# From Complexity to Confidence: Core-to-Edge Secure AI Factories with Consistent Operations
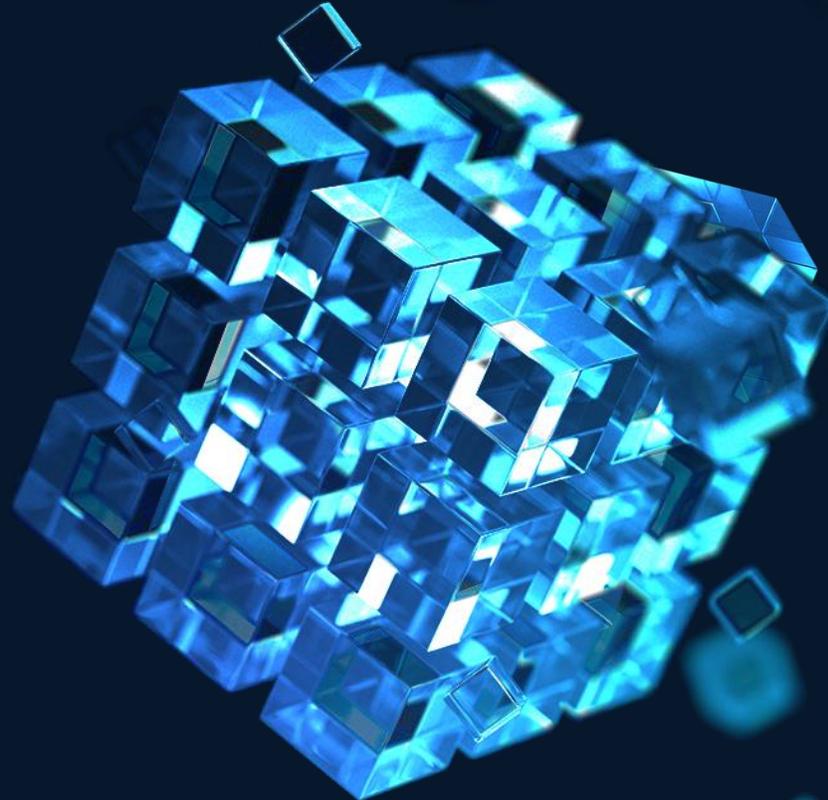
**Kevin Wollenweber**
SVP/GM, Data Center and Internet Infrastructure

# Where does complexity live?

## Hardware is only part of the story

Validation cycles

Compliance overlays

GPU cost

Hardware spend

Incident debugging

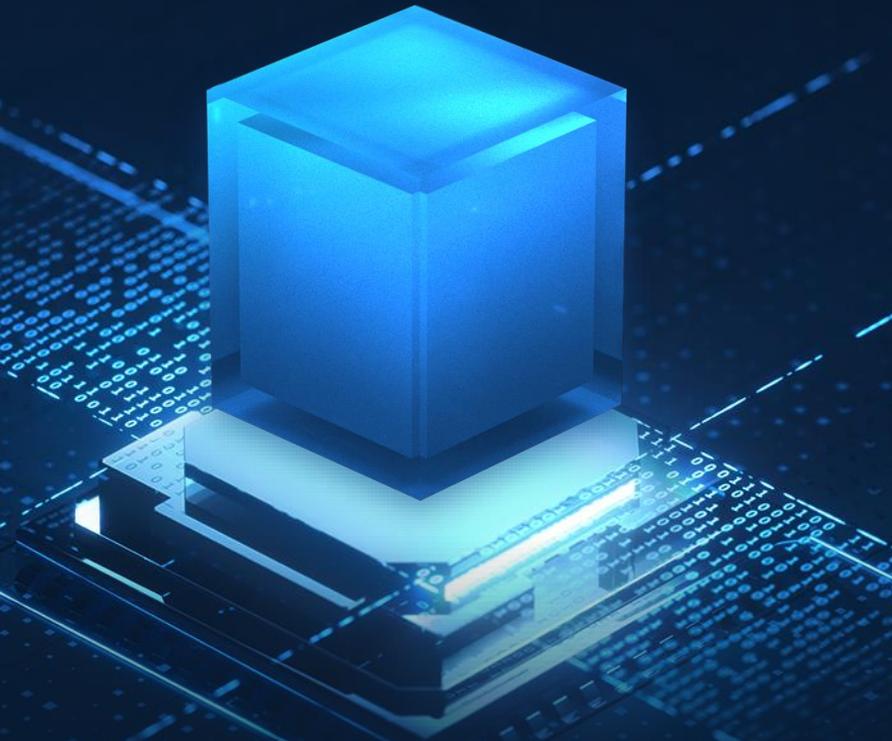Energy inefficiencies

Engineering hours

Tuning time

# Cisco is on a journey to making AI simpler

# A full stack approach
# changes your starting point

## Enables to start from a use case vantage point

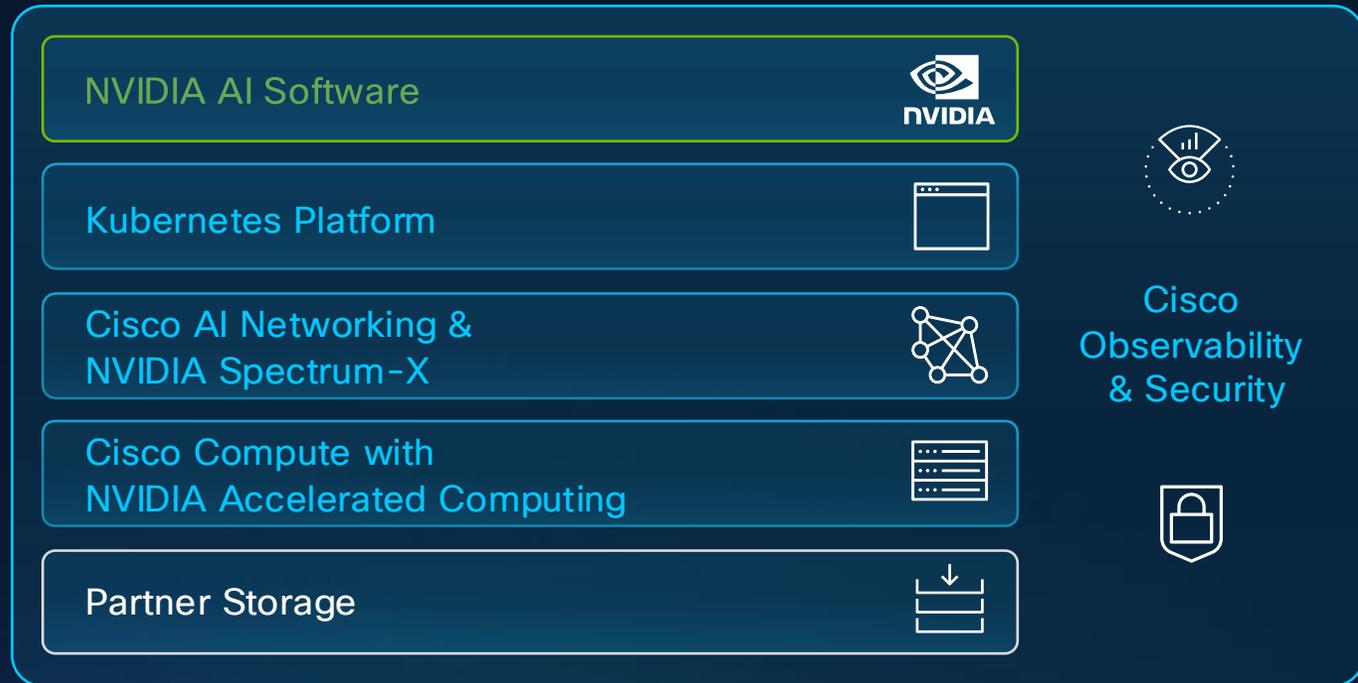| Pre-validated architectures | Operationalize from day one | Reduced integration cycles | Accountability |

# Cisco Secure AI Factory with NVIDIA

## Secure. Scalable. Simple.

NVIDIA AI Software

Kubernetes Platform

Cisco AI Networking &
NVIDIA Spectrum-X

Cisco Compute with
NVIDIA Accelerated Computing

Partner Storage

Cisco Observability & Security

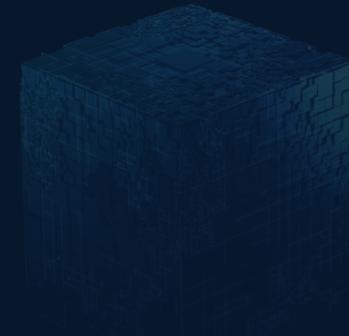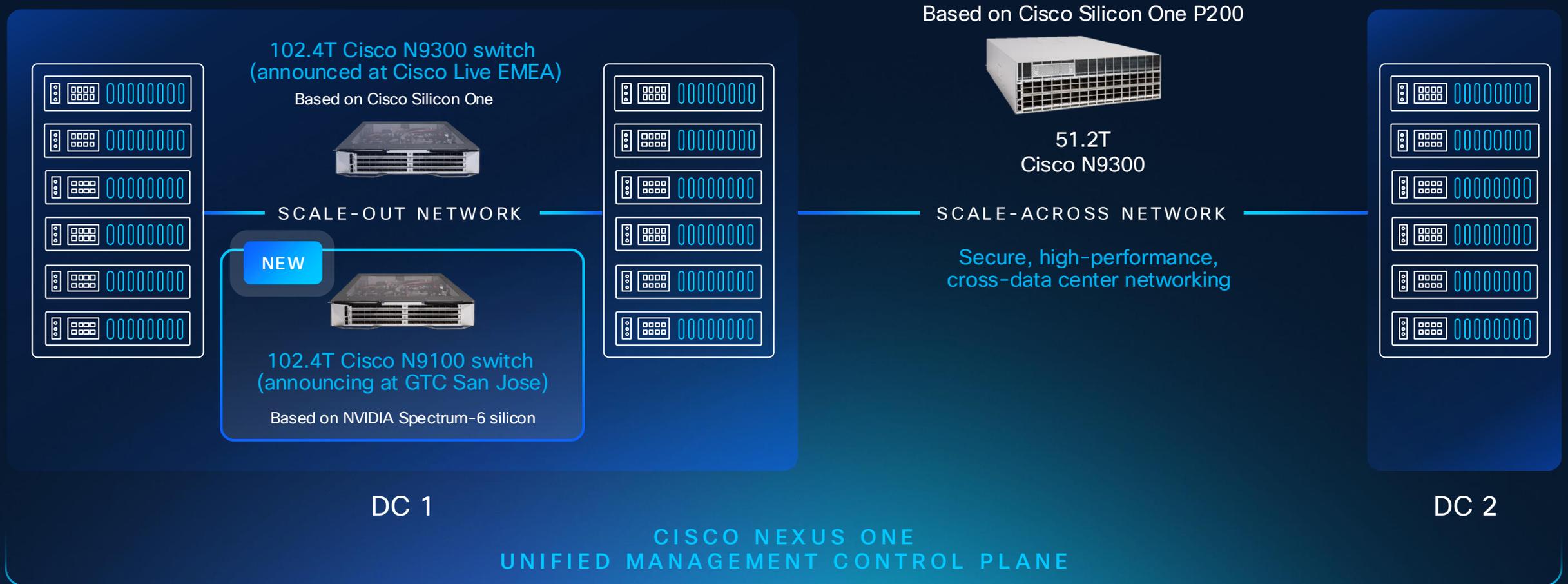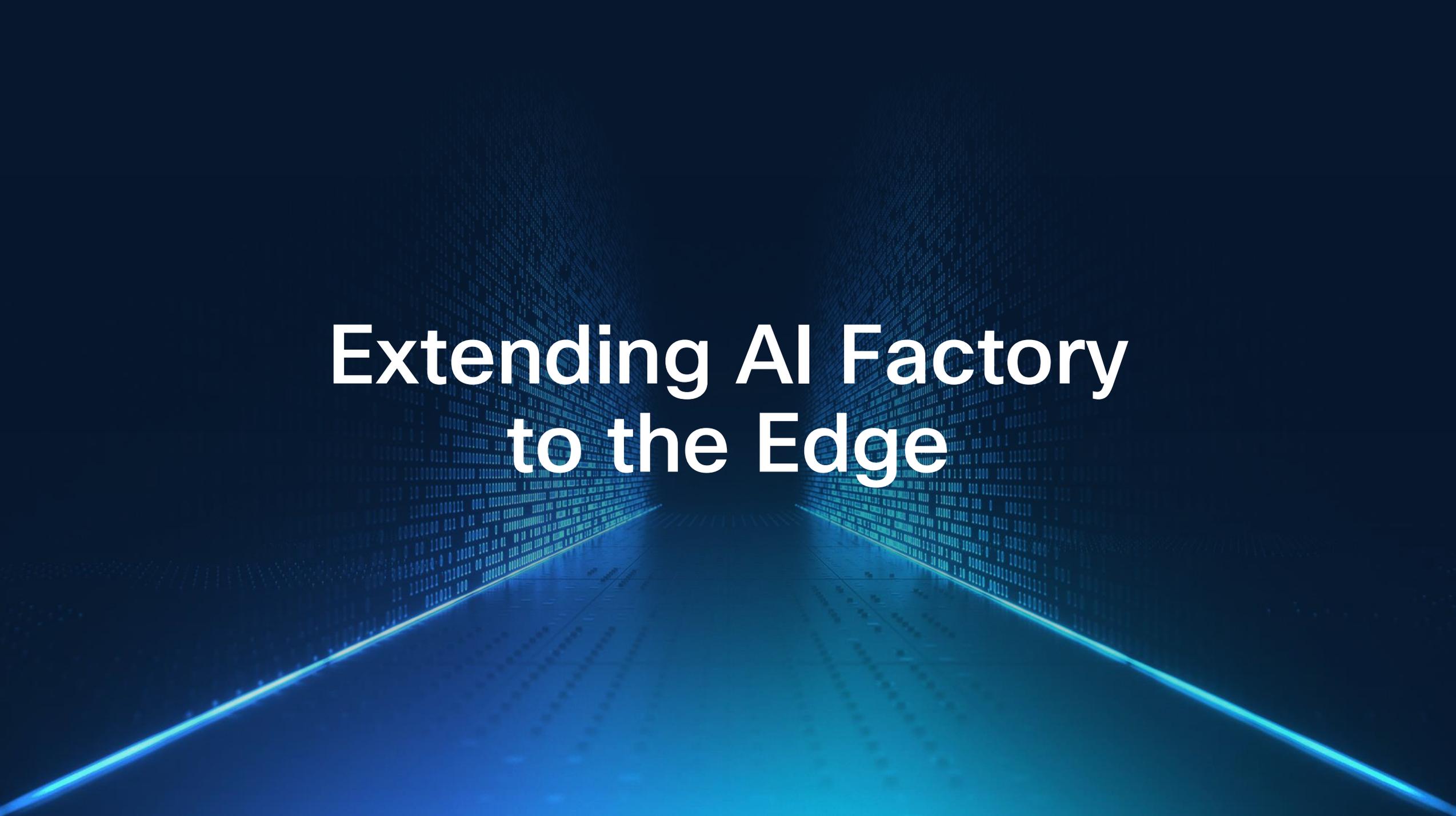| Cisco Enterprise Reference Architecture (ERA) | Cisco Enterprise Reference Architecture (ERA) | Cisco Cloud Reference Architecture (CRA) | Cisco Cloud Reference Architecture (CRA) |
|---|---|---|---|
| < 1024 GPUs | NVIDIA ERA Compliant  < 1024 GPUs | 1K~32K GPUs | NVIDIA NCP RA Compliant  1K~32K GPUs |

# Scalable

Secure

Simple

# Sovereign & Neocloud AI at Giga-Scale

### Enabling customer need for scale-out and scale-across networking

### Choice of NCP and Cisco CRA Compliant architectures

Based on Cisco Silicon One P200

102.4T Cisco N9300 switch
(announced at Cisco Live EMEA)

Based on Cisco Silicon One

51.2T
Cisco N9300

SCALE-OUT NETWORK

SCALE-ACROSS NETWORK

**NEW**

Secure, high-performance,
cross-data center networking

102.4T Cisco N9100 switch
(announcing at GTC San Jose)

Based on NVIDIA Spectrum-6 silicon

DC 1

DC 2

CISCO NEXUS ONE
UNIFIED MANAGEMENT CONTROL PLANE

# Extending AI Factory to the Edge

# AI Factories at the Telco Edge for distributed inferencing

## Cisco Secure AI Grid powered by NVIDIA

AT&T

Data Center

Core
Network

Access
Network

Data Center

**Real-time
Inference at the EDGE**

AI Software w/ NVIDIA AI Enterprise

**Cisco Mobility Services**
**Reasoning Models and Inferencing**

Kubernetes Platform

Cisco AI Networking

Cisco Compute w/
NVIDIA Accelerated Computing

Cisco Security

Splunk Observability

**Service-ready infrastructure**

End Point Swarm
5G/6G Devices

S I M P L E ,  S E C U R E  I N T E L L I G E N T  C O N N E C T I V I T Y

# Cisco Intersight Unified Fleet Management

## Broad Spectrum of Cisco accelerated computing from Core to Edge

### TRAINING

**Dense GPU Server**

C845A M8

GPUs supported:
NVIDIA RTX Pro 6000,
H100, H200,
AMD MI210, L40S

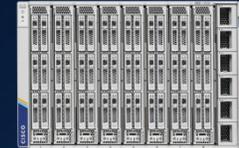### OPTIMIZATION & INFERENCING

**Rack Server**

C245 M8 /
C240 M8

GPUs supported:
NVIDIA RTX Pro 6000,
H100, H200,
A16, L40S, L4

**Modular Servers**

X-Series

GPUs supported:
NVIDIA RTX Pro 6000,
H100, H200,
A16, L40S, L4

### INFERENCING

**Unified Edge**

GPUs supported:
NVIDIA L40S

**NEW** | NVIDIA RTX Pro 4500 GPUs
Price, performance efficiency for inferencing and data analytics

Data center (Core)     Regional Edge     Edge

# Extending Choice of AI Software

Broader choice of tooling for AI practitioners to speed up development and delivery of AI applications

**NEW**

NVIDIA AI Enterprise

Red Hat OpenShift, Nutanix Kubernetes Platform, Upstream Kubernetes

Red Hat AI Factory software

NVIDIA AI Enterprise

Red Hat AI Enterprise*

Data center (Core)

Regional Edge
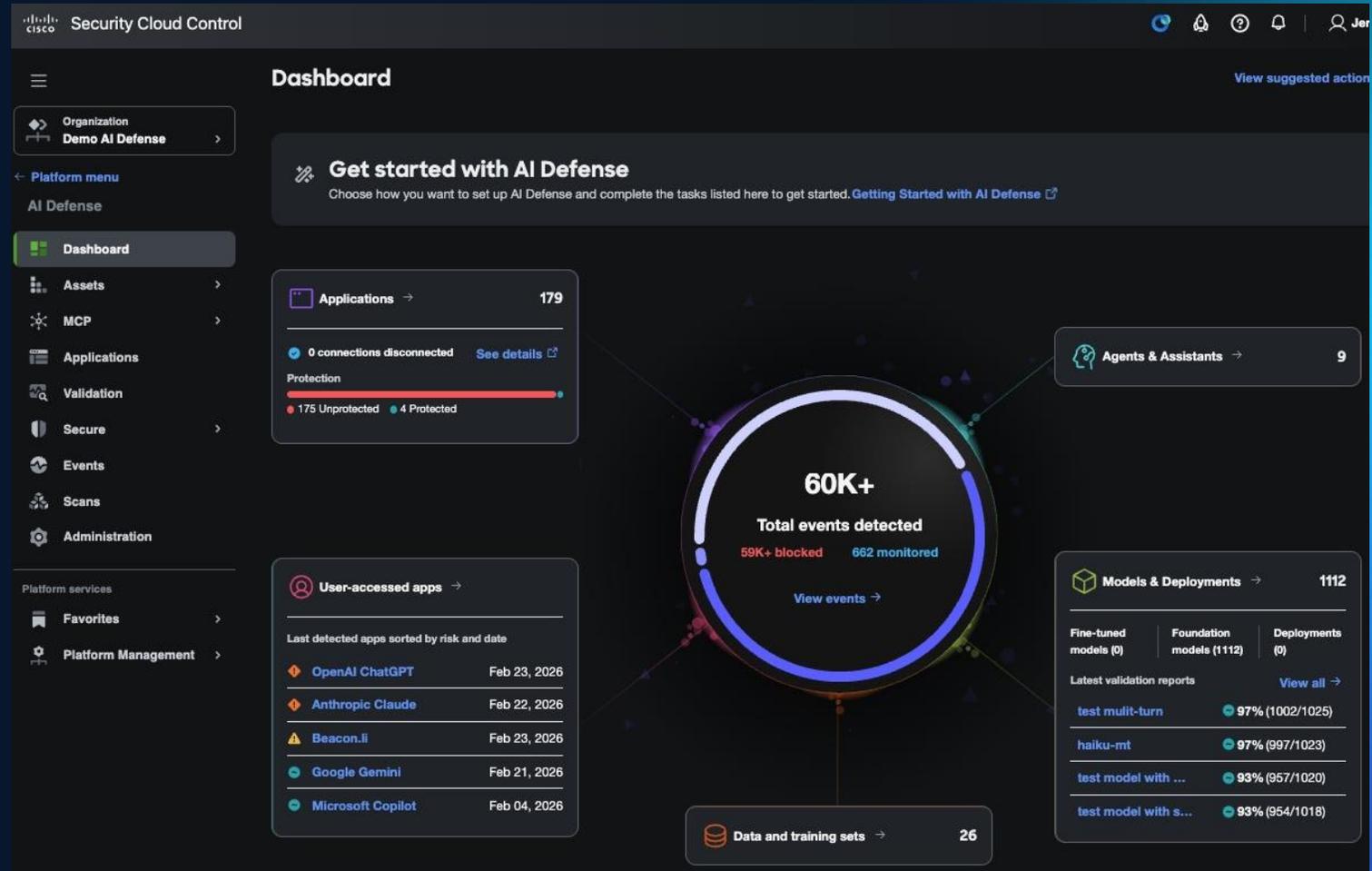
Edge

* Includes Red Hat OpenShift Container Platform

**Safeguard AI models, agents, and applications from safety and security risks with Cisco AI Defense**

Identify vulnerabilities

Mitigates threats in real time
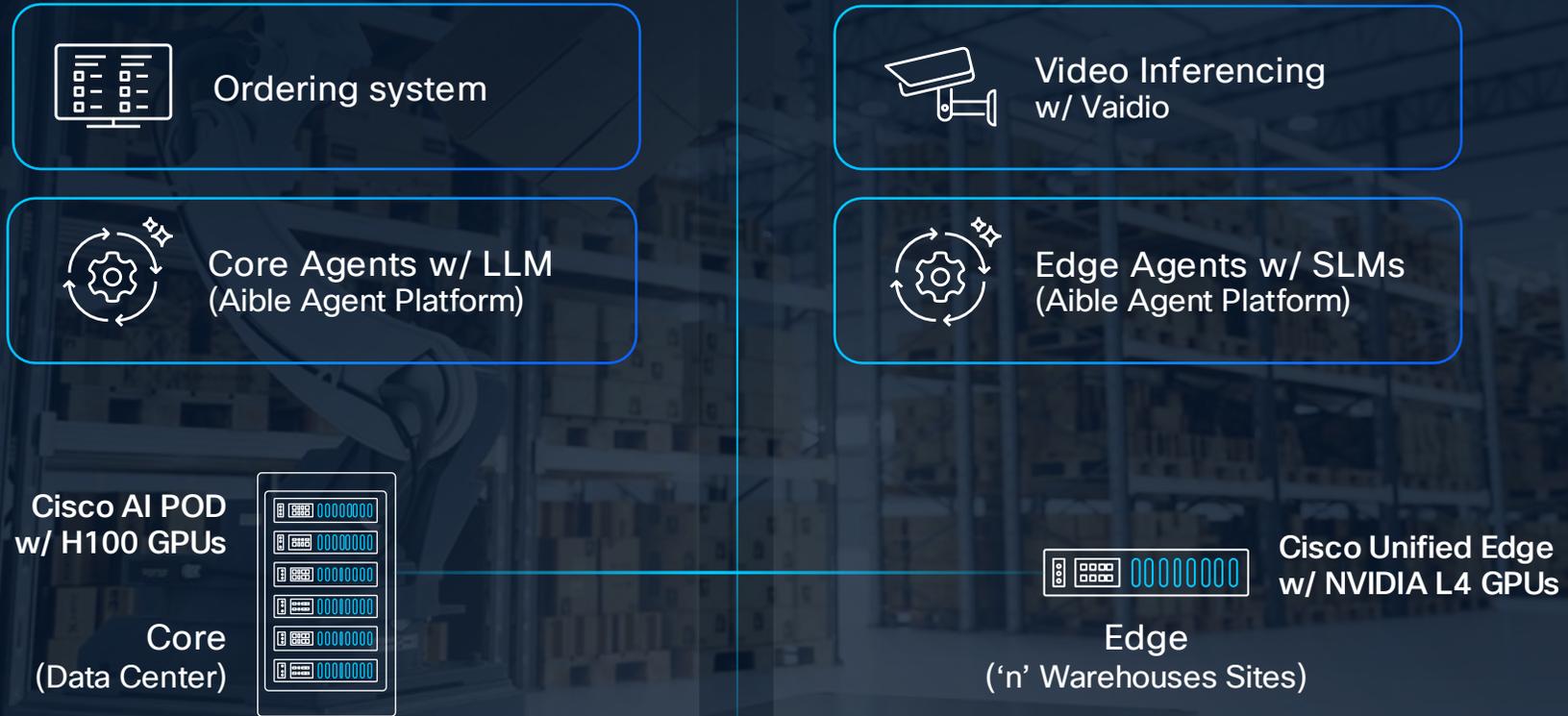
# Security close to every workload
# Hybrid Mesh Firewall

**NEW**

Cisco Hybrid Mesh Firewall extends security policy enforcement on NVIDIA BlueField DPUs on AI servers to ensure protection without negatively impacting CPU & GPU performance.
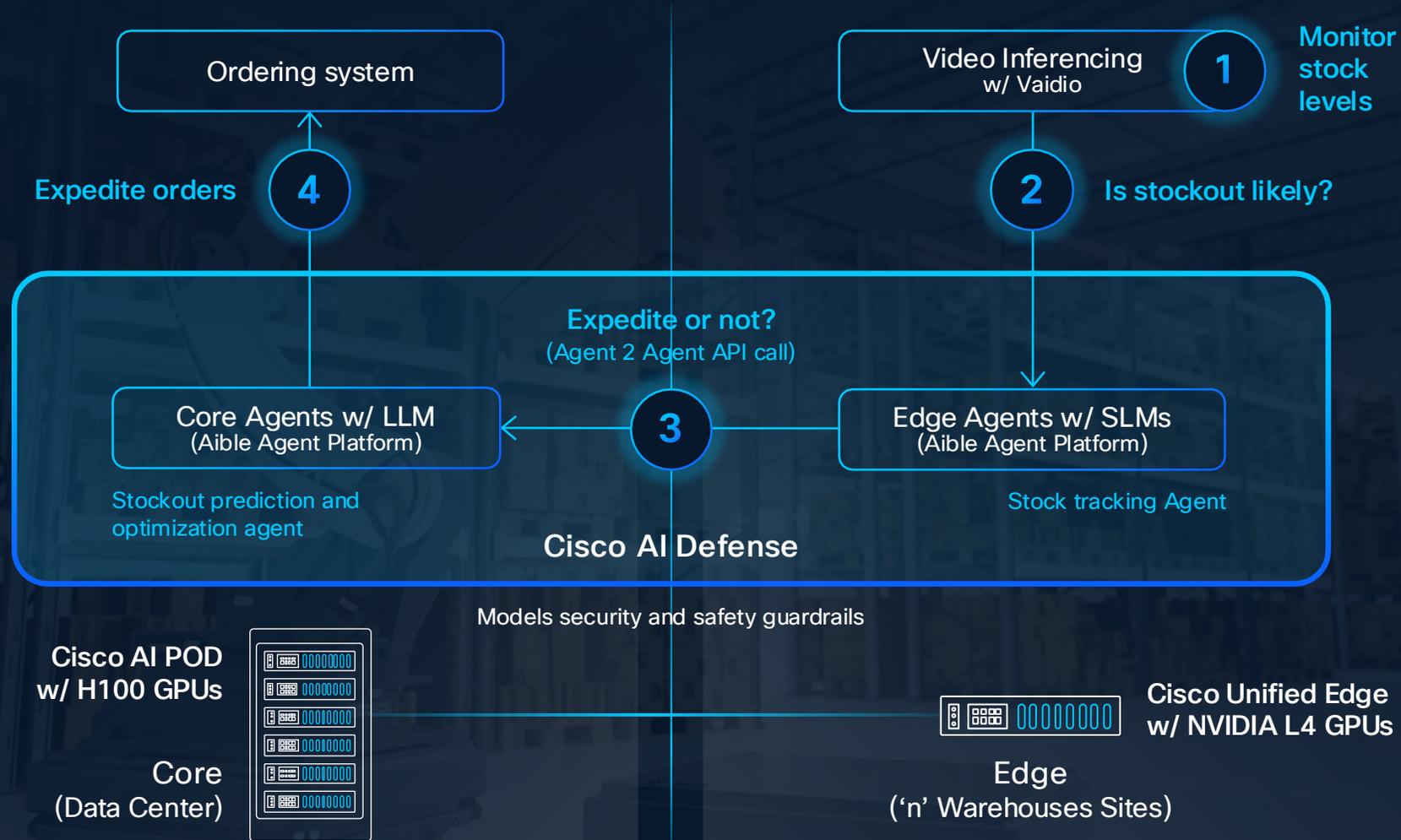
Hybrid Mesh Firewall

**NEW**

eBPF

### Cisco Firewalls
Physical | Virtual | Cloud | FWaaS

### 3rd Party Firewalls

### NVIDIA Bluefield DPU Firewall on AI Servers

### Smart Switches

### Workload Agents

Define policy once, enforce everywhere

# New solution to showcase secure multi-agent system use case

Use Case: Intelligent Warehouse with secure multi-agent operations

Ordering system

Video Inferencing
w/ Vaidio

Core Agents w/ LLM
(Aible Agent Platform)

Edge Agents w/ SLMs
(Aible Agent Platform)

Cisco AI POD
w/ H100 GPUs

Cisco Unified Edge
w/ NVIDIA L4 GPUs

Core
(Data Center)

Edge
('n' Warehouses Sites)

# New solution to showcase secure multi-agent system use case

Use Case: Intelligent Warehouse with secure multi-agent operations

**Ordering system**

**Video Inferencing**
w/ Vaidio

**1**  Monitor stock levels

**4**  Expedite orders

**2**  Is stockout likely?

**Expedite or not?**
(Agent 2 Agent API call)

**Core Agents w/ LLM**
(Aible Agent Platform)

**3**

**Edge Agents w/ SLMs**
(Aible Agent Platform)

Stockout prediction and optimization agent

Stock tracking Agent

**Cisco AI Defense**

Models security and safety guardrails

**Cisco AI POD
w/ H100 GPUs**

**Cisco Unified Edge
w/ NVIDIA L4 GPUs**

**Core**
(Data Center)

**Edge**
('n' Warehouses Sites)

# Simple

Scalable

Secure

# Unified operations with Cisco Nexus One

FULLY INTEGRATED STACK

## Nexus Dashboard
On premises

## Nexus Hyperfabric
Cloud management

SILICON

Cisco Silicon One
NVIDIA Spectrum-X Ethernet

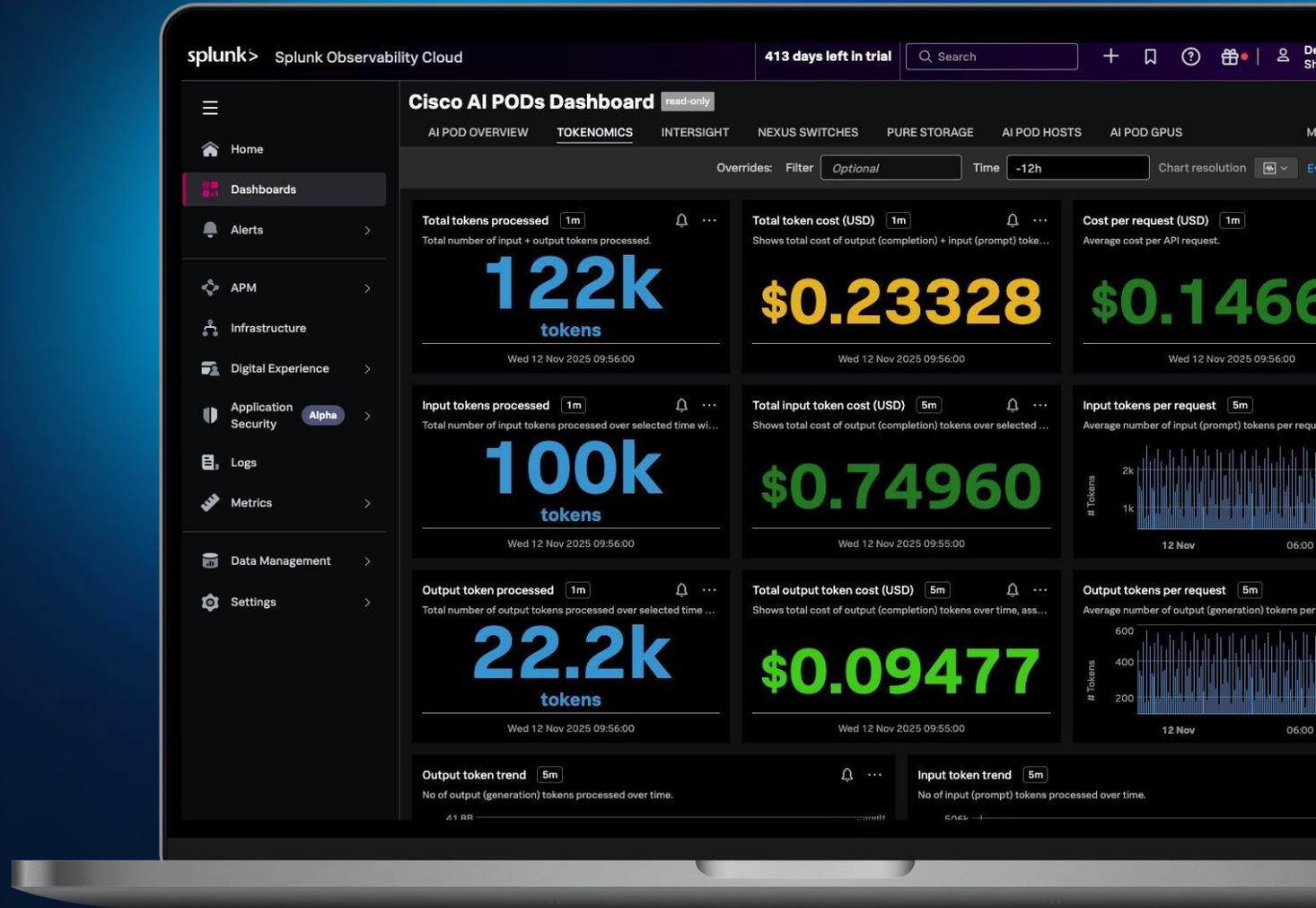SYSTEMS

Cisco N9000
Systems

OPTICS

Cisco Optics

SOFTWARE

Cisco NX-OS
Cisco ACI, SONiC

SECURITY AND OBSERVABILITY

# Splunk Observability for full stack monitoring and reporting

## Tokenomics

**Total tokens processed, Estimated token cost, etc...**

# AgenticOps for Data Center with AI Canvas

## Cisco Deep Network Model

Domain-Specific LLM · 20% higher precision on networking tasks

**Next frontier of operations**

Contextual data center operations with AgenticOps capabilities for troubleshooting

**Precision reasoning**

Powered by purpose-built deep network models

**Breaking silos**

Collaborate and correlate across entire infrastructure using natural language



Generative UI Interface

Shared Workspace

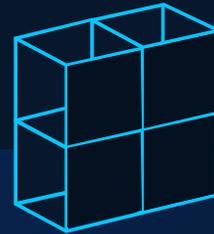AI Assistant

MCP Server

**Cisco Nexus One**
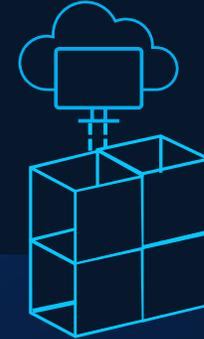
# Flexible deployment options

## Build your own

Buy and deploy individual products, as needed

## AI POD with on-prem network management

Buy & deploy full stack

**Modular, pre-validated**

Backed by CVDs

NVIDIA ERA compliant

## AI POD with cloud based network management

Buy and deploy full stack

Nexus Hyperfabric for cloud-managed networking or turnkey physical infrastructure

NVIDIA ERA and NCP RA compliant

Faster time to first intelligent and ROI with professional services from Cisco and partners, globally

# Fireside chat

**CISCO**

## Kevin Wollenweber

**SVP/GM, Data Center and
Internet Infrastructure**

Cisco

**❤ CVS**

## Saul Mankes

**VP, Digital Platforms and
Infrastructure**

CVS