

Red Hat  
**Summit**

**Connect**

# It's time to simplify AI infrastructure and operations

Stefano Gioia, EMEA Solution Eng., Cloud & AI  
Cisco

Milan, 19 Nov 2024



# Agenda

- Cisco AI-Ready Data Center
- The value of the Cisco & Red Hat Partnership
- Call to Action

# By the end of this session...

- Become familiar with the Cisco AI-Ready Data Center
- Understand the different options available within AI PODs
- Know the value of the Cisco & Red Hat Solution

# Cisco & Red Hat Go-to-Market Solutions

Backed by Cisco Validated Design and Solution Support

Application Platform  
Modernization



VM  
Options

AI Ready  
Infrastructure



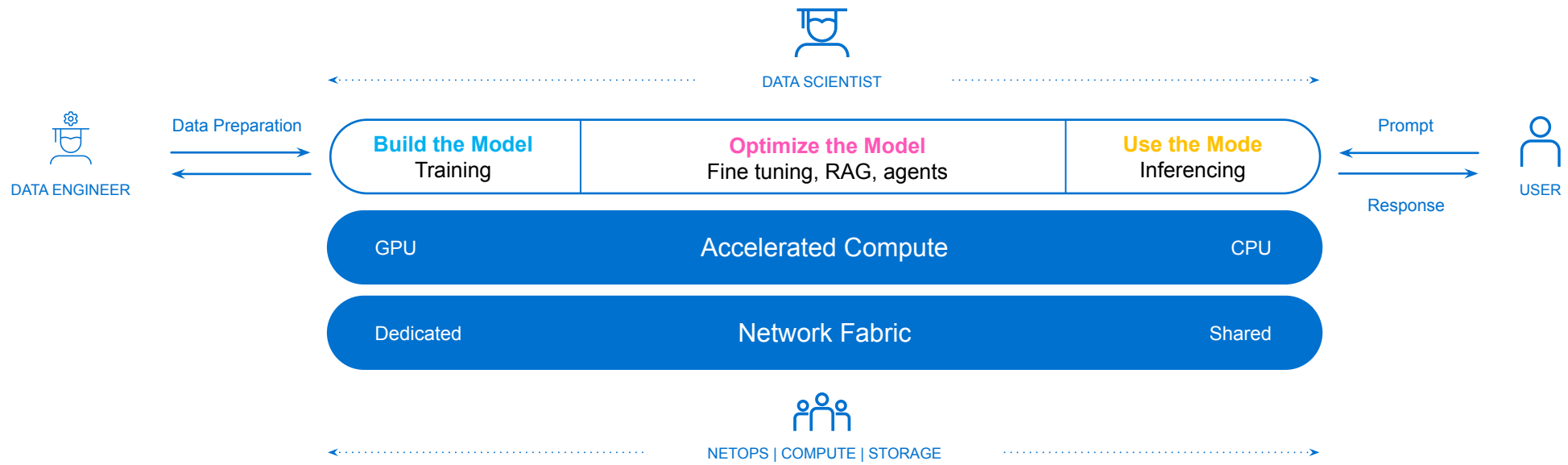
AI Ready

Edge Workload  
Platforms

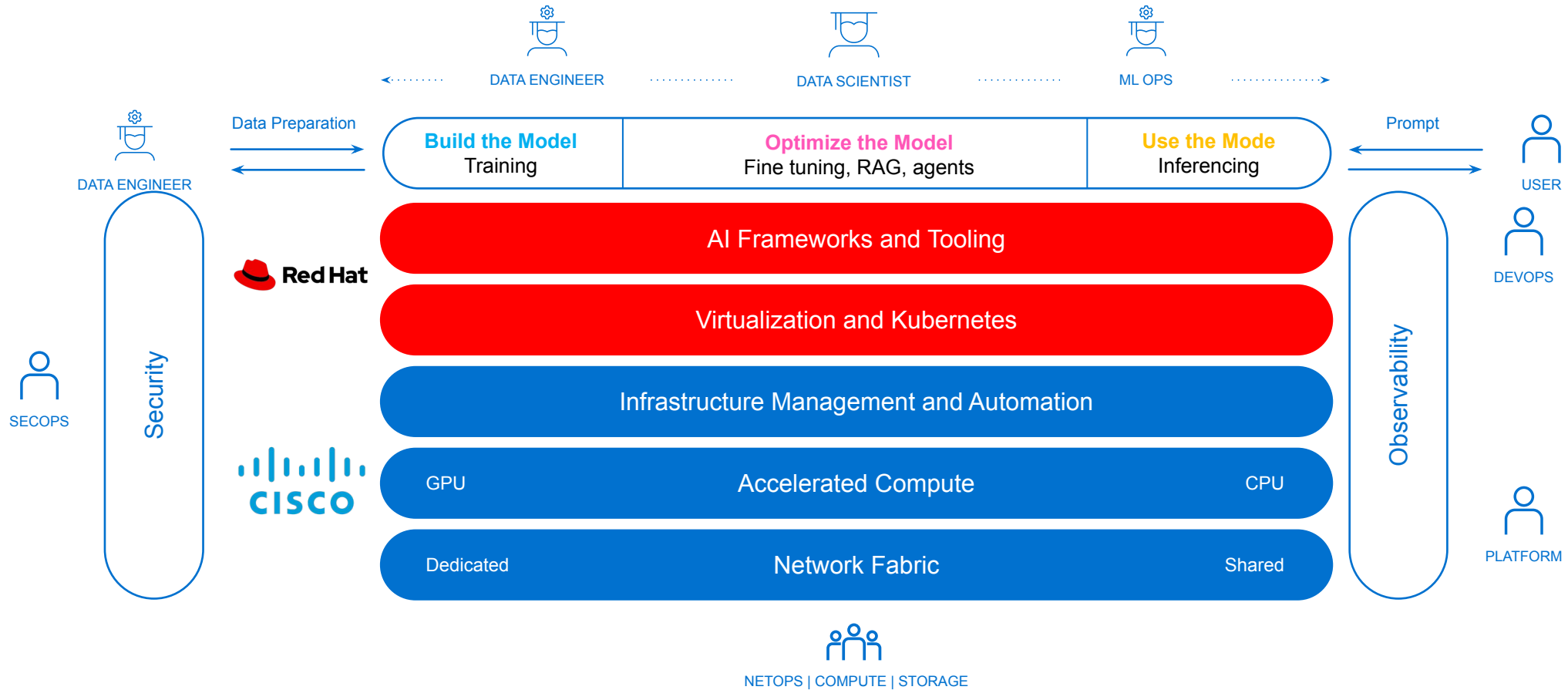


Edge

# Generative AI – Full Stack Systems

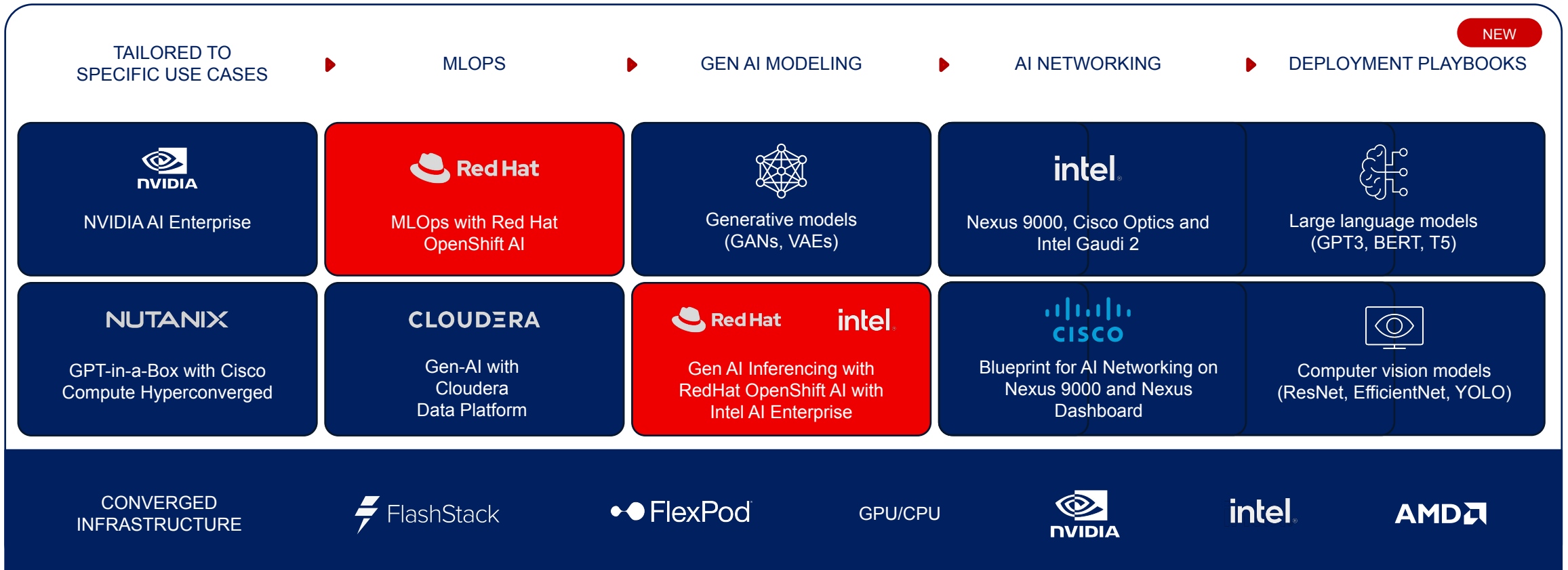


# Generative AI – Full Stack Systems



# Solutions to Simplify and Automate AI Infrastructure

Cisco Validated Designs for the Data Center

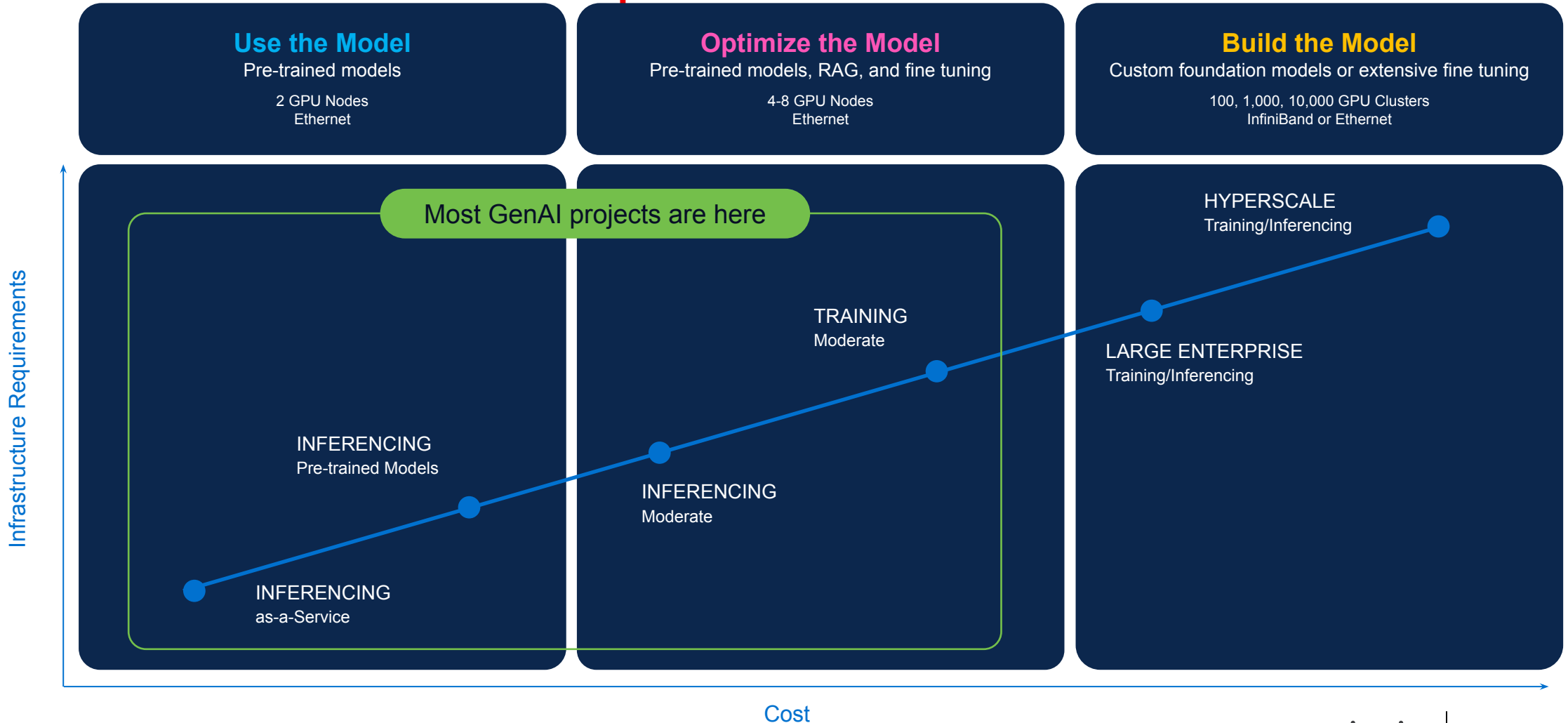


# Every organization's AI approach and needs are different





# Generative AI is a Spectrum



# Every organization's AI approach and needs are different



# Introducing Cisco UCS C885A

Building high-density GPU servers to the Cisco UCS family and to Cisco's AI solution portfolio

Discover data-intensive use cases like model training and deep learning



UCS Accelerated  
UCS C885A M8

Nvidia HGX with  
8 Nvidia H100/H200 GPUs  
AMD Mi300X  
2 AMD 4<sup>th</sup> Gen  
EPYC™ Processors

# Every organization's AI approach and needs are different



# Introducing Cisco AI PODS

## Faster time to value with pre-configured bundles

Deploy AI with confidence

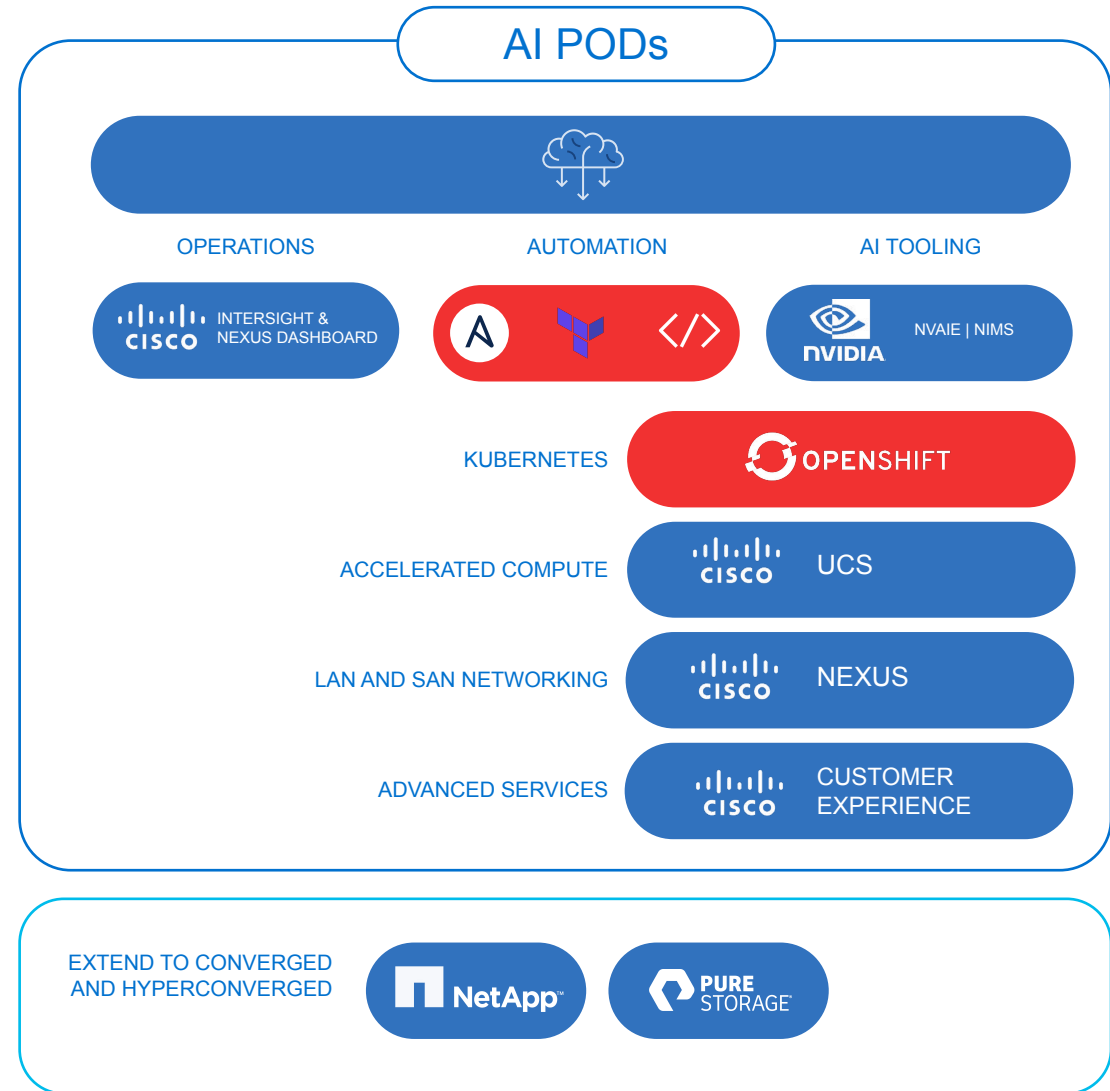
Orderable, validated AI-Ready infrastructure stacks

Fully supported stack including Cisco and 3<sup>rd</sup> party components

AI Advisor tool for configuration guidance

COMING SOON

## Cisco AI-Ready Infrastructure Stacks



# AI PODs Use Cases

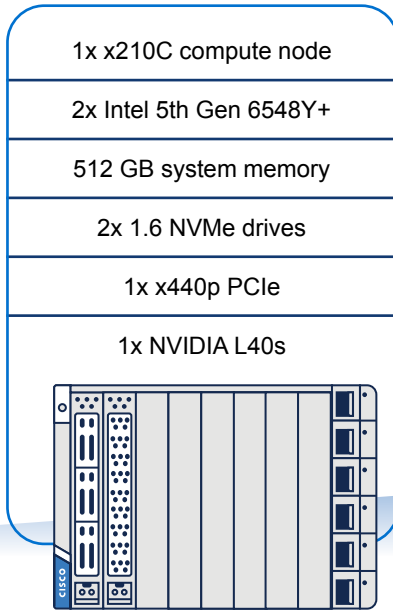
Typical use case

Sizing Example

Hardware specification

Data Center & Edge Inferencing

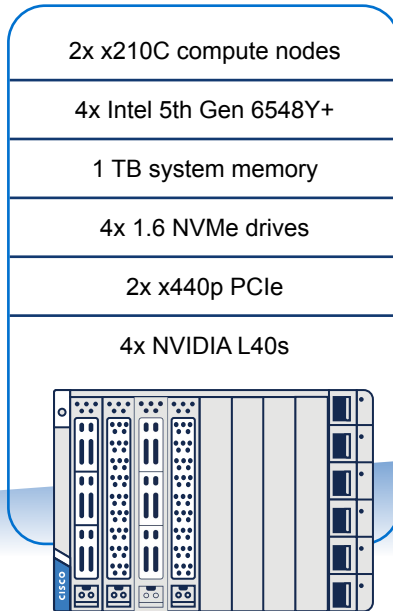
(Llama-2 7B GPT 2B)



RAG Augmented Inferencing

Contextual Accuracy

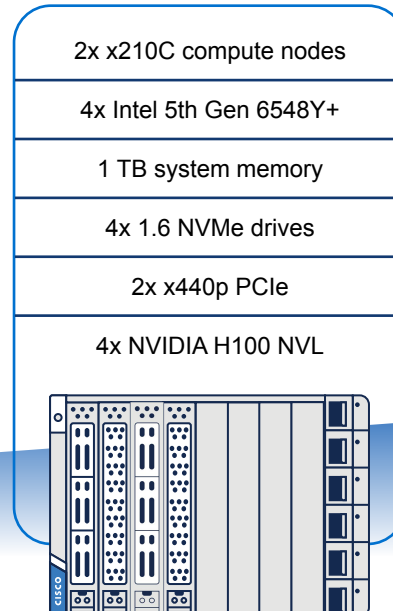
(Llama-2 13B OPT 13B)



Scale Up for High Performance

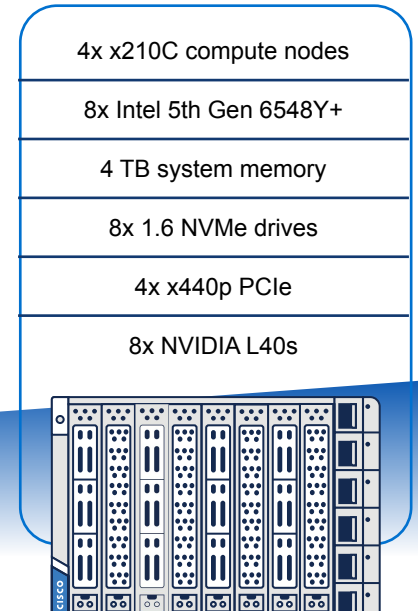
Expanded Context

(Code Llama 34B Falcon 40B)



Scale Out for Large Deployments

Multi-Model Deployments  
High Concurrency

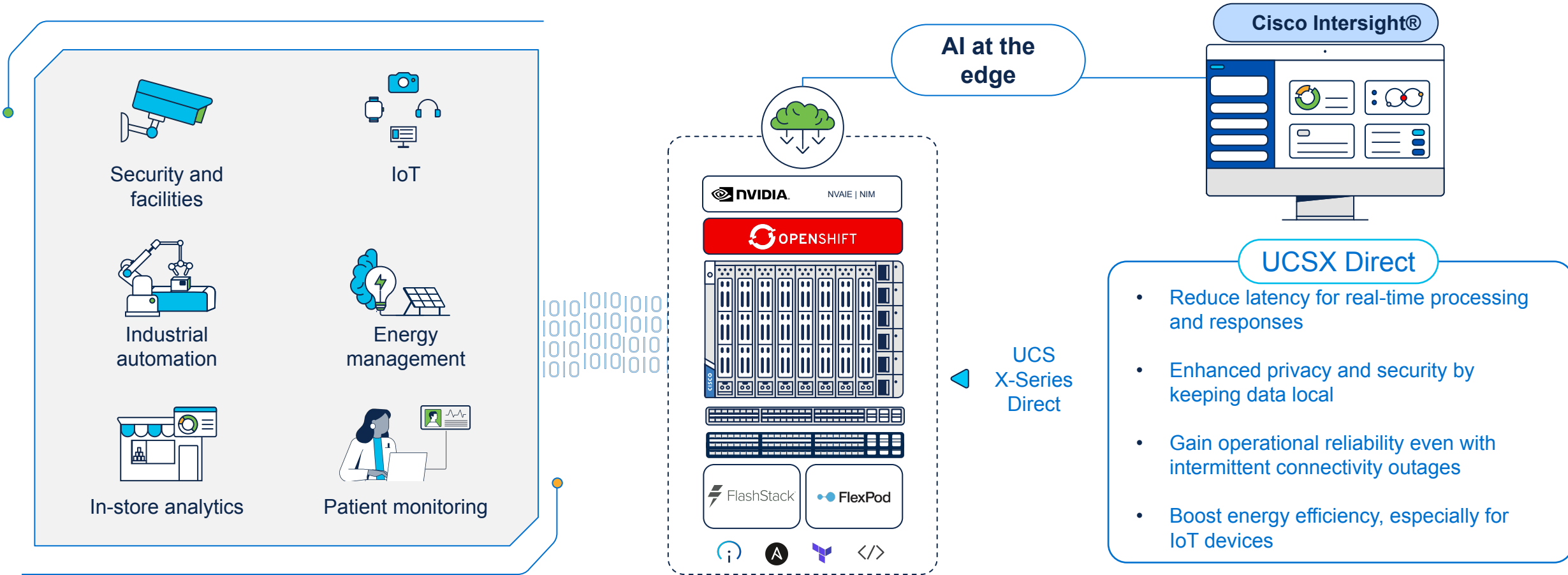


Performance and Scale

# Every organization's AI approach and needs are different



# Inferencing at the Edge



“It’s time to simplify AI infrastructure and operations” – Cisco System



# Cisco AI Pods Benefits

**1** Simplified  
Purchasing  
Experience

**2** Seamless  
Deployment  
and Integration

**3** Efficient and  
Scalable  
Operations

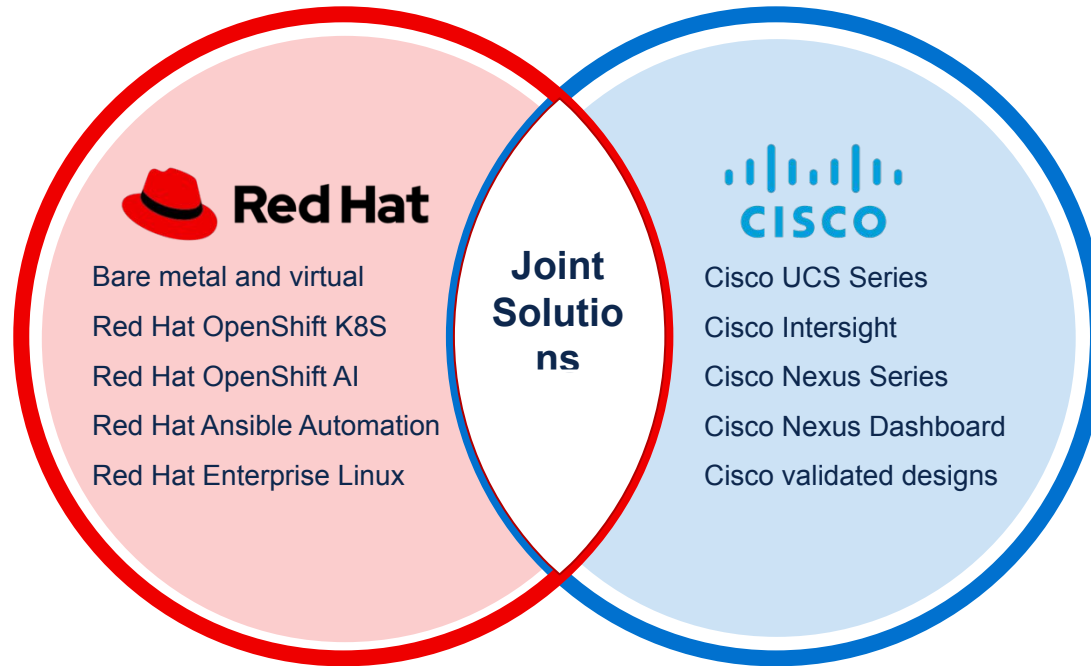
# Cisco + RedHat Partnership

## Open Cloud Infrastructure

platform built on open-source innovation

## Accelerated time to value

with turnkey experience and integrated automation  
For VMs and C



## Simplified Operations and Support

with Cloud managed infrastructure and Cisco Solution Support across Red Hat on converged infrastructure stacks

## Reduced Risk

with Cisco Validated Designs, delivering tested architectures for standardized, repeatable deployments.

Operate across hybrid multicloud

More choice and flexibility

20+ Cisco Validated Designs

Consistent app dev experience

Increased sustainability

# Cisco Validated Design with Red Hat



FlexPod Datacenter with Red Hat OCP Bare Metal  
Manual Configuration with Cisco UCS X-Series Direct

Updated: September 19, 2024

Bias-Free Language Contact Cisco

Save Download Print

Published: September 2024

In partnership with:

Table of Contents

- Table of Contents
- About the Cisco Validated Desi...
- Executive Summary
- Solution Overview
- Deployment Hardware and Sof...
- Network Switch Configuration
- NetApp ONTAP Storage Confi...
- Cisco Intersight Managed Mod...
- OpenShift Container Platform I...
- Deploy a Sample Containerize...
- About the Authors
- Appendix
- Feedback

Red Hat OpenShift is now supported on bare metal UCS/UCSX servers! Cisco CVD's will no longer require VMware to remove layers and simplify our solutions.

Just released:

- [Flexpod Datacenter with Red Hat OpenShift Bare Metal](#)
- More OpenShift on bare metal CVD's coming soon:
- FlexPod with OpenShift Virtualization
- FlexPod with OpenShift AI
- FlashStack with OpenShift Virtualization
- FlashStack with OpenShift AI
- [Datacenter Cisco Validated Design Center](#)

Red Hat  
**Summit**

**Connect**

Q&A





**Connect**

Thank you

