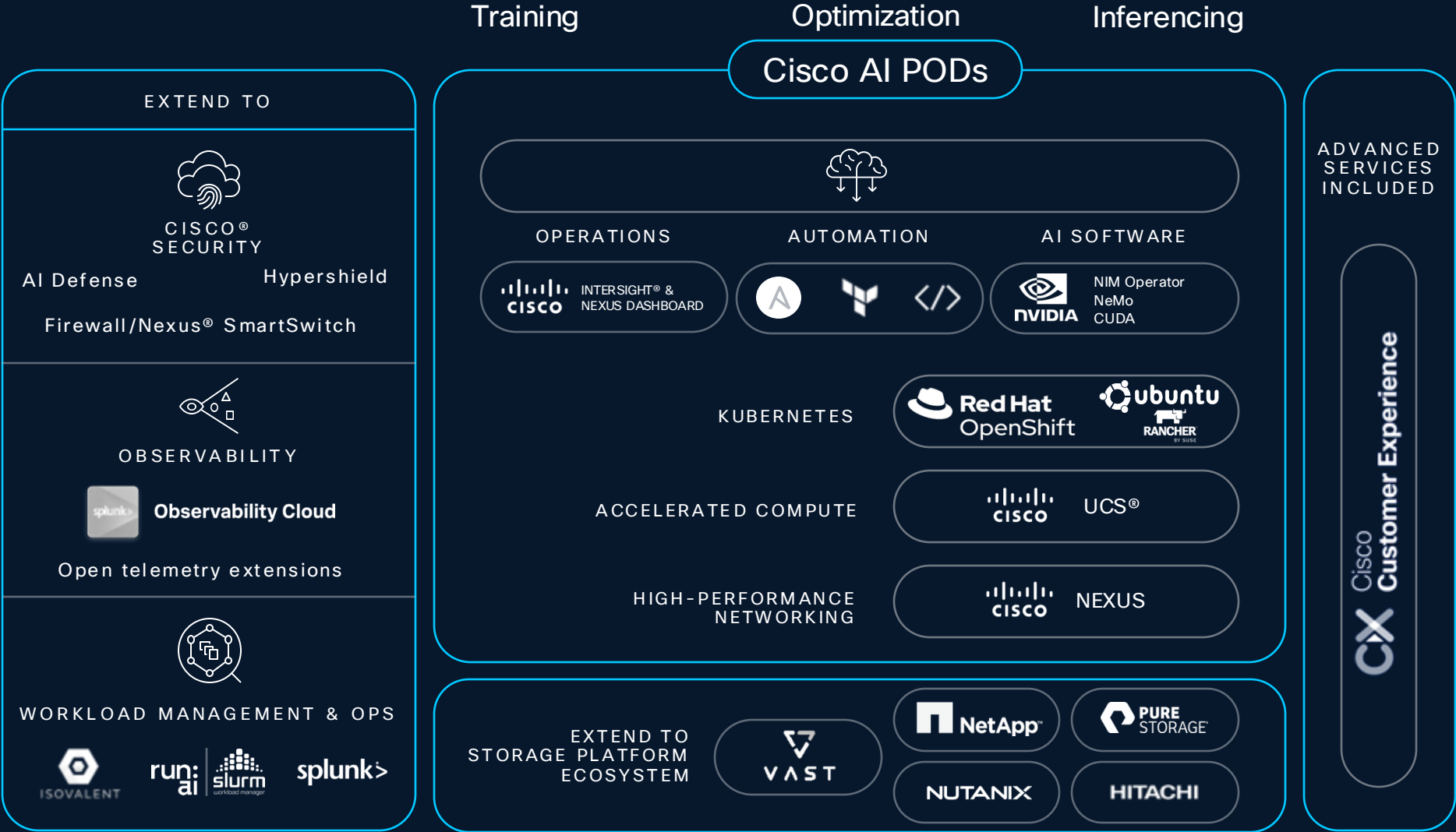
Fully supported stack
including Cisco and 3rd
party components

**Cisco CX
Success Track**

Orderable, use case
driven AI-ready
infrastructure stacks

**Inferencing.
Optimization.
Training.**

Incremental, atomic-level
-or- fabric-based
cluster scale



Cisco AI PODs

Expanded design portfolio

Full AI lifecycle use-cases support

Based on Cisco Validated Designs

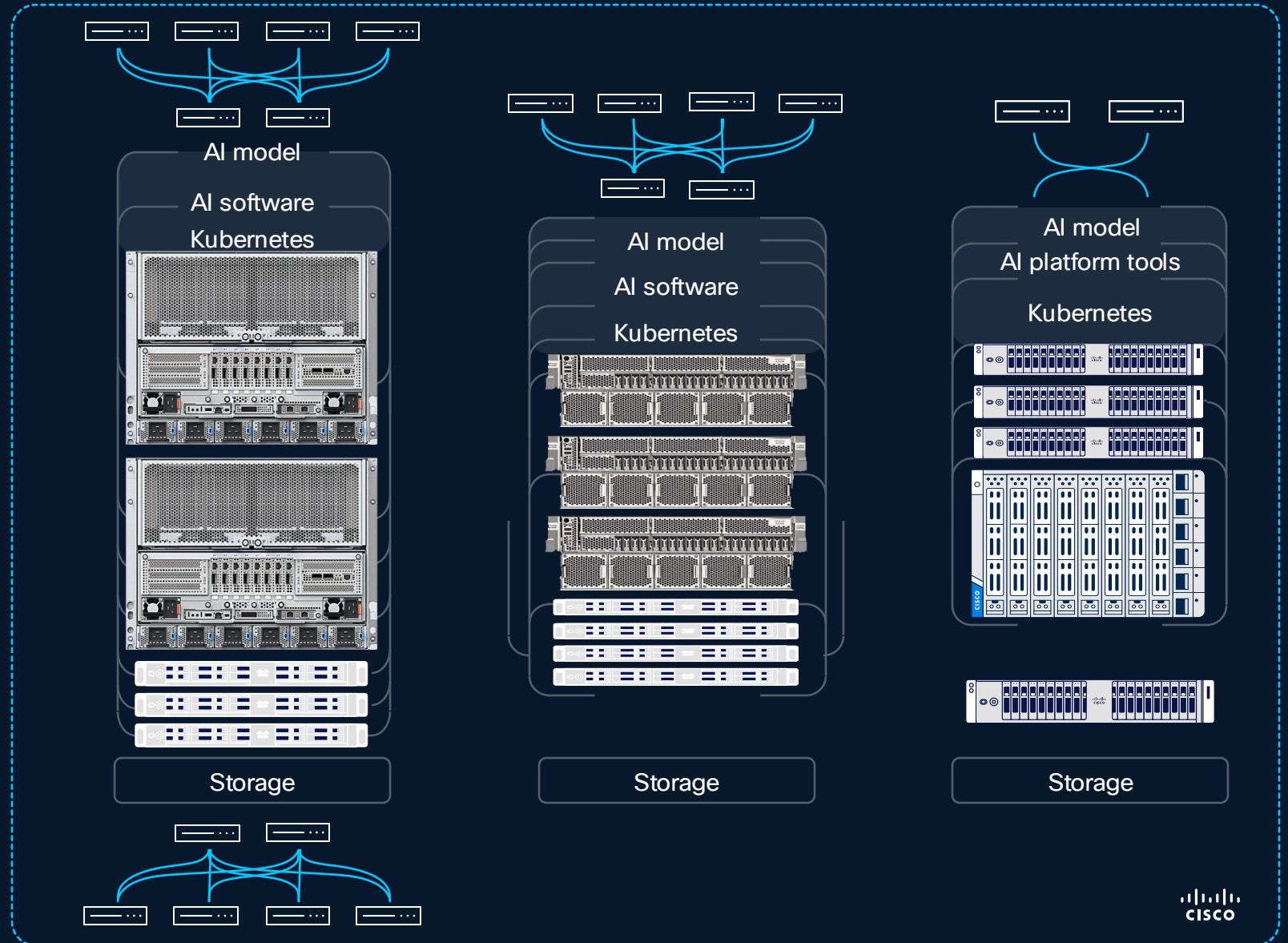
Pre-packaged AI-stack hardware +
software platforms & automation

Latest NVIDIA and AMD GPU Compute

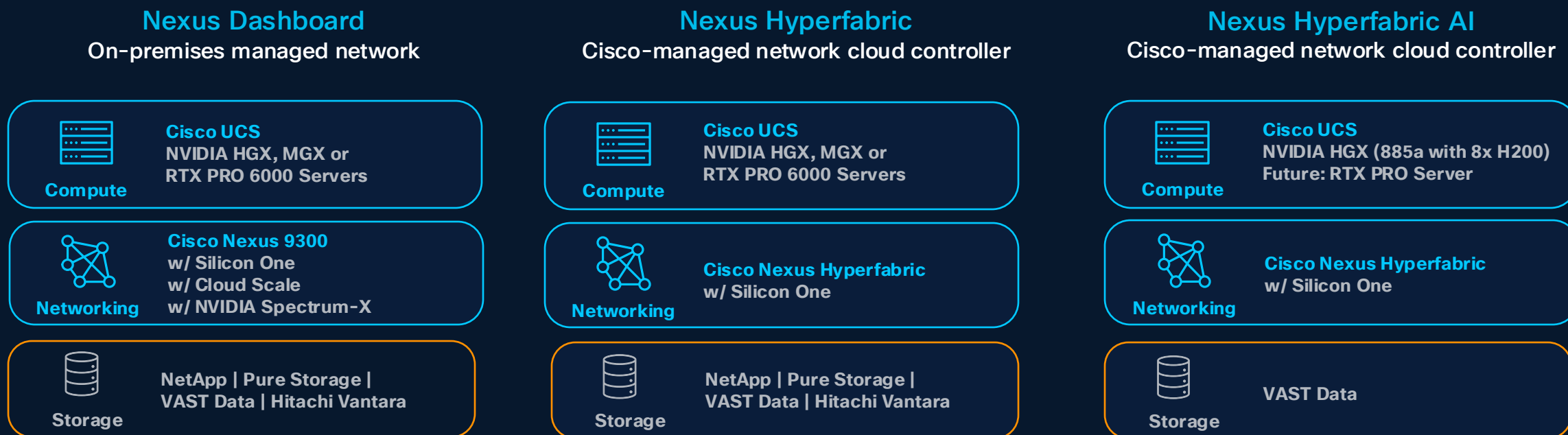
Adding VAST Data Storage

Nexus 9000 Switches

Cisco Validated Architecture Training, optimization and inferencing



Cisco AI PODs deployment options



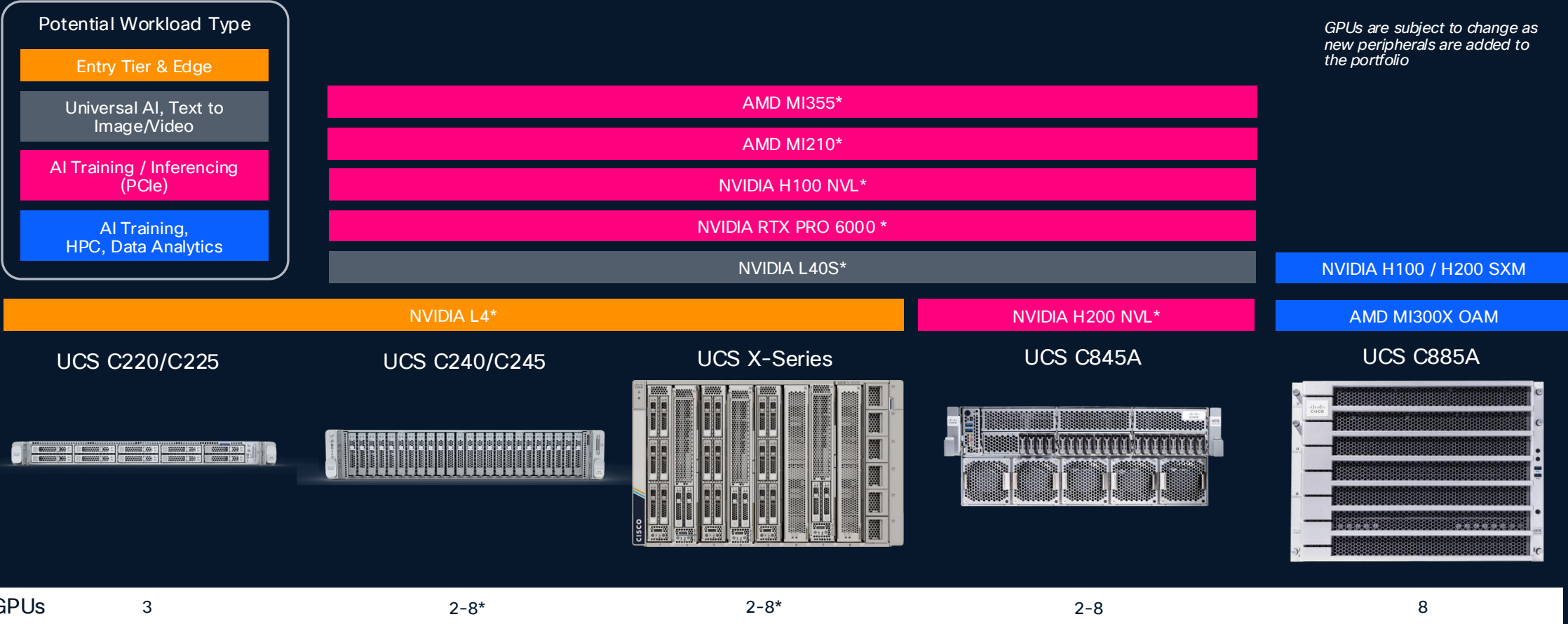
Customizable



Prescriptive

AI Platform Considerations: UCS GPU Options

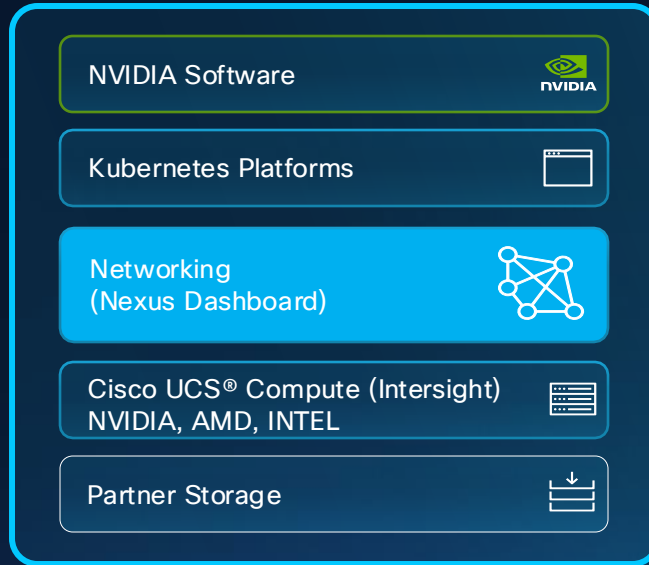
GPUs are subject to change as new peripherals are added to the portfolio



* 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

Cisco AI PODs: Flexible Operating Models

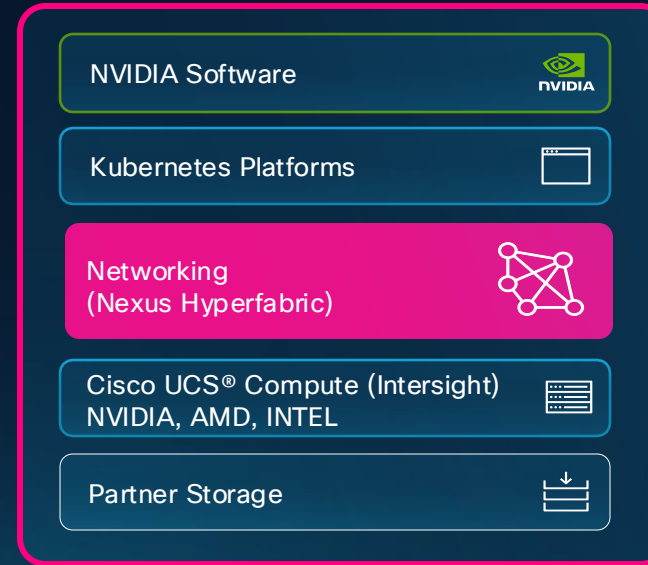
AI POD w/ On-prem management



Modular, pre-validated infrastructure:

- Full stack, buy & deploy
- Nexus Dashboard: On-prem networking management

AI POD w/ Cloud management

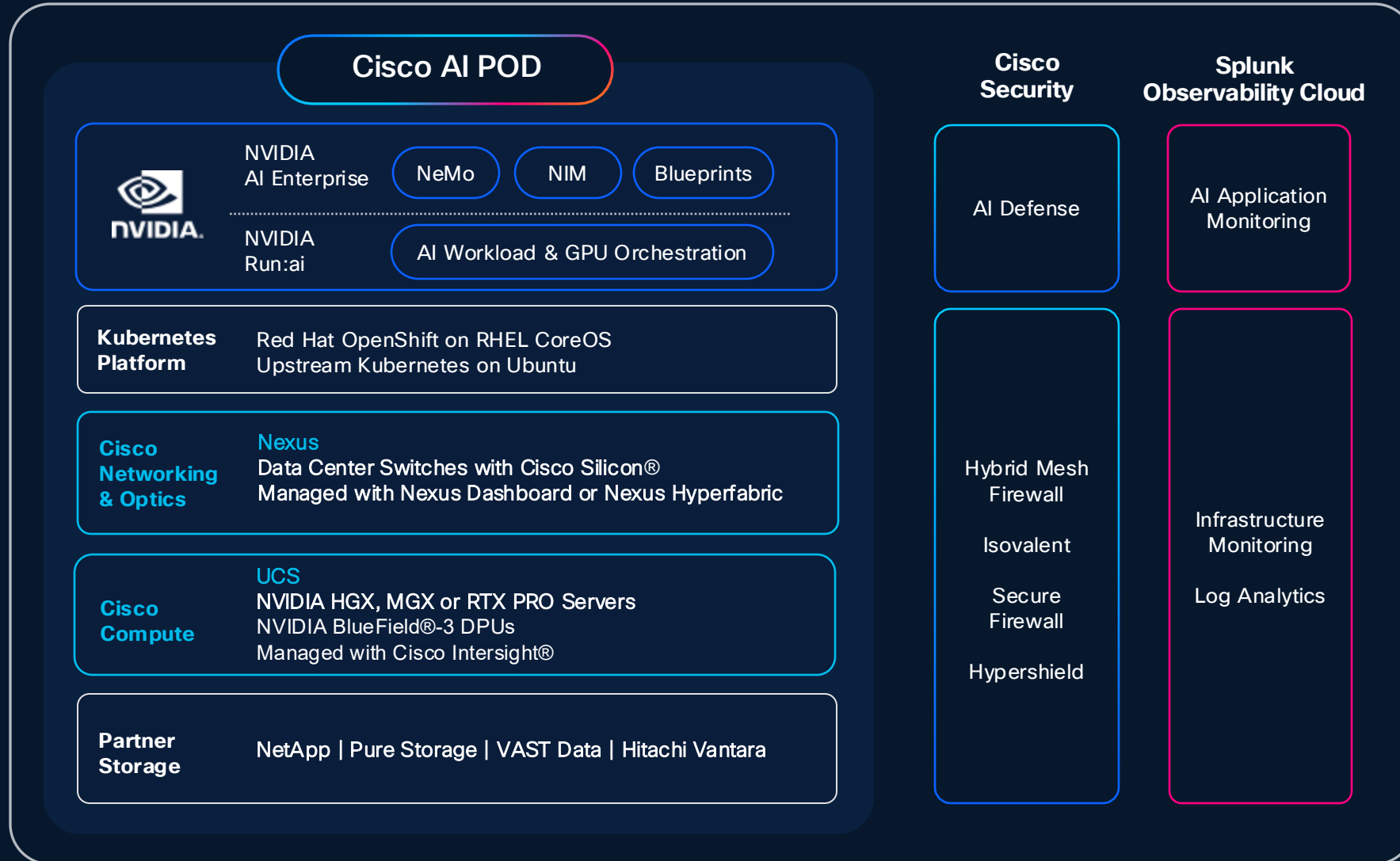


Turnkey infrastructure:

- Full stack, buy & deploy
- Nexus Hyperfabric: Cloud-managed Networking
- Nexus Hyperfabric AI: Cloud-managed physical infrastructure

Cisco Secure AI Factory with NVIDIA

Delivering trusted AI outcomes



Bringing observability to the Factory

Cisco AI Defense on AI POD - Summary



End-to-end security



Continuous visibility and governance

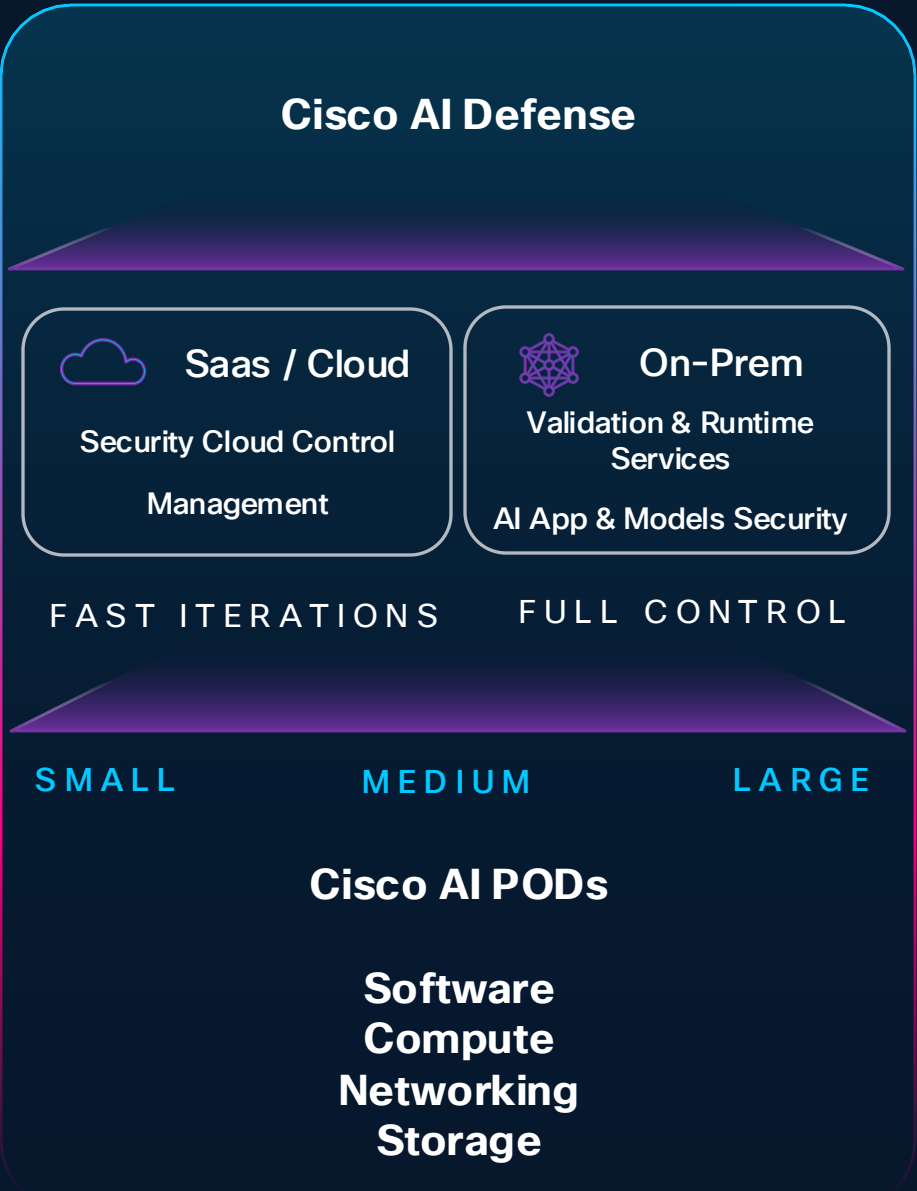


Sovereign ready


- DATA CONTROL
- SECURITY
- USE CASES
- MODEL AND CUSTOMIZATION

Beta
Available


Controlled Availability
Expected in 90 days



Full Application Security



AI Application and Model Validation



AI Runtime Application Protection

Cisco AI PODs

Included in Cisco Secure AI Factory with NVIDIA

Why Cisco AI PODs?



Security-first architecture enables safe enterprise AI



Unmatched performance AI infrastructure enables efficient model training, customization, and inferencing



Pre-validated AI infrastructure stack for simplified deployment drastically reduces set-up time

Typical edge deployment

Retail

User Edge

B Gateways /
access points

On Premises Edge
(Light edge and connectivity)
Meraki/Catalyst Access Points
Far Edge



A Sensors and
smart cameras
Device Edge
Meraki Smart Cameras

C Edge server,
networking and
security
On Premises Edge
(Local heavy compute,
networking)
Cisco Unified Edge
Catalyst
Security, Observability

Inventory Management

A Track items on shelves
Send data to gateways

B Aggregate raw data
Do local pre-processing
Perform light computation*
Send data to edge servers

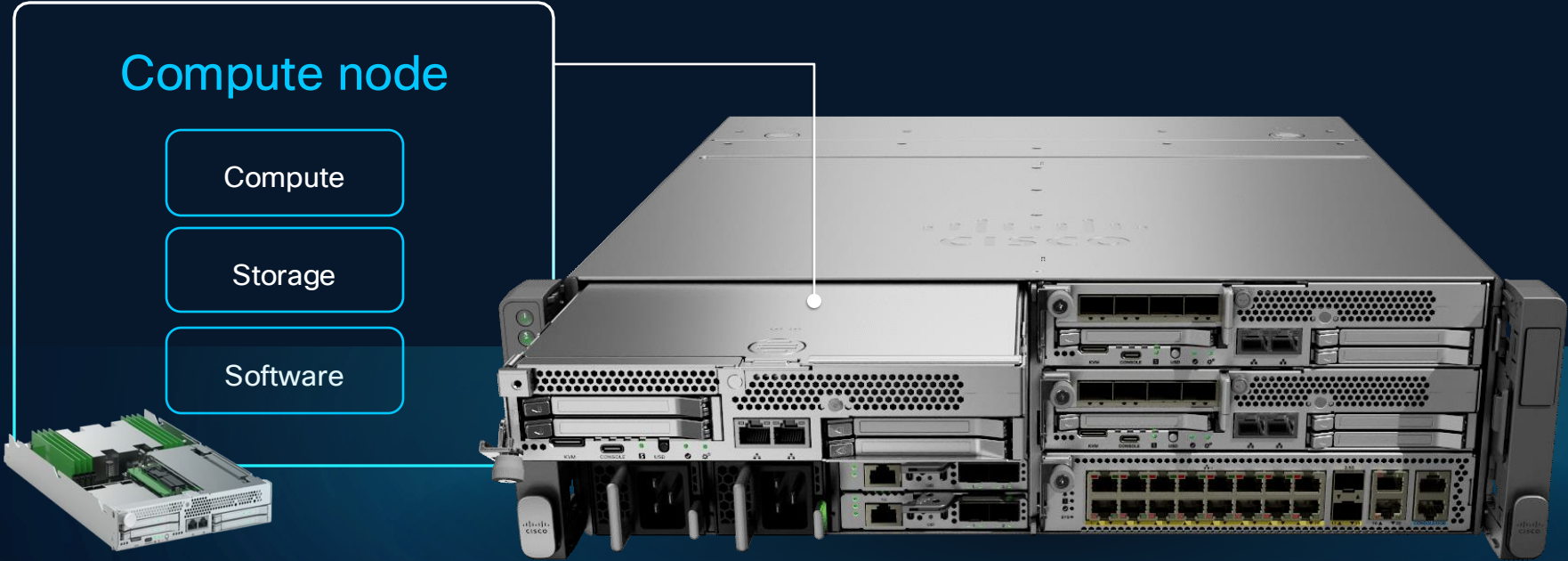
C Run heavier AI analytics /
models locally
Send real-time alerts on
low inventory, out-of-stock

* Gateways can also run lightweight AI models and send real-time alerts

Cisco Unified Edge: Future-Ready Performance

Future-ready performance

Integrates compute, networking, storage, and security



NUTANIX

Nutanix Stack for Edge



Red Hat Edge Stack



VMware Edge Compute Stack



Intel Tiber



Virtual Networking

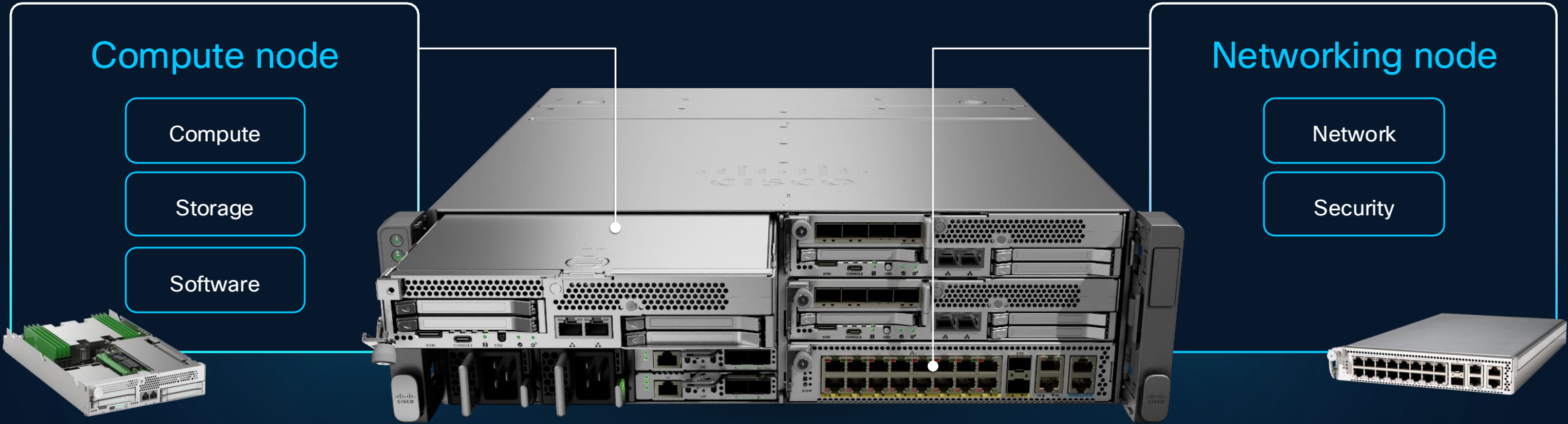


SUSE Edge

Cisco Unified Edge

Fully validated, full-stack environment that integrates advanced network, compute, storage and security

Future-ready performance



Nutanix Stack for Edge

Red Hat Edge Stack

VMware Edge Compute Stack

Intel Tiber

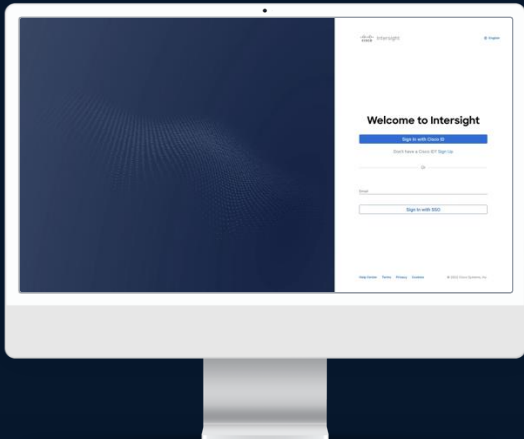
Virtual Networking

SUSE Edge

Day 0: Zero touch provisioning

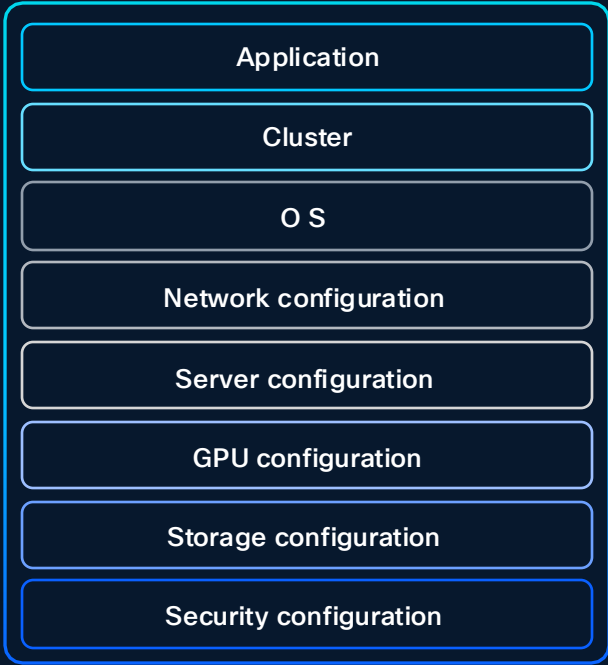
Consistent, repeatable AI-ready infrastructure deployments across multiple sites

Unified operations



Cisco Intersight

Cisco edge blueprints



Golden configurations base on
Cisco validated designs

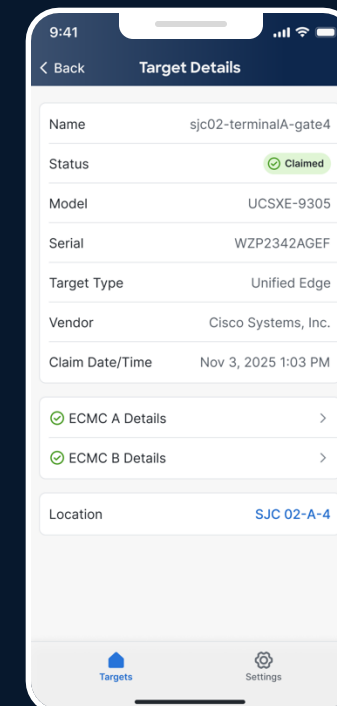
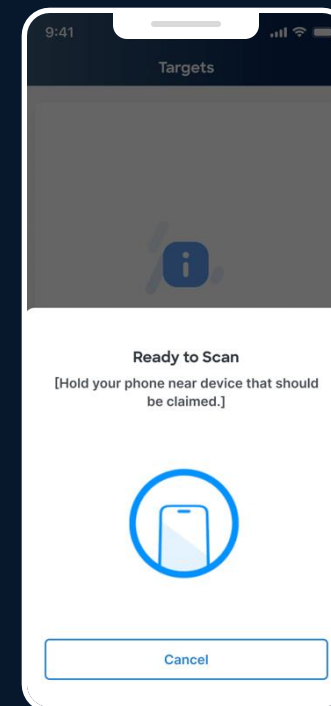
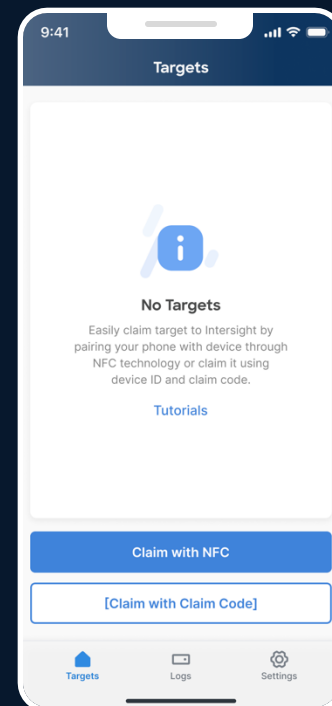
Faster deployments
with fewer errors

Supported by
Cisco TAC

Available in
a marketplace

Day 1: NFC-based claim via mobile app

Unified operations



Fewer truck rolls – no need to onboard systems centrally

No technical expertise needed to claim device

Pre-claimed capabilities accessible to authorized on-site staff

Day 2: Deployment at scale

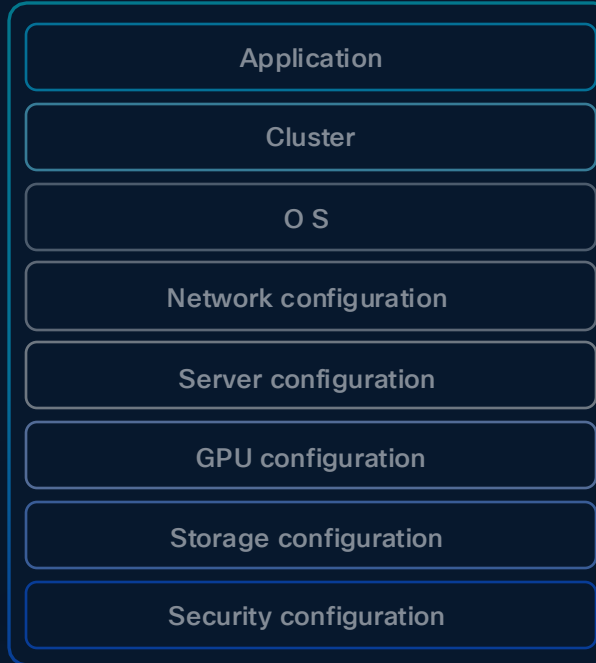
Consistent, repeatable AI-ready infrastructure deployments across multiple sites

Unified operations

Cisco Edge blueprints (Retail, manufacturing, healthcare)



Cisco Intersight

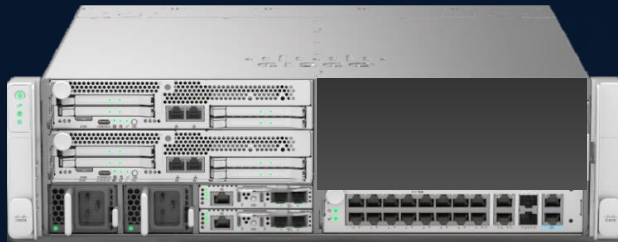


Day N: Seamless, scalable operations

Fleet management at global scale

Unified operations

Deploy a new service



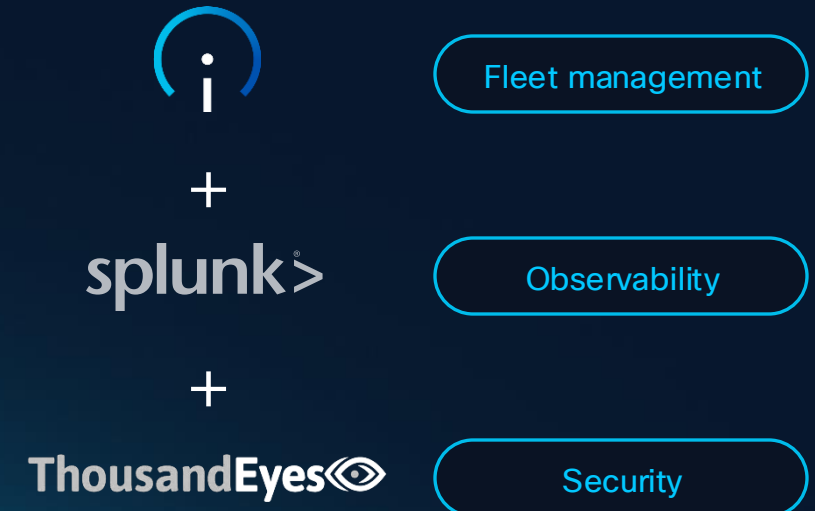
- 1 New AI application requires GPU
- 2 Install a new compute node with GPU
- 3 No forklift upgrade

Avoiding unscheduled down-time



- 1 Intersight identifies predictive maintenance item
- 2 Replacement node sent out automatically
- 3 Easy on-site install

End to end visibility & observability



Thank you



