



Modernizing the DC to Support Al

Building, Observing, and Securing the Next Gen Infrastructure

Cloud and Al Infrastructure Specialist SE

Jason Powell Eliot Ngwa

Cisco powers how people and technology work together across the physical and digital worlds



Al-Ready Data Centers



Secure Global Connectivity

Digital Resilience





Every organization's Al approach and needs are different

Build the Model | Training

Optimize the Model | Fine-tuning and RAG

Use the Model | Inferencing

















And applications have evolved

From simple and deterministic to complex and adaptive

Pre-Al Period

- Deterministic logic
- · Predictable behavior
- Static query / response

Large Language Models (LLMs)

- Natural language input / output
- No workflows / ~1 span

Retrieval-Augmented Generation (RAG) with LLMs

- Hybrid and multi-cloud environments / Microservices
- Linear workflows / ~10s of spans

Agentic Applications

- Autonomous task completion
- · Unpredictable behavior
- Multi-step actions

Past Present Future

The emphasis in AI infrastructure is shifting from training massive models to optimizing inference

Brad Lightcap



"The industry is realizing that inference—not just training—is where the economic value of Al is created."

Sarah Guo



"We're seeing a fundamental forking of Al scaling: the old model was 'bigger is better' for training, but now inference efficiency is becoming the real battleground."

Jonathan Ross

"The shift to inference marks a new phase of Al's power—moving from potential to real-world impact at scale."

Aidan Gomez

scohere

"Training pushed Al forward, but inference is where the real-world performance gaps are being solved now."

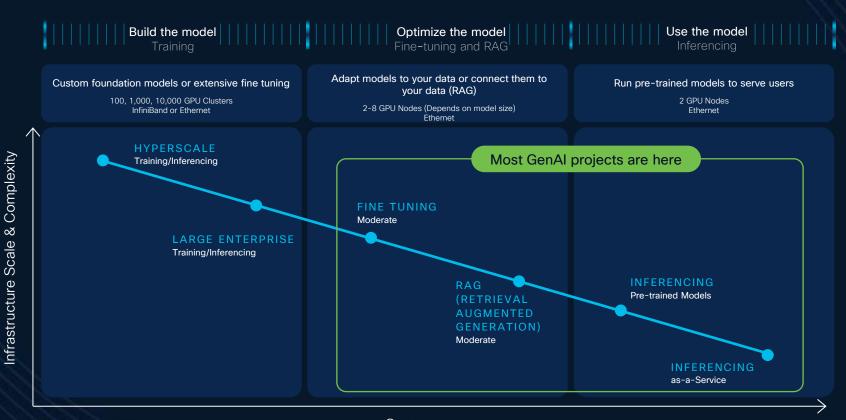
Arvind Jain glean

"For enterprise AI, it's all about inference—we're not training new models, we're optimizing the ones we already have to deliver results instantly."

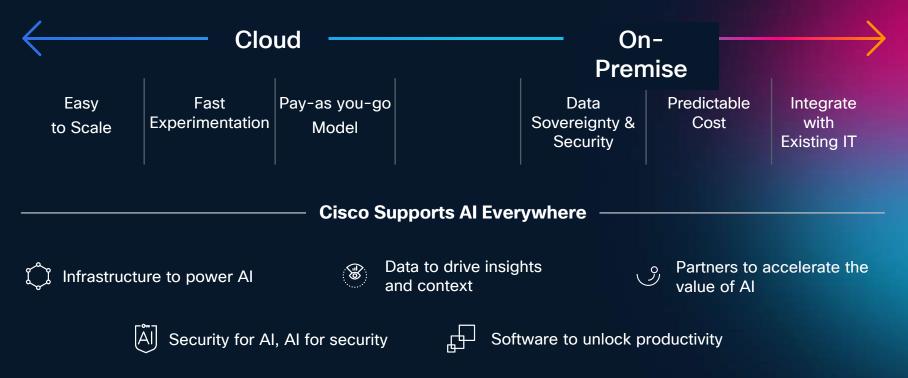
Jonathan Ross

"We're shifting from 'Can Al work?' to 'How can we scale it effectively?' it's all about inference now."

Al Workload Spectrum: Training, Fine-Tuning, or Inference



Where It Runs: Cloud vs On-Premise



Al is forcing a rethink of the tech stack

>80%

not Al ready - need hardware acceleration

20%

of CIOs have a mature data platform strategy

>48%

of network, compute, and storage infrastructure is 10+ years old

Network

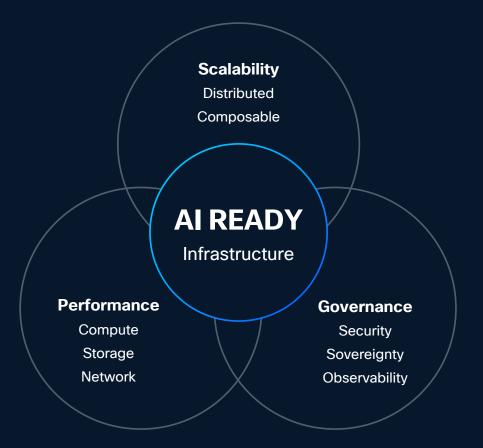
speed, fabric designs, and performance

>48%

investigating VMware alternatives

40%

increase in security breaches year over year due to old infrastructure



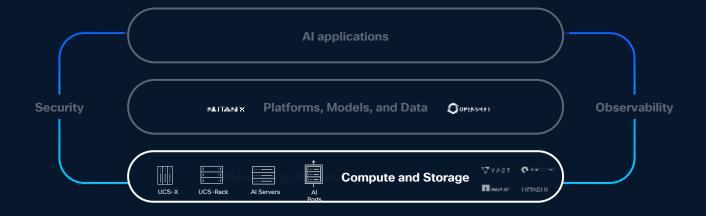
Sources: IDC; Gartner; Forrester; TechRadar; CIO.com, Accenture Tech Vision (2025)



Compute for Al

Unified Computing Systems

Deliver resources in any location from the cloud



Cloud Managed

Al-ready infrastructure

Train, fine tune and inference on accelerated servers and integrated full stack systems

Simplify at Scale

Unified infrastructure operations for faster time to value and easier lifecycle management

Hybrid Multcloud

Modernize with validated converged and hyperconverged platforms that support distributed applications

Cisco Al Compute Portfolio

Unified approach to accelerated Al compute

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



Build the model | Training

Optimize the model | Fine-tuning and RAG

Use the model | Inferencing

Cisco Al Compute Portfolio

Dense GPU Al Servers

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



Cisco UCS® C885A M8 Rack Server

NVIDIA HGX platform with 8 NVIDIA H100 NVL/H200 NVL and M300X GPUs 2 AMD 4th Gen/5th Gen EPYC processors For fine tuning and scaleable inference use cases



Cisco UCS® C845A M8 Rack Server

NVIDIA MGX platform with 2/4/6/8 NVIDIA H100 NVL/H200 NVL/L40S GPUs 2 AMD 5th Gen EPYC processors

Bringing High-Density GPU Servers to the Cisco UCS Family

Built for LLM training, deep learning, fine-tuning, and HPC

UCS Accelerated | UCS C880A M8



AVAILABLE NOW



2 CPUs

Intel Xeon 6th Gen Scalable Processor

NVIDIA HGX with 8 GPUs

NVIDIA B300 with NVL8 Air Cooled

Networ

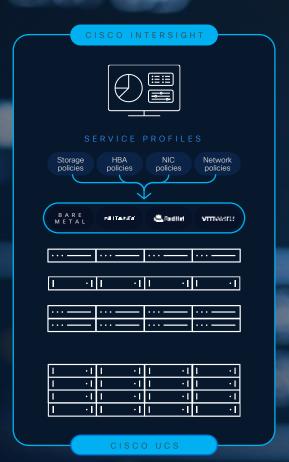
- (8) NVIDIA ConnectX-8 GPU Board Integrated (E-W)
- (2) NVIDIA BF3 B3220, NVIDIA BF3240, NVIDIA ConnectX-7 (N-S)

Power

(12) 50V 3200W (N+N redundancy)

Computing that's elastic, responsive and intelligent

Delivering industry differentiation with Cisco's SaaS-managed, fabric-based data center



Elastic

Shape infrastructure to any workload, depending on business needs

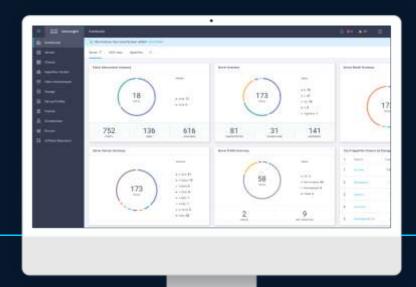
Intelligent

Connect and scale from the data center to the edge with 360° visibility

Responsive

Adapt to the future with unparalleled simplicity, performance and resiliency

Work smarter and faster with a simplified, unified operating model



Cisco Intersight

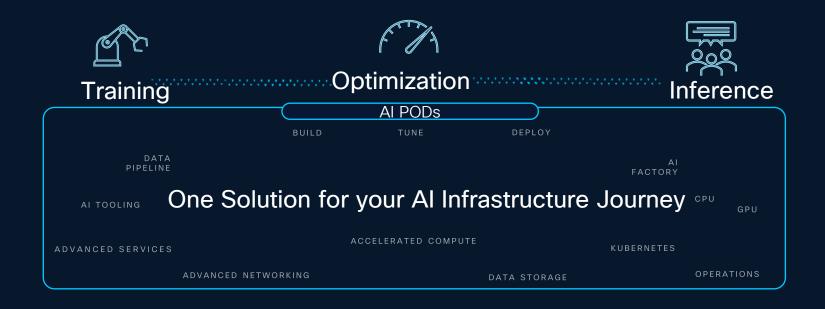
See your global on-premises, cloud, and edge environments **Connect** your infrastructure operations across compute and storage

Secure operations with built-in advisories and continuous risk mitigation

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

Cisco Al Pods

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



Cisco Al PODs

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

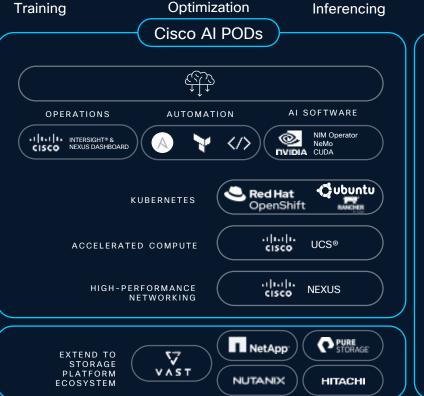
Cisco CX Success Track

Orderable, use case driven Al-ready infrastructure stacks

Inferencing.
Optimization.
Training.

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







ıı|ıı|ıı cısco

Cisco Validated Designs - Reducing Risk

What a Validated Design Guide provides









Reliability

CVDs are extensively tested. You can confidently set performance expectations when you deeply your solution.



Using a CVD reduces both the risk that products won't work together and the risk that they won't perform as promised.

Comprehensiveness

CVDs provide everything from system designs to configuration instructions to a bill of materials (BOM).

24-hour support

Because CVD solutions are guaranteed to work as specified, we offer 24-hour support options for any issues that might arise.

Easily Deploy New Systems With Expert Guidance Consistent, Successful Technology Deployments Addressing Business Initiatives

CVDs to simplify and automate AI infrastructure



CVD blueprint for AI networks



Best performing AI/ML networks. focus on application performance



Dynamic congestion avoidance



Automation for day-2 operations



Intelligent buffer, low latency, telemetry and RoCFv2



One IP network for both front-end and back-end



Validated designs for network and ecosystem partners





NVIDIA AI

NUTANIX



Red Hat

Red Hat

OpenShift Al

GPT-in-a-box Gen-Al with on Nutanix Cloudera Data Platform



FlexPod

intel.



NEW

CVD playbooks for deploying common Al models



Large language models (GPT3, BERT, T5)



Computer vision models (ResNet, EfficientNet, YOLO)



Generative models (GANs. VAEs)







Developer Cloud

As your needs become more specific, so do our recommendations



Your Trusted Al Advisor

Helping you optimize your Al infrastructure

Workload sizing



GPUs recommendations



Multiple LLM model support



End-to-end guidance



Bot assistance and BOM generation

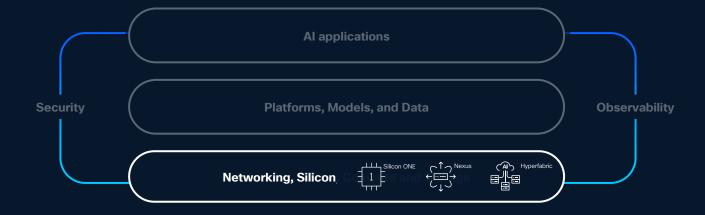


Expanded model selection



Data Center Networking

Industry leading data center fabrics coupled with operational simplicity



Fabric Options

Choice of Fabric

Cisco Nexus and Hyperfabric connect and protect the most demanding workloads, powered by Silicon One

Simplified Operations

Choose on-premise or cloud managed operational model that delivers operational insights, efficiency and sustainability

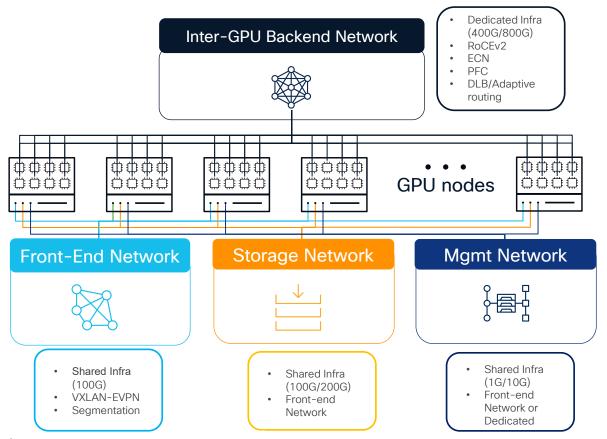
Validated Designs

Design, deploy and operate with codeveloped reference architectures and best practice from Cisco and NVIDIA



Datacenter Networking for Al

Multiple networks for Al infrastructure



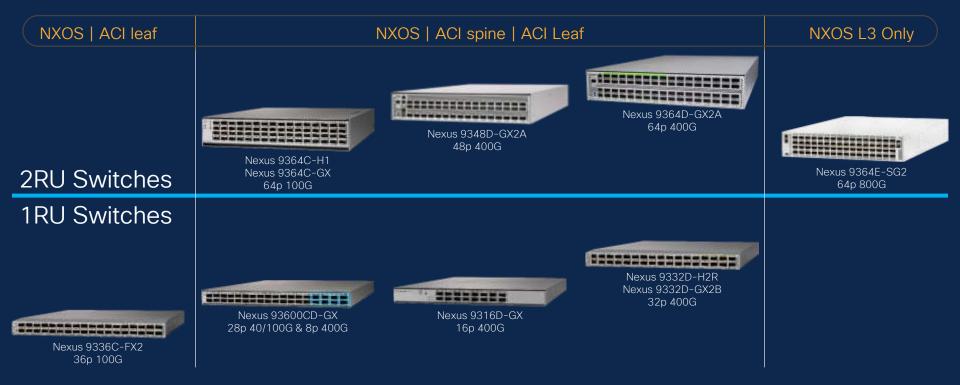
Al Infrastructure with Cisco Nexus 9000 Switches



Cisco Data Center Networking Blueprint for Al/ML Applications

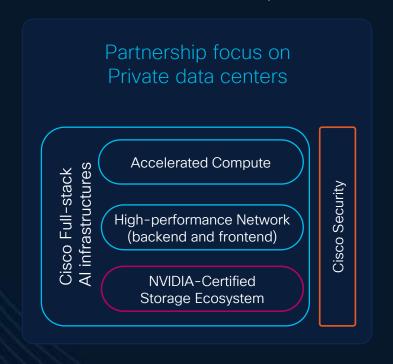


Cisco Nexus 9300 Series - 100G/400G/800G Fixed Switches



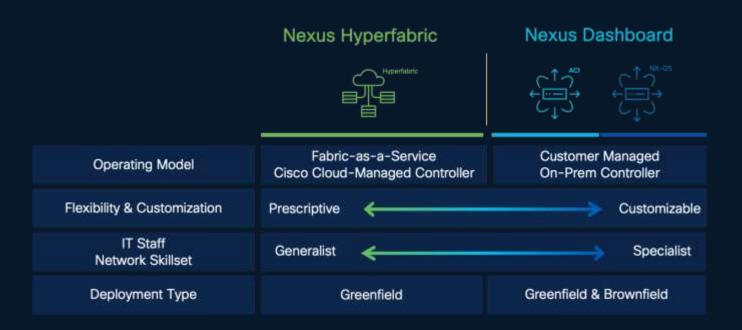


Cisco and NVIDIA expand partnership to accelerate Al adoption in the enterprise





Managing Al Infrastructure



Greenfield: new fabrics not being managed by Nexus Dashboard



Cisco Secure Al Factory with NVIDIA

