What Goes Into an Al Cluster

Discover how Cisco accelerates AI adoption with secure, scalable, and high-performance AI-ready datacenter infrastructure — Integrating Cisco's advanced technologies and NIVIDIA's innovations to simplify deployment and support of diverse AI use cases.



Nate Reid Al Solutions Engineer

October 2025

Nate Reid

https://www.linkedin.com/in/natereid/

- Cisco Solutions Engineer with 25+ years of industry experience in tech.
- Blend of experience; including virtual infrastructure, cybersecurity, DevOps, MLOps, container orchestration, and machine learning/Al.
- Prior to Cisco, worked for an Al spinout from ETH Zurich. Focused on computer vision and LLM robustness.
- Detroit based, working with our Central region customers and partners.

Any information provided in this document regarding future functionalities is for informational purposes only and is subject to change including ceasing any further development of such functionality. Many of these future functionalities remain in varying stages of development and will be offered on a when-and-if available basis, and Cisco makes no commitment as to the final delivery of any of such future functionalities. Cisco will have no liability for Cisco's failure to deliver any or all future functionalities and any such failure would not in any way imply the right to return any previously purchased Cisco products.

Agenda

- 01 Al Infrastructure Considerations
- 02 Cisco Al PODs
- 03 Al Workload Orchestration
- 04 Operations and Automation
- 05 MLOps
- 06 Cisco Al PODs Extensibility

But first... a quick level-set on some terminology

- Foundational Model
- Pre-Training
- Fine-tuning
- Inferencing
- Token
- Context
- Parameters
- TTFT

- Transformer
- Attention Head
- Embedding
- RAG (Retrieval-Augmented Generation)
- MCP (Model Context Protocol)
- MLOps (Machine Learning Operations)
- Guardrails
- Al Factory
- ...



But first... a QUICK level-set on some terminology

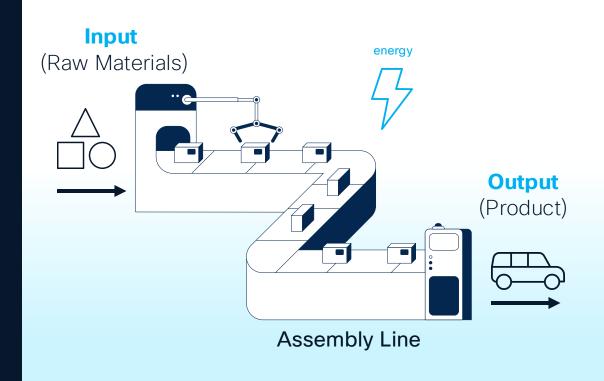
- Foundational Model
- Pre-Training
- Fine-tuning
- Inferencing
- Token
- Guardrails
- Al Factory



The Factory

Repeatable, scalable business capability

- Mass Production Efficiency
- Quality Control
- Supply Chain Integration

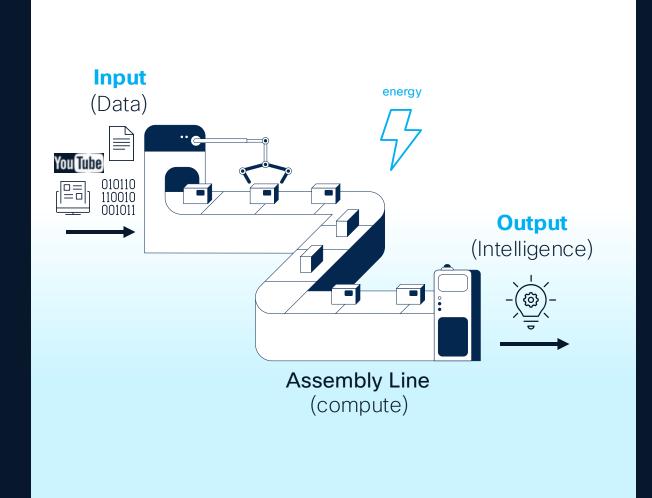


The Al Factory

The Al Factory intuition is to unify hardware + software into a single system. Outcome: Resource sizing is not standalone – compute, network, and storage must be co-designed to avoid bottlenecks.

- Mass model/token production
- Model performance control

Data and workflow integration

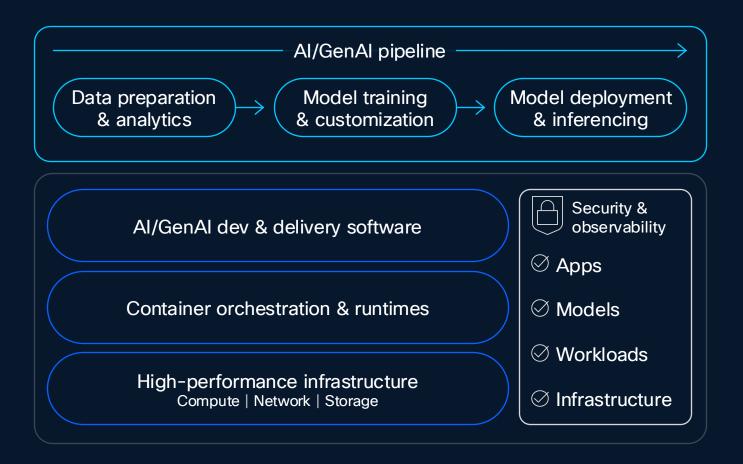


Agenda

01	Al Infrastructure Considerations
02	Cisco Al PODs
03	Al Workload Orchestration
04	Operations and Automation
05	MLOps
06_	Cisco Al PODs Extensibility

Enterprises need capable infrastructure to operationalize Al

Al applications are different, and they are driving demand for new architectures, security mechanisms, and lifecycle management



Faster time to business value

Mitigate risks

Simplify deployment

Reliability, availability, flexibility

Enterprise Al Infrastructure Requirements



Customizing foundational models

Training LLMs from scratch is costprohibitive for the average enterprise



Support multiple, smaller workloads

Enterprises can have many use cases spread across different LOBs, each using an LLM (worst case)



Integrate into existing data centers with ease

Al-enabled enterprise applications often need data, applications and other resources in existing data centers

Enterprise Al Infrastructure Requirements



Operational ease and consistency

Existing DC operational model to manage Al backend/frontend fabrics; simplify with fabric + server templates or IaC



Consistency and simplicity at scale

Building block approach using a spineleaf architecture to scale-out with consistency and predictability



Fit-to-purpose Al infrastructure

Performant infrastructure without compromising on choice; established technologies

Common Al Challenges



Unclear business objectives & priorities

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



Complex Al infrastructure deployment

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



Security vulnerabilities

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



Network performance challenges

Model training and inferencing can generate a lot of traffic, slowing networks. Multi-node training demands significantly more from the network

Cisco Al-Ready Data Center

Transform data centers to power AI workloads anywhere

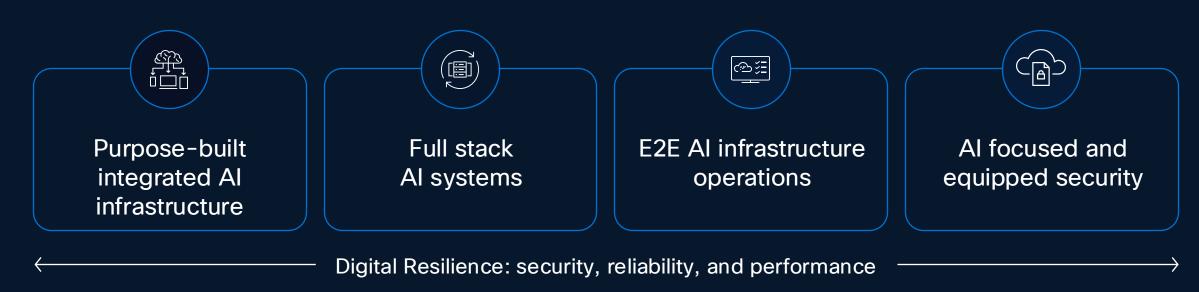
Solving key Al Challenges...







...by leveraging critical components:



Agenda

- 01 Al Infrastructure Considerations
- 02 Cisco Al PODs
- 03 Al Workload Orchestration
- 04 Operations and Automation
- 05 MLOps
- 06 Cisco Al PODs Extensibility

Because no two organizations' Al requirements are identical, Cisco Al PODs deliver flexible, pre-validated infrastructure to meet diverse Al needs, from large-scale training to edge inferencing, with standardized, secure, and scalable solutions.



Large-scale training clusters



Extensive model optimization



Inferencing at scale and edge

Where **Cisco Al PODs** shine:

Large deployments with dedicated GPU backend networks

Multi-workload architectures and mixed use-cases to run Al models alongside enterprise apps

Inferencing for large Al models and high user concurrency at scale

Al Practitioners / MLOps

IT Infrastructure & Operations



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

Deploy Al with confidence

Cisco CVD, NVIDIA ERA

Fully supported stack including Cisco and 3rd party components

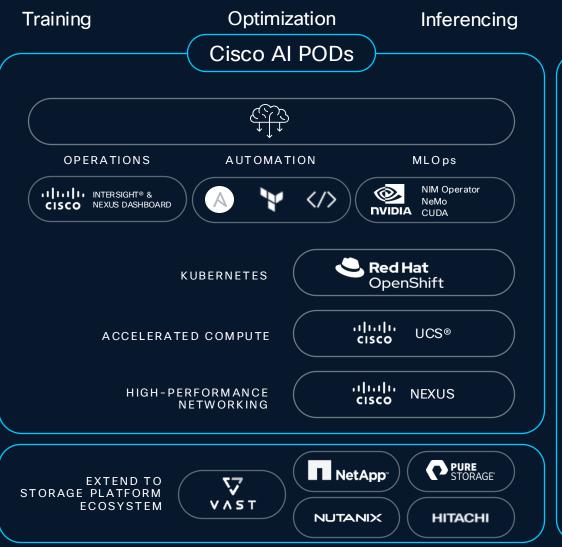
Cisco CX Success Track

Orderable, use case driven Al-Ready infrastructure stacks

Inferencing.
Optimization.
Training.

Incremental, atomiclevel -or- fabric-based cluster scale







ADVANCED

SERVICES

INCLUDED

Cisco Customer Experience

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

Deploy Al with confidence

Cisco CVD, NVIDIA ERA

Fully supported stack including Cisco and 3rd party components

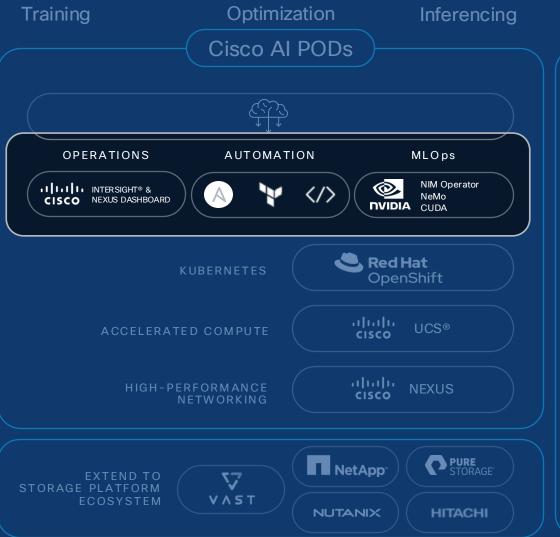
Cisco CX Success Track

Orderable, use case driven Al-Ready infrastructure stacks

Inferencing.
Optimization.
Training.

Incremental, atomiclevel -or- fabric-based cluster scale







ADVANCED

CX Cisco Customer Experience

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

Deploy Al with confidence

Cisco CVD, NVIDIA ERA

Fully supported stack including Cisco and 3rd party components

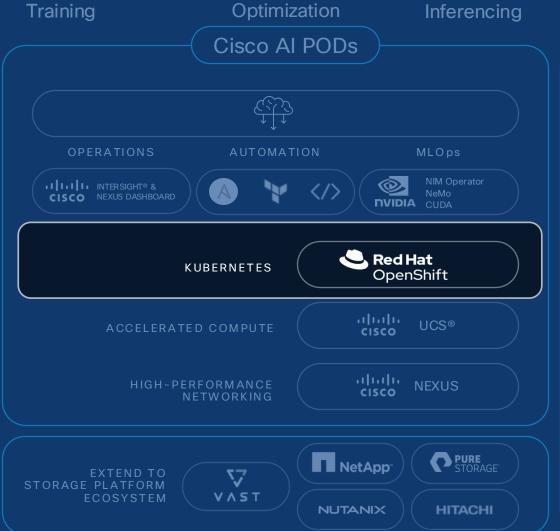
Cisco CX Success Track

Orderable, use case driven Al-Ready infrastructure stacks

Inferencing.
Optimization.
Training.

Incremental, atomiclevel -or- fabric-based cluster scale







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

Deploy Al with confidence

Cisco CVD, NVIDIA ERA

Fully supported stack including Cisco and 3rd party components

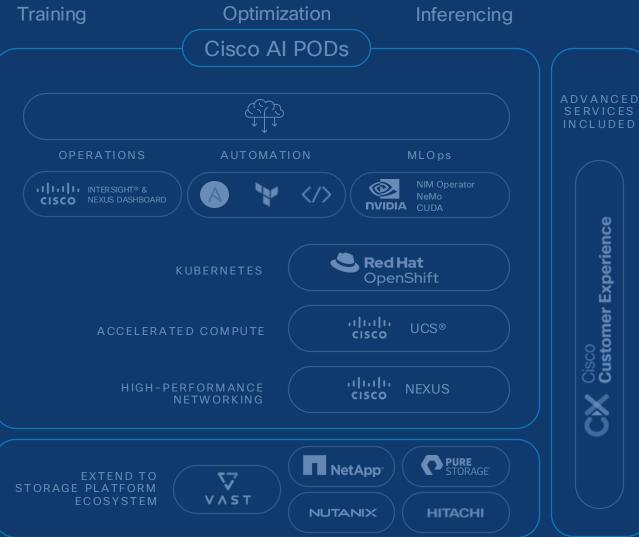
Cisco CX Success Track

Orderable, use case driven Al-Ready infrastructure stacks

Inferencing.
Optimization.
Training.

Incremental, atomiclevel -or- fabric-based cluster scale

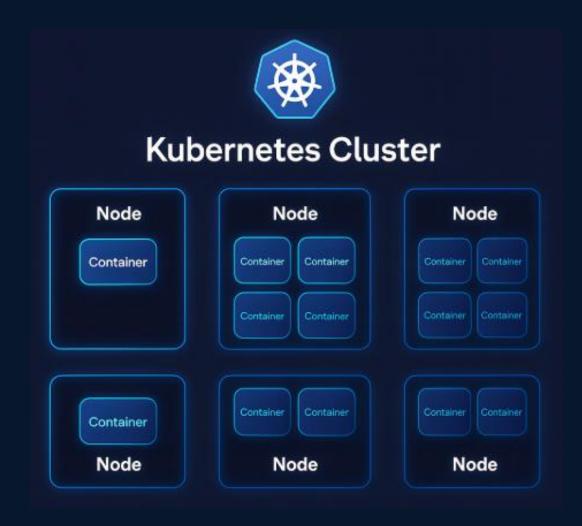




Agenda

- 01 Al Infrastructure Considerations
- 02 Cisco Al PODs
- 03 Al Workload Orchestration
- 04 Operations and Automation
- 05 MLOps
- 06 Cisco Al PODs Extensibility

Kubernetes



- Automation & orchestration CI/CD
- Robust ecosystem of Al tooling
- GPU scheduling
- Portability
- Scalability
- Fault tolerance & reliability

Kubernetes with Red Hat OpenShift

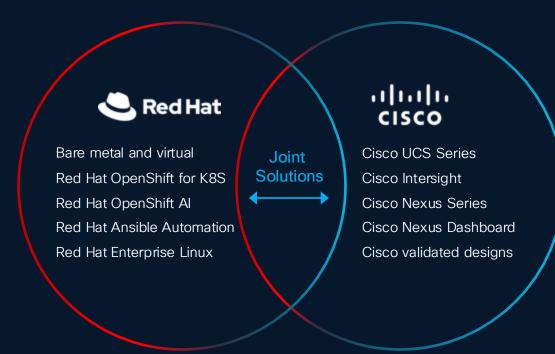
on-premises and cloud

Open Cloud Infrastructure

platform built on open-source innovation

Accelerated Time to value

with turnkey experience and integrated automation



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

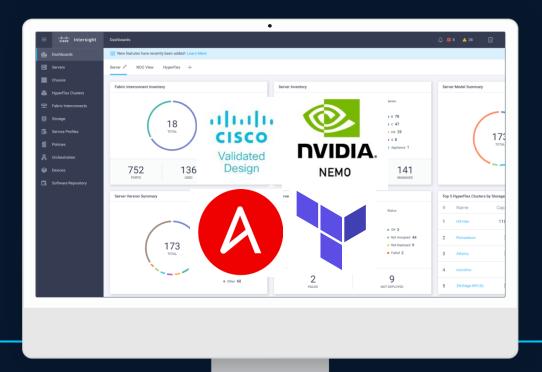
Increased sustainability

Agenda

- 01 Al Infrastructure Considerations
- 02 Cisco Al PODs
- 03 Al Workload Orchestration
- **04** Operations and Automation
- 05 MLOps
- 06 Cisco Al PODs Extensibility

Cisco Intersight

Unified Operating Model



Intersight Dashboard

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

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

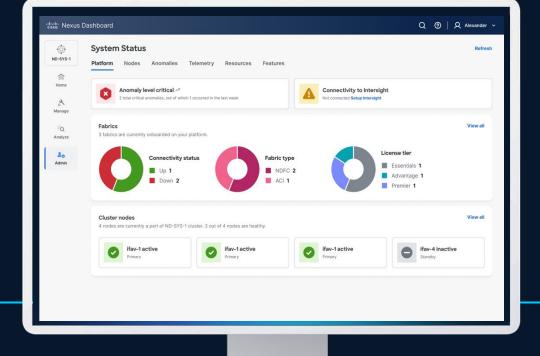
Consistent operational model globally, from DC to edge, at cloud scale

Secure operations with built-in advisories and continuous risk mitigation

Cisco Nexus Dashboard

Simplify Data Center Network Operations

Common policy across NX-OS and ACI fabrics



Nexus Dashboard

Configure, operate and analyze your network from one place across data center networks

Minimize downtime through increased visibility and resolve problems with fix recommendations

Track Power and Cooling by surfacing your networks impact on KWh and CO2 emissions

Accelerate innovation with built in infrastructure-as-code and popular automation tooling integration

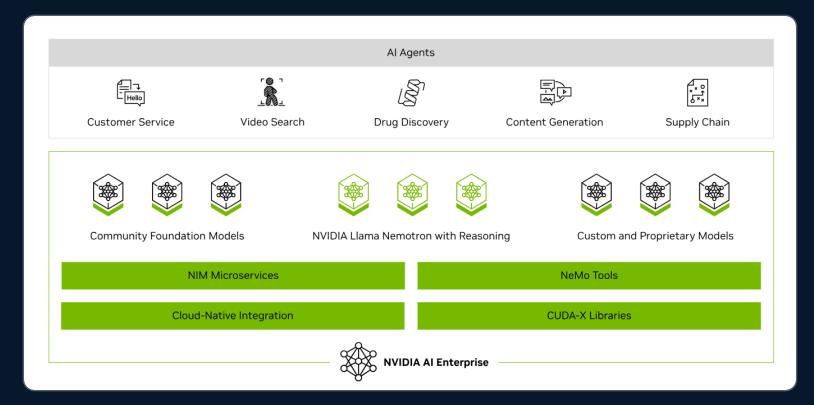
Agenda

- 01 Al Infrastructure Considerations
- 02 Cisco Al PODs
- 03 Al Workload Orchestration
- 04 Operations and Automation
- 05 MLOps
- 06 Cisco Al PODs Extensibility

NVIDIA AI Enterprise

Delivering building blocks for enterprise Al

Production-ready software for agentic Al



The NVIDIA AI Enterprise tools on Cisco AI-PODs provide support for each step in the training, optimization, and deployment of AI agents.



Deploy the latest state-of-the-art Al models

Explore the NVIDIA NIMs catalog of enterprise-ready, performance-optimized models for efficient inference and reasoning.



Build and manage data flywheels with NeMo

Discover powerful, ready-to-use model training, evaluation, and guard railing tools and RAG building blocks for optimizing agentic AI.



Customizable blueprints for your use case

Reference workflows for building fast, high-performance, and secure agentic systems using the latest machine learning best practices.



What will my support experience be?

Support in a multi-vendor solution

"If something breaks, which support team do I call?"

"We don't have the resources to manage multiple product support teams."

"Even minor maintenance changes can cause serious issues."

"Our support experience with some vendors has been inconsistent."

Cisco Solution Support

One Service, Broad Coverage

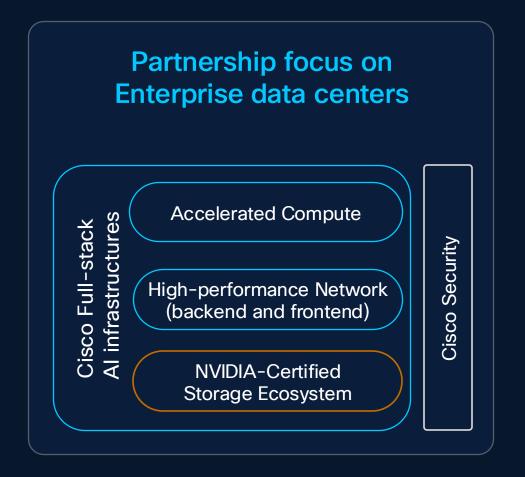
Service Features	Cisco Solution Support
Global 24x7 product-level technical support	•
24-hour access to Cisco® online resources	•
Hardware replacement (2- and 4-hour, next business day)	•
Network management / operating system software updates and upgrades	•
Proactive diagnostics and immediate alerts on devices through Cisco Smart Call Home	•
Web-based user community for self-service support of smart capabilities	•
Cisco software application support	•
Primary point of contact with solution-level expertise	•
Accountability for issue resolution, no matter where it resides	•
Coordination between Cisco TAC and solution partner product support teams	•
Case management from first call to resolution	•

Agenda

- 01 Al Infrastructure Considerations
- 02 Cisco Al PODs
- 03 Al Workload Orchestration
- 04 Operations and Automation
- 05 MLOps
- 06 Cisco Al PODs Extensibility

Expanded partnership to accelerate Al adoption in the enterprise





Cisco is now included in NVIDIA SpectrumTM-X Ecosystem, and a validated partner for NVIDIA Enterprise Reference Architectures



Cisco Secure Al Factory with NVIDIA (Secure Al Factory) enables enterprises to accelerate the adoption of Al use cases



Jointly deliver Secure Al Factory with NVIDIA SpectrumTM–X Silicon integrated Cisco Silicon One® switches.

Cisco

Security

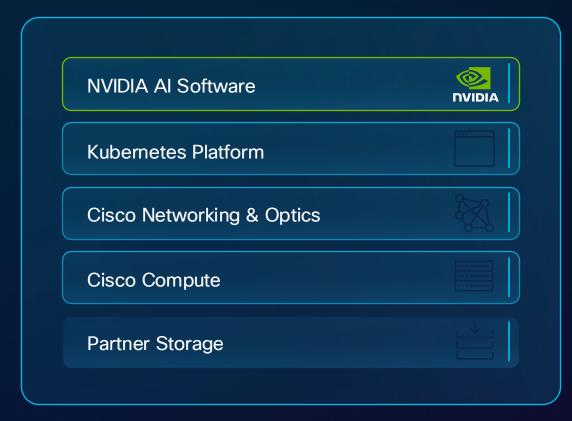
Cisco Secure Al Factory with NVIDIA

What is it?

Reference architecture

Validated solutions and turnkey offerings

Differentiated with Security and Observability



Security-first architecture enables safe Enterprise Al



Security at all layers of the stack

Securing the Applications

Cisco Al Defense—Robust testing and runtime security of LLMs and generative Al applications, integrated with NVIDIA Al.

Securing the Workloads

Cisco Hypershield—Protection against adversary lateral movement and proactive vulnerability mitigation without the need for patching, all from a single management interface, integrated with NVIDIA AI.

Integration with NVIDIA Bluefield-3's DOCA AppShield for intrusion detection in Al-focused VMs and containers.

Future

Securing the Infrastructure

Cisco Hybrid Mesh Firewall—Unified security management and consistent and pervasive policy across multiple enforcement points.

Cisco Isovalent: Enhanced visibility into cloud native interactions, enabling smooth policy definition and enforcement across software defined networks.

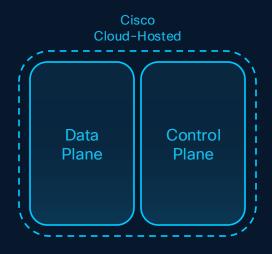
To include management of NVIDIA BlueField®-3 DPUs for enabling Al Cluster perimeter firewalls.

Future

Cisco Secure Al Factory - Layered Defense



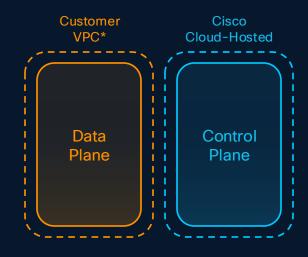
Al Defense Deployment options for every situation



SaaS

Fully hosted and managed in the cloud

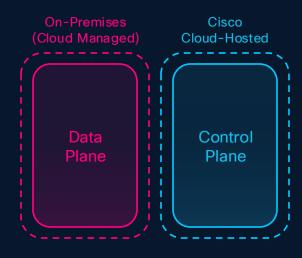
Best for customers looking for a simple, flexible deployment with zero infrastructure to manage



VPC

Virtual private cloud environments with a cloud-hosted control plane

Best for customers looking to balance data control and compliance with cloud scalability



On-Premises

Combines Cisco UCS hardware with a cloud-hosted control plane

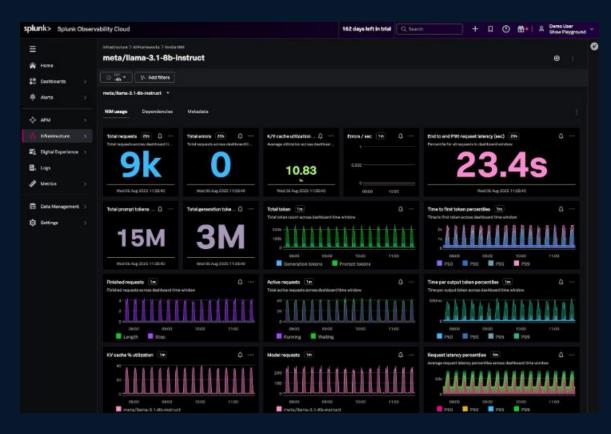
Best for customers that want to manage Al workloads themselves rather than relying on hyperscalers

illiili cisco



= Better Together

Splunk acts as the **observability + analytics layer** for Cisco Al PODs. While Intersight manages lifecycle, Splunk provides **cross-domain correlation**, **security**, **and operational intelligence**



- Al Infrastructure Observability
- AI/ML Cluster Monitoring
- Security & Compliance
- Capacity & Planning
- Networking & Storage Visibility



VAST Data on Cisco Al PODs

Simplified Full Stack Orderability

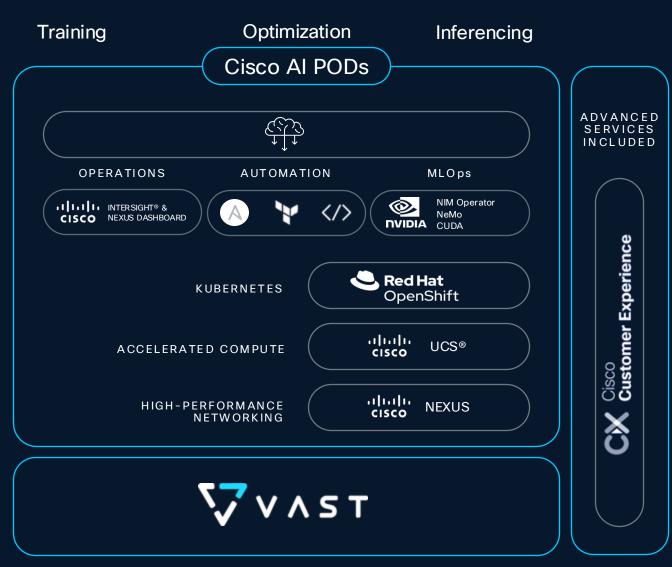


Faster time to value with pre-configured bundles

Deploy AI with confidence

Orderable, validated Al-ready infrastructure stacks

Al Advisor tool for configuration guidance



The V A 5 T Al Operating System



Data**Engine**

Scalable, Event-Driven Computing Triggers, Functions, Containers



Sync**Engine**Data Router Migrations





Agent**Engine**Agent Tools & Orchestration
Agents, Tools, Observability



Data**Store**

Universal Storage & Data Protection Infrastructure NFS, SMB, S3, NVMe/TCP, VAAI, CSI, vLLM



Data**Base**

Transactional, Analytical, & Vector Database Kafka, Python, Parquet, SQL



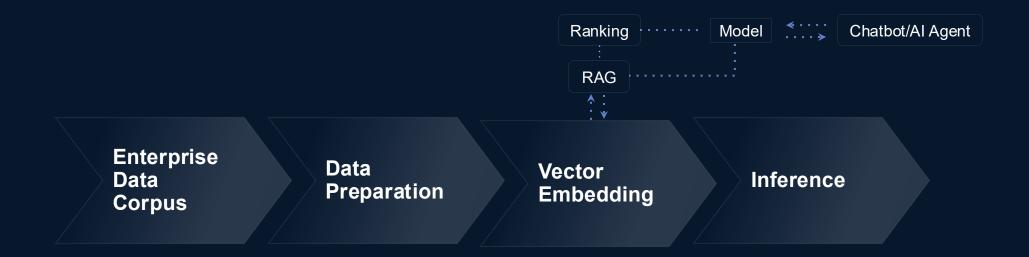
Data**Space**

Globally-Distributed Data Computing



Understanding Enterprise Al Pipelines and RAG

Cisco, NVIDIA, and VAST working as a single solution



Summary - Cisco Al PODs

Deploy Al with confidence

Cisco CVD, NVIDIA ERA

Fully supported stack including Cisco and 3rd party components

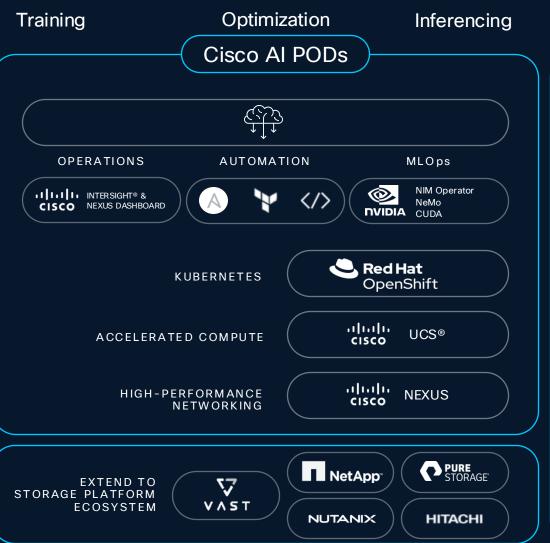
Cisco CX Success Track

Orderable, use case driven Al-Ready infrastructure stacks

Inferencing.
Optimization.
Training.

Incremental, atomiclevel -or- fabric-based cluster scale







Cisco Differentiation



The Security

Security-first architecture enables safe enterprise Al



The Infrastructure

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



The Assurance

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

Thank you

ıı|ıı|ıı CISCO Making Al work for you.



