CISCO Engage | Tech Day

# *Accelerate Your AI-Ready Data Center*

*With Cisco CX Services*

CISCO

*Vamsidhar A*
*Business Development Manager - AI Factory*
*vamsa@cisco.com*

*Dec 17th  2025*

# What We'll Cover

## Agenda

### The AI Horizon

- Are we Ready for AI

- Cisco AI Ready Data Center

- Cisco's AI Services portfolio

- AI Ready Enterprises

- Future of AI

CISCO

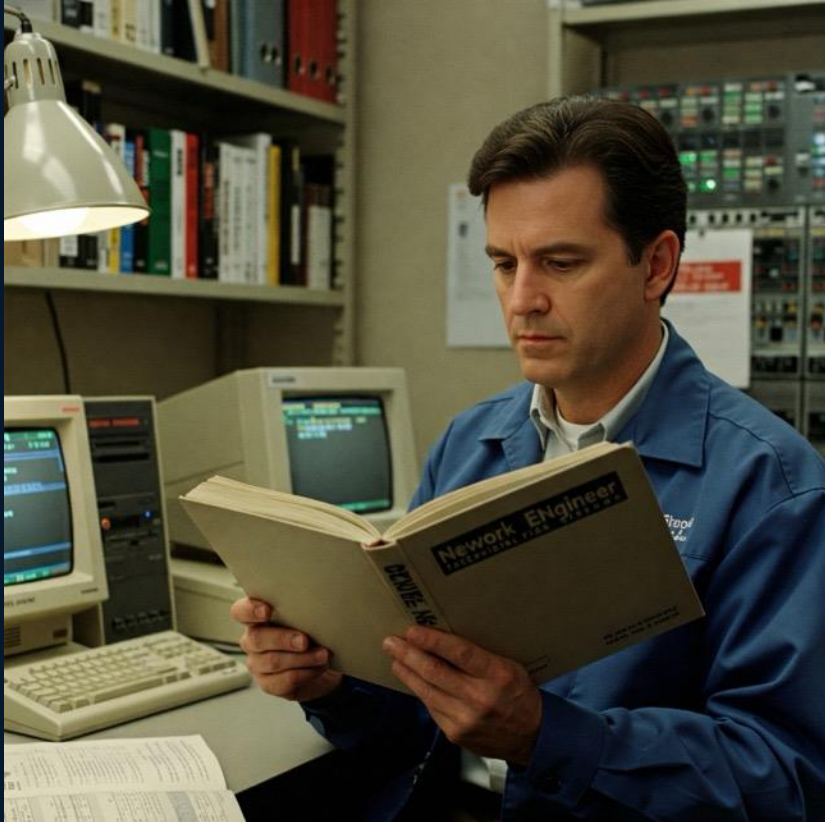# Are we Ready for Future !

# World has shifted
# 3 Years ago!

How Powerful is your Brain !

* Ref 30th Nov 2022

# In the 90s, before the internet, we had to read books

# Fundamental Shift on How we work

*With internet and search engines*



image generated by AI

What we thought AI would look like

Will Robots Pamper us !

AI: The Next Fundamental Transformation in How We Live and Work

# Are we Ready for AI

# AI is Everywhere

## $15.7T
*Potential contribution to global economy by 2030*

## $300B
*Global spending on AI by 2026*

## 75%
*Of large enterprises will rely on AI-infused processes by 2026*

*Healthcare and Life Sciences*
Diagnosis
Drug discovery
Personalized medicine

*Financial Services*
Fraud detection
Risk assessment
Trading

*Retail*
Personalization
Inventory optimization
Virtual agents

*Manufacturing*
Predictive maintenance
Quality control
Demand forecasting

*Agriculture*
Yield optimization
Automated irrigation
Pest prediction & prevention

*Transportation*
Route optimization
Autonomous vehicles
Predictive maintenance

*Energy*
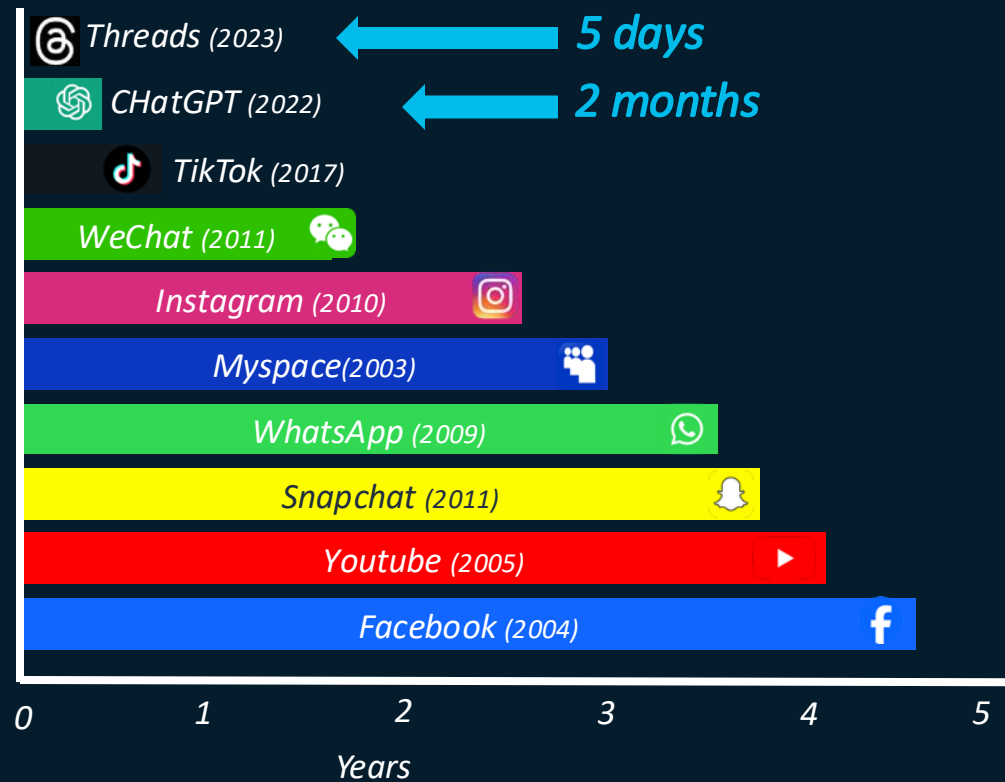Distribution optimization
Fault prediction
Demand forecasting

*Public Sector*
Smart cities
Security
Services improvement

*Sources: PWC, IDC*

#CiscoConnect

# Pace of Adoption

## Time to Reach 100M users



Threads (2023) ← 5 days

CHatGPT (2022) ← 2 months

TikTok (2017)

WeChat (2011)

Instagram (2010)

Myspace (2003)

WhatsApp (2009)

Snapchat (2011)

Youtube (2005)

Facebook (2004)

0   1   2   3   4   5

Years

visualcapitalist.com

**86%** of companies not fully ready to integrate AI into their businesses.

#CiscoConnect

# Every organization's needs are different



Build the model
Training

Optimize the model
Fine-tuning and RAG

Use the model
Inferencing

| Text | Code | Image | Speech | Video | 3D | Other |
|------|------|-------|--------|-------|-----|-------|
| Marketing (content) | | | | | | |
| Sales (email) | | | | | | Gaming |
| Support (chat/email) | Coding generation | Image generation | | | | Music |
| General Writing | Coding documentation | Consumer / Social | | | | Audio |
| Note taking | Web app builders | Media / Advertising | | | | Healthcare / Biology / Chemistry |
| Other | Text to SQL | Design | Voice synthesis | Video editing/ generation | 3D modeling | |

\* **Multimodal LLMs**

#CiscoConnect

# AI and Infrastructure



**Data Preparation**

Preparing structured or unstructured data to create a training data set for the model

High storage requirement for ETL, data cleansing and optimized for AI retrieval

**Training**

A selected model learns from the training data set and builds relationships

Compute intensive often with GPU acceleration and high-speed low latency network

**Inference**

When prompted the model interprets new, unseen data and creates a response

Lower compute requirements, GPU acceleration and network demands.

Prompt

Response

User

**IO Intensive**

**Compute Intensive**

**Latency Sensitive**

*Platform | DevOps | SecOps | Infrastructure*

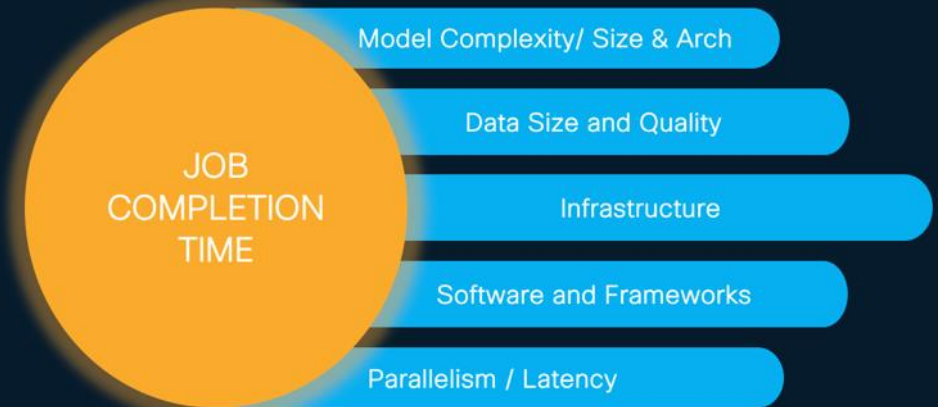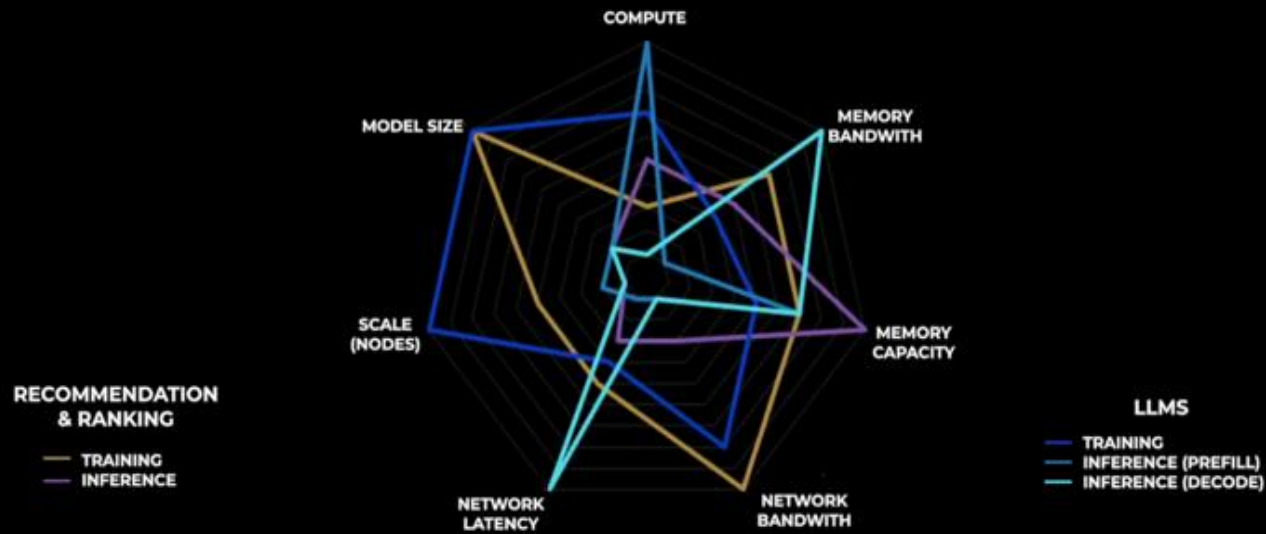## Model performance tightly coupled to infrastructure

# AI/ML Requirements
## Training vs Inference

# AI Ready Data Center

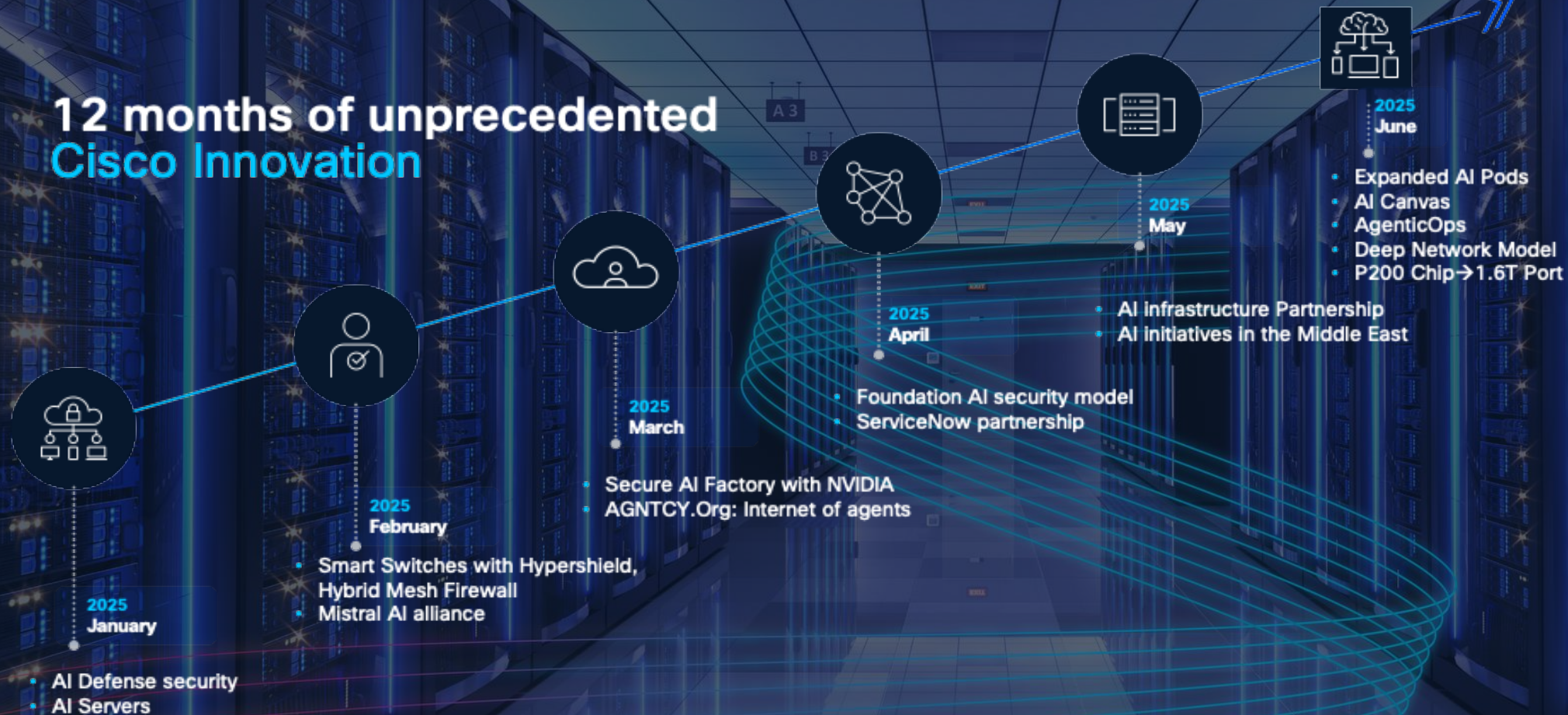The world is experiencing the *largest expansion of data centers* in history

*Cisco is foundational* to the world's data center buildouts

Enterprises | Neoclouds | Service Providers | Hyperscalers

# The AI-ready data center

**12 months of unprecedented Cisco Innovation**

**2025 January**
- AI Defense security
- AI Servers

**2025 February**
- Smart Switches with Hypershield, Hybrid Mesh Firewall
- Mistral AI alliance

**2025 March**
- Secure AI Factory with NVIDIA
- AGNTCY.Org: Internet of agents

**2025 April**
- Foundation AI security model
- ServiceNow partnership

**2025 May**
- AI infrastructure Partnership
- AI initiatives in the Middle East

**2025 June**
- Expanded AI Pods
- AI Canvas
- AgenticOps
- Deep Network Model
- P200 Chip→1.6T Port

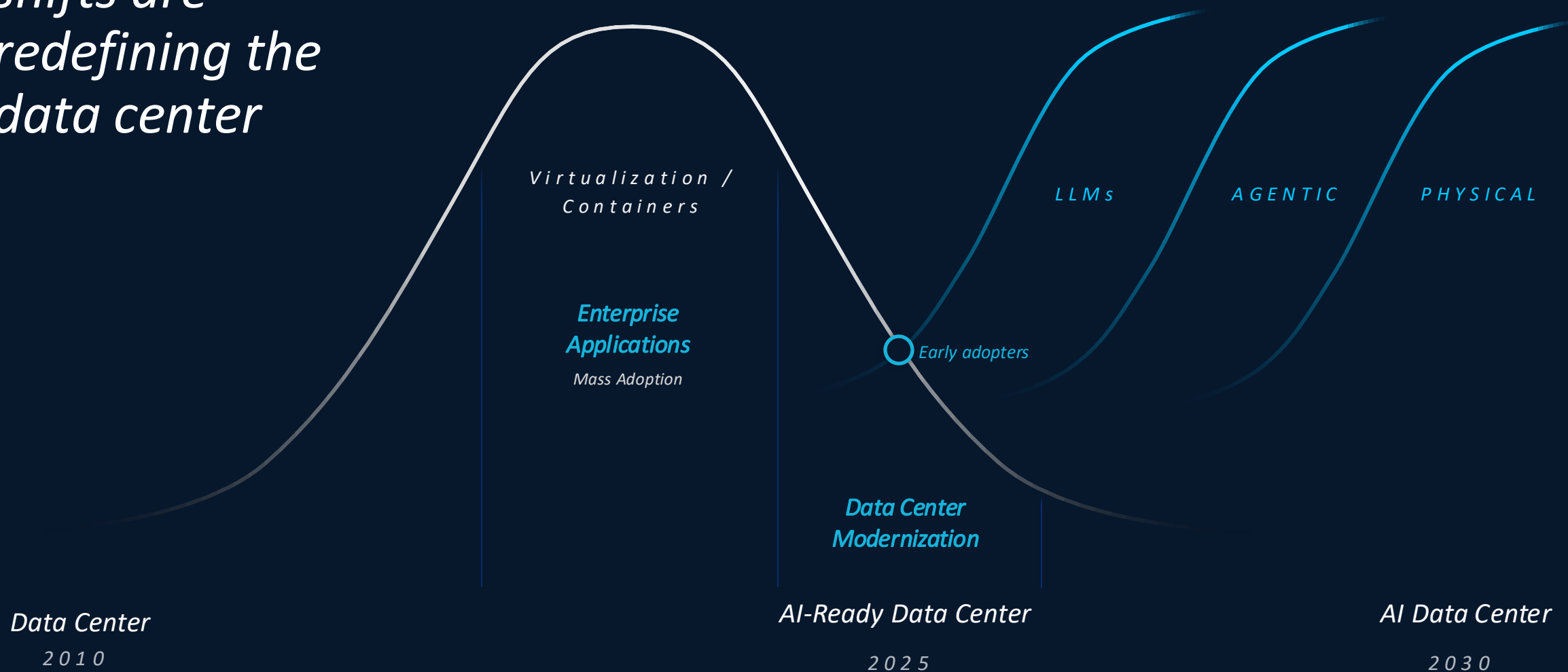*AI is network bound*

*Private data center "re-acceleration"*

*Common foundation for AI safety*

*Hyper-distributed security*

*What we know*

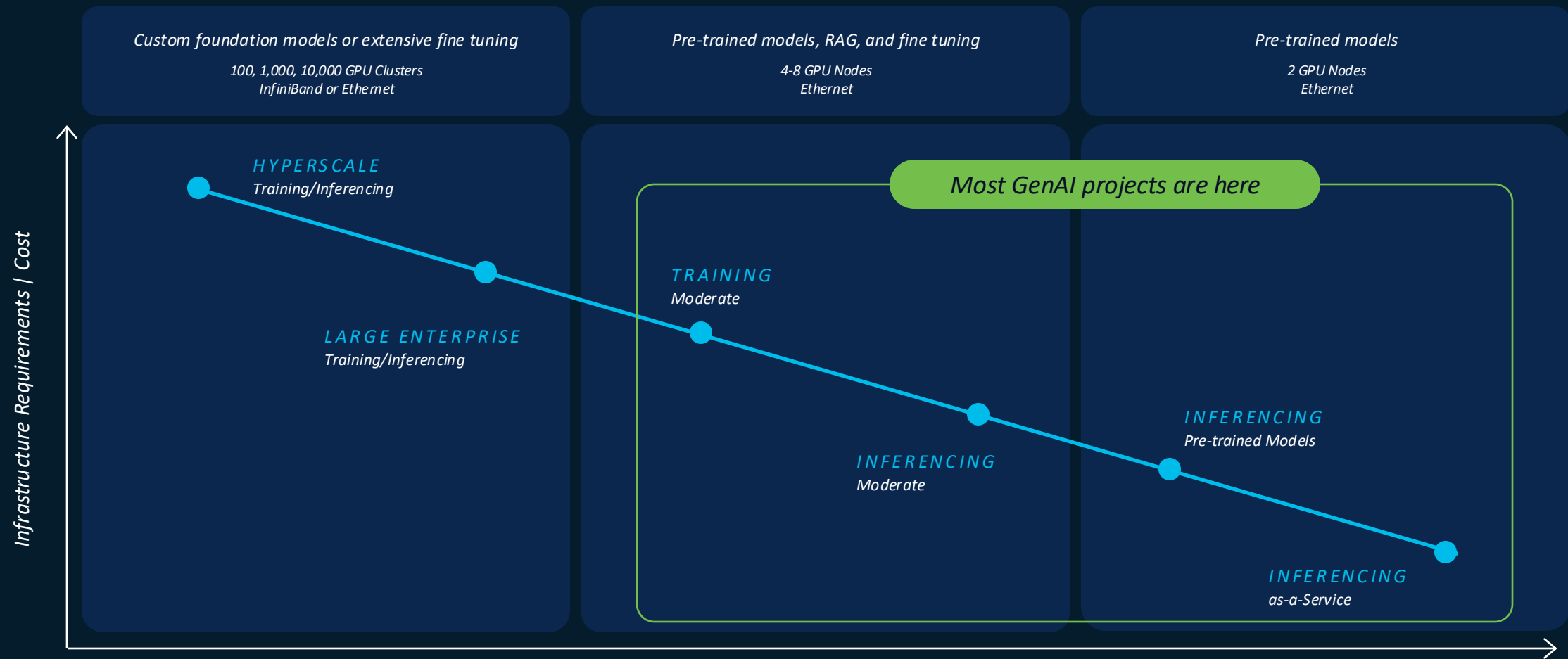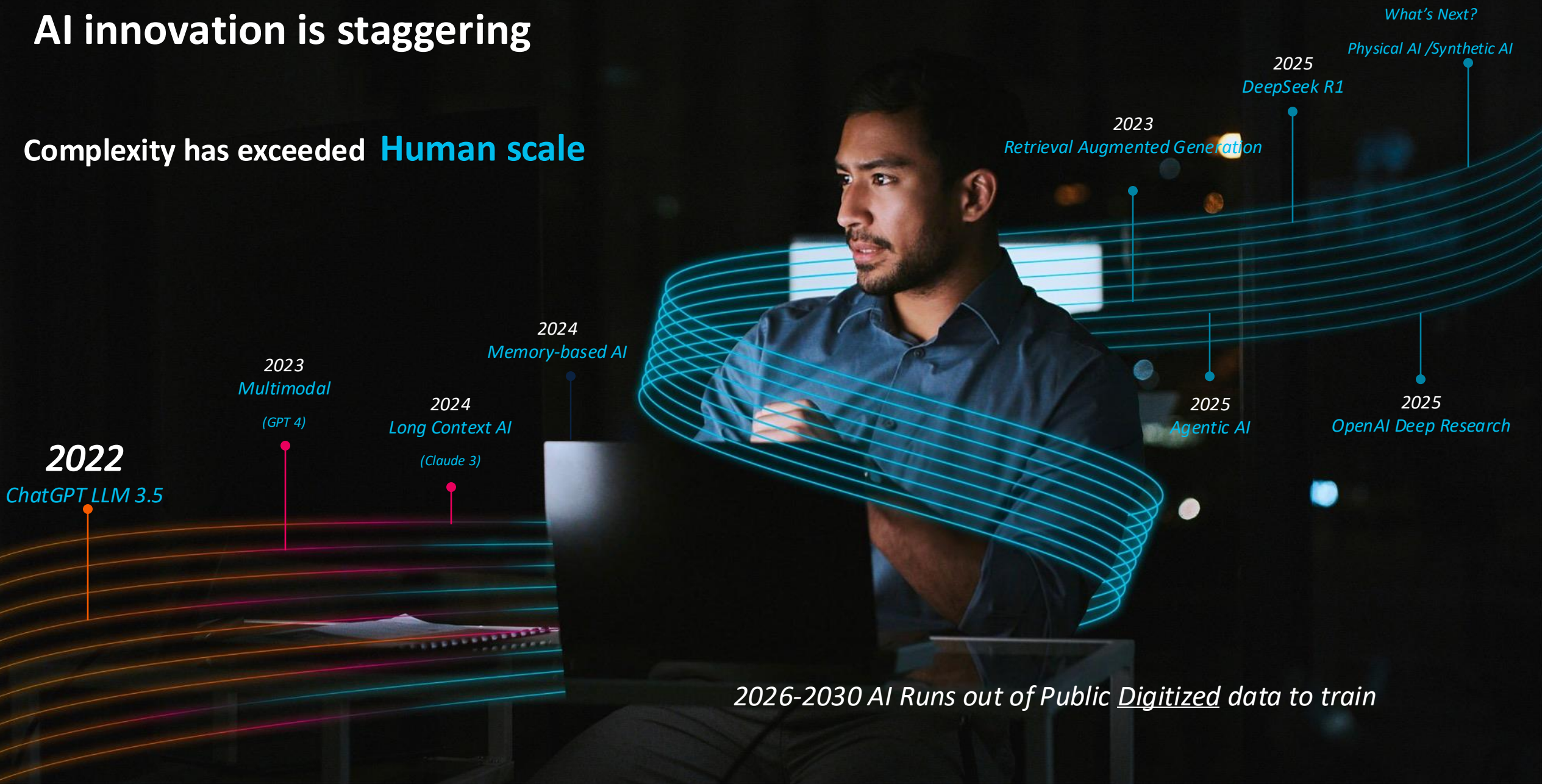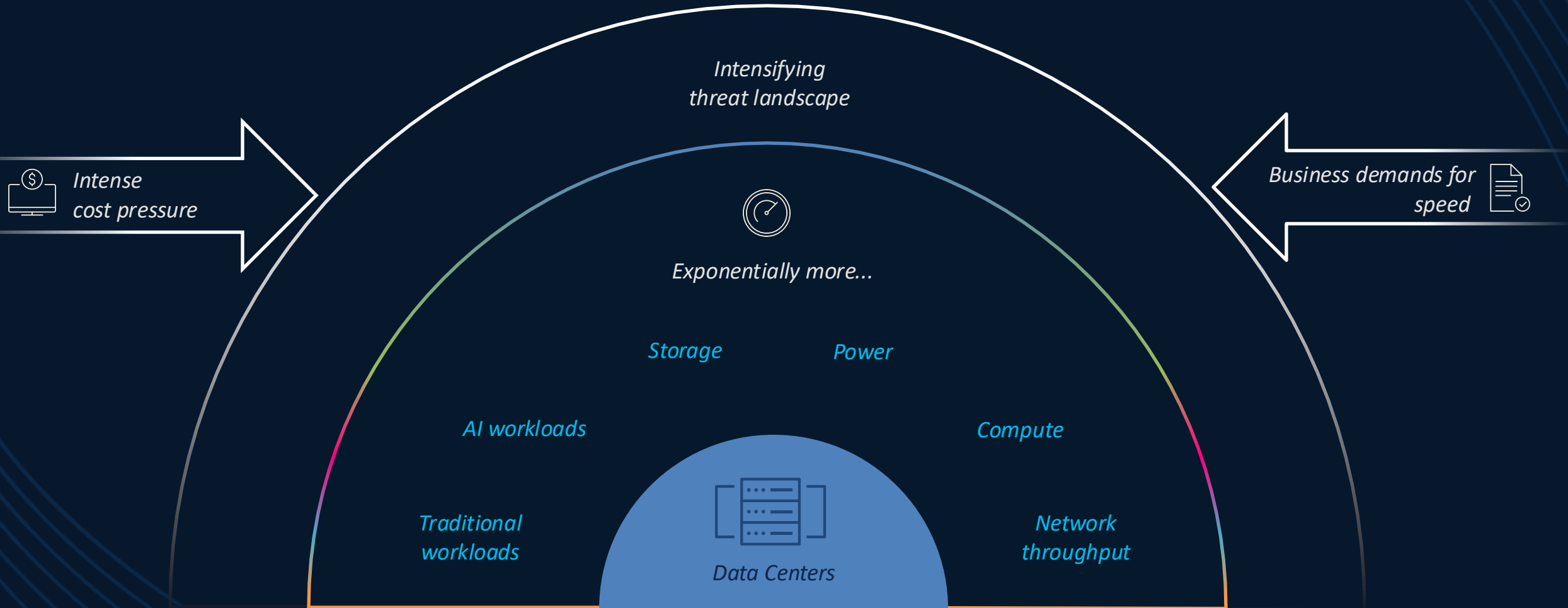# Architectural shifts are redefining the data center

8B to 80B

*Virtualization / Containers*

*Enterprise Applications*

Mass Adoption

○ Early adopters

*LLMs*    *AGENTIC*    *PHYSICAL*

*Data Center Modernization*

*Data Center*

*2010*

*AI-Ready Data Center*

*2025*

*AI Data Center*

*2030*

CISCO

# AI Use Case Spectrum

| | | |
|---|---|---|
| *Custom foundation models or extensive fine tuning* | *Pre-trained models, RAG, and fine tuning* | *Pre-trained models* |
| *100, 1,000, 10,000 GPU Clusters*<br>*InfiniBand or Ethernet* | *4-8 GPU Nodes*<br>*Ethernet* | *2 GPU Nodes*<br>*Ethernet* |

**Most GenAI projects are here**

*Infrastructure Requirements | Cost*

**HYPERSCALE**
Training/Inferencing

**LARGE ENTERPRISE**
Training/Inferencing

**TRAINING**
Moderate

**INFERENCING**
Moderate

**INFERENCING**
Pre-trained Models

**INFERENCING**
as-a-Service

The pace of
AI innovation is staggering

Complexity has exceeded **Human scale**

*What's Next?*
*Physical AI /Synthetic AI*

*2025*
*DeepSeek R1*

*2023*
*Retrieval Augmented Generation*

*2024*
*Memory-based AI*

*2025*
*Agentic AI*

*2025*
*OpenAI Deep Research*

*2023*
*Multimodal*
*(GPT 4)*

*2024*
*Long Context AI*
*(Claude 3)*

**2022**
*ChatGPT LLM 3.5*

*2026-2030 AI Runs out of Public Digitized data to train*

# Reimagine the data center for the AI era
While managing cost pressures, intensifying threats, and business demands

Intensifying
threat landscape

Intense
cost pressure

Business demands for
speed

Exponentially more...

Storage          Power

AI workloads                    Compute

Traditional
workloads

Network
throughput

Data Centers

# AI Infrastructure Stack
## Compute, Network, Security, Observability

| Data Infrastructure | Prepare Data | Train Model | Evaluate Model | Deploy, Inference & Improve |

Network, Compute: Infrastructure Optimized for GenAI Compute

Security: Data Integrity, Attack Prevention, Policy Management

Observability: Model and Infra Monitoring

**A Multi-architecture play**

# Cisco's AI Ready DC

AI-Ready Data Centers

SECURE GLOBAL CONNECTIVITY

Digital Resilience

< < < < < < Accelerated by Cisco AI > > > > > >

# Cisco's AI Ready DC Portfolio

**AI-Ready Data Center**

*Digital Resilience*

## Distributed AI infrastructure

- UCS Accelerated Compute
- Nexus 9K Series
- Cisco Optics
- Cisco Silicon One
- Cisco 8K series (hyperscalers)
- Storage (partners)

## Full stack systems*

- Nexus HyperFabric AI
- AI PODs
- Nutanix GPT
- NVIDIA & AMD GPU
- Virtualization (partners)
- AI SW Tools (partners)

## AI infrastructure operations

- Intersight
- Nexus Dashboard
- Nexus HyperFabric
- NVIDIA AI Enterprise & NIMs
- Splunk
- AppDynamics
- ThousandEyes
- Accedian(PCA)

## Cloud protection

- Cloud Protection Suite: Hypershield, Secure Workload, Multicloud Defense
- AI Defense & Perimeter Firewall

CISCO

**Build the model**
Training

**Optimize the model**
Fine-tuning and RAG

**Use the model**
Inferencing

| Security | Enterprise Applications | Gen AI | Observability |
| --- | --- | --- | --- |
| | AI Frameworks and Tooling | | |
| | Application Platforms | | |
| | Compute and Data | | |
| | Networking & Storage | | |

Data center          Edge          Colocation          Public cloud

# Cisco AI-Ready Data Center
## Ecosystem Stack

**Operations**

**Security**

Firewall | Secure workload | Robust Intelligence | Splunk

splunk>

**AI frameworks**
NVIDIA NGC
intel Developer Cloud

**AI management tools**
NVIDIA
PyTorch

**Virtualization and Kubernetes**
OPENSHIFT

**Infrastructure management**
Nexus Dashboard
Intersight

**AI Infrastructure**
NVIDIA    intel    AMD    COHESITY
FlashStack    FlexPod    VAST    NUTANIX

Nexus® Series | Nexus Dashboard          UCS Series | Intersight®

**Observability**

Observability platform | Splunk

**Sustainability**

splunk>

Data center    Edge    Colocation    Public cloud

Cisco Connect

# Cisco AI PODs

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

# Cisco AI POD

## Scaling Units



**Scale Unit – Type 1**

**32**

Leaf – 2 x Cisco Nexus 9332D-GX2B

8 x 400GbE per node

4 x Cisco UCS C885A M8

SU = 4 Nodes / 32 GPUs

**Scale Unit – Type 2**

**64**

Leaf – 2 x Cisco Nexus 9364D-GX2A

8 x 400GbE per node

8 x Cisco UCS C885A M8

SU = 8 Nodes / 64 GPUs

**Scale Unit – Type 3**

**128**

Leaf – 2 x Cisco Nexus 9364E-SG2

8 x 400GbE per node

16 x Cisco UCS C885A M8

SU = 16 Nodes / 128 GPUs

**Note:**
- **Scale Unit Type: A pair of leaf switches + UCS nodes(leaf uplinks excluded)**
- **Non-blocking design for max GPU performance**

# Transform Cisco AI PODs into a GPU cloud

*Deliver sovereign and enterprise AI clouds*

*Experience delivery:*

**Training-aaS**    **Optimization-aaS**    **Inferencing-aaS**

*Cluster vending:*

| SELF-SERVICE GPU CONSUMPTION | GPU SLICING & POOLING | MULTI-TENANT CLUSTERS |
| --- | --- | --- |

**RUN:Ai**

**RAFAY**

### Cisco AI PODs

*OPERATIONS*    *AUTOMATION*    *AI SOFTWARE*

| CISCO — INTERSIGHT® & NEXUS DASHBOARD | A   </> | nVIDIA — NIM Operator NeMo CUDA |
| --- | --- | --- |

*Environment Mgt*

*KUBERNETES* — Red Hat OpenShift   ubuntu RANCHER BY SUSE

**FULL AI-STACK ORCHESTRATION**

*ACCELERATED COMPUTE* — CISCO UCS®

*HIGH-PERFORMANCE NETWORKING* — CISCO NEXUS

*K8s & Cloud Mgt*

*Governance and control:*

*EXTEND TO STORAGE PLATFORM ECOSYSTEM* — VAST · NetApp · PURE STORAGE · NUTANIX · HITACHI

CISCO

# Product view: Cisco Secure AI Factory with NVIDIA

## Cisco AI POD

### NVIDIA AI Enterprise
- NeMo
- NIM
- Blueprints

### NVIDIA Run:ai
- AI Workload & GPU Orchestration

### KUBERNETES PLATFORM
Options: Red Hat OpenShift, Nutanix NKP, Upstream Kubernetes

### CISCO AI NETWORKING
Cisco Data Center Switching
Options for Network Management: Nexus Dashboard or Nexus Hyperfabric |
Cisco Isovalent Enterprise Networking

### CISCO COMPUTE
Options: NVIDIA HGX, MGX or RTX PRO based UCS Servers
NVIDIA BlueField®-3 DPUs
Managed with Cisco Intersight®

### PARTNER STORAGE
Options: NetApp, Pure Storage, VAST Data, Hitachi Vantara, Nutanix NUS

## AI Security

### Cisco AI Defense

### Cisco Hybrid Mesh Firewall
Isovalent Runtime Security

Hypershield

Secure Firewall

### Splunk Enterprise Security

## AI Observability

### Splunk Observability

AI Infrastructure Monitoring

# AIML Network
*Larger Scale*

Cisco Nexus Dashboard

*non-blocking fabric*

Spines
Nexus 9364D GX2

*32-port 400G*

Leaf
Nexus 9364D GX2

*32-port 400G*

32 nodes UCS 885, 256 GPUs

*Smaller GPU clusters can use a single-switch network*

# Security at every layer of the stack

Cisco AI Defense

Cisco Hypershield



**Secure what AI is doing**

**Secure where AI is running**

**Only Cisco unifies networking, compute, security, and observability to deliver AI-ready data centers.**

# *Cisco Services*

*Customer Experience CX*

## The empowering force to help **you** achieve AI Ready DC

cisco *Connect*

# Cisco Services
## CX Bridge the Gap

**Lifecycle Services**

### From: Product

- Proposed Architecture
- Software Licenses
- Buying Programs
- Products
- Features / Capabilities

### The: Gap

Define
Document
Support
Integrate
Provision
Troubleshoot
Optimize
Install
Run
Adopt
Configure
Design
Policies
Maintain
Learn
Implement
Migrate
Manage
Train

### To: Outcomes

- Cost Reduction
- Improve Brand
- Enhance Security
- Reduce Risk
- Improve Experience

CISCO Connect

# CX and One Cisco
CX Offers to AI-Ready Data Centers

*Together with our partners, we get you to outcomes faster*



Empower Me ◄

► Do It For Me

▲ Advise Me

▲ Do It With Me

▼ Guide Me

Discovery Phase

Deployment Phase

Development Phase

Qualification

Ideation

Business Case

Proof of Concept

Pilot/
Prototype

Production

Postproduction

Expansion and
Scaling

# Cisco CX Unifies : AI Cluster Build & Scale

**The empowering force to help you achieve AI Ready DC**

*AI Platform – Plan & Design*

*Implementation, Test & Validate*

*AI Infrastructure Optimization & Support*

*Infrastructure Fine Tuning*

*AI Infra Consulting Support*

CISCO SERVICES

< < < < < < *Accelerated by Cisco CX* > > > > > >

# AI Ready Data Center - CX Services

*AI-Ready Data Centers*

### AI infrastructure Readiness

- AI Infrastructure Readiness

### AI POD Inference

- AI POD Adoption Service

### Private AI Foundation

- Customize Private AI for Training and Inference
- AI Ready Full Stack Infrastructure
- AI ML Network Services

### AI Digital Resilience

- AI Observability
- AI Security

CISCO

# AI Ready Data Center - CX Services

**Comprehensive Approach**

| Design | Implement | Fine Tune | Operations |
|--------|-----------|-----------|------------|

## AI Compute & Networking

*AI Infrastructure Readiness*
**Review and Analyze existing Infra**
**Identify gaps for AI Infra**

*AI POD Deployment Services*
**Leverage standardized architecture for AI POD deployment and adoption for Inference**

*Private AI – Compute & Network Services*
- **Infra Readiness, POCs**
- **Customized approach for a combined AI Compute & Networking env based on use cases**
- **Work alongside with teams to Fine tune & Optimize the AI Ready infrastructure**

## AI Security

*Securing AI*
**Comprehensive approach to protecting the development, deployment and use of AI applications.**

*Secure AI Ready DC Services*
- **Safeguard AI systems against evolving threats by implementing comprehensive strategies that address data integrity, model protection**
- **Advanced measures such as zone firewalling and segmentation to ensure robust protection against data breaches and cyberattacks**

## AI Observability

*AI Observability*
*Cisco's AI Observability services provide a comprehensive framework for ensuring resilience and performance with real-time insights, secure operations, and actionable data*

**AI Operations**
**Optimize AI infrastructure while preventing degradations and downtime**

CISCO Connect

# AI Ready Data Center
*Defining Infrastructure Requirements*

## Requirements

- *What is the use case?*
- *Am I training? Fine tuning? Inferencing?*
- *Am I using private data?*
- *Do I need a dedicated network?*
- *Who is responsible for management?*
- *Where will the workload live?*

## Considerations

- *Cost*
- *Accuracy*
- *Model Size*
- *User Experience (response time)*
- *Concurrent users*

- *Data Fidelity*
- *Power Requirements*
- *Cooling and Power Management*
- *Software Ecosystem*
- *Data Management Strategy*

# AI Compute Stack
## CX Services

### AI POD for Inferencing

Large language models ▶

AI tooling ▶ — **NVIDIA** NVAIE | NIMS

Kubernetes ▶ — **OPENSHIFT**

Operations ▶ — **CISCO** Nexus Dashboard and Intersight

Accelerated compute ▶ — **CISCO** UCS

**CISCO** Nexus

**NetApp**  **PURE**STORAGE

Automation ▶

### Private AI for Training/Inferencing

**Ecosystem**

intel AMD NVIDIA NUTANIX vmware Red Hat kubernetes
NetApp PURESTORAGE COHESITY veeam COMMVAULT Google Cloud Azure aws RoCE v2

| | |
|---|---|
| **DevOps** | **Applications** |
| | **Red Hat OpenShift**          **NVIDIA AI Enterprise** |
| **AI Ops** | **AI Powered Cyber Security, Observability & Data** |
| | **Infrastructure Management, Automation, & Monitoring** |
| **Compute** | |
| **Network** | |

# Storage Connectivity

- Support for standard IP-based storage protocols leveraging OS/software clients

- Connects to Shared storage for Loading (reading), Checkpointing (writing) model parameters



**NFS** — Structured data or file system

**Object** — Unstructured data or object storage

**Other** — Other IP-based data sources

*Note: Most customers leverage a dedicated storage network, but storage traffic can be aggregated with other north-south networks as well*

**PARTNERSHIPS WITH**

PURESTORAGE®    NetApp

# Every AI/ML Network is different



*Build the model*
**Training**

*Optimize the model*
**Fine-tuning and RAG**

*Use the model*
**Inferencing**

**From Training to Inferencing: A Songwriting Analogy**

# AI/ML Network - CX Services
## Need to be best Optimized

| Training requirements | Inference requirements | General requirements |
|---|---|---|

### Training requirements

**High bandwidth**

*Movement of massive datasets across the network demands high bandwidth networks*

**Non-Blocking Lossless**

*Network inconsistencies can affect the accuracy and training time of AI models*

**Congestion management**

*Detect potential congestion and redistribute network traffic accordingly*

### Inference requirements

**Low latency**

*Real-time AI applications require extremely low latency*

### General requirements

**Visibility**

*Comprehensive visibility tools for real-time monitoring, issue detection, and troubleshooting*

**Scalability**

*Dynamic nature of AI workloads mean that the network needs to be agile and scalable*

---

**AI/ML Infrastructure Service**  |  **AI Pods Service**  |  **Front End Network**

---

| Latency | Network Losses | Bandwidth/Scale | Load Balancing | Visibility |
|---|---|---|---|---|

# Optimizing AI/ML Cluster

Parallelism Techniques  - Distributed Training and Inference
About running AI smarter and more efficiently

Model Parallelism

## Data Parallelism

| Model Copy | Model Copy | Model Copy | Model Copy |
|---|---|---|---|
| W | W | W | W |
| GPU1 | GPU2 | GPU3 | GPU4 |
| G | G | G | G |

Dataset

Uses *All-Reduce* to sum all gradients together for the full batch, update weights, learn, repeat

## Pipeline

| Layer 0 | Layer 1 | Layer 2 | Layer 3 |
|---|---|---|---|
| W | W | W | W |
| GPU1 | GPU2 | GPU3 | GPU4 |
| G | G | G | G |

*Send Recv* *Send Recv* *Send Recv*

It splits a *model* into stages or groups of layers across multiple GPUs.

The dataset goes from one GPU to the next using *send/receive* operations.

## Tensor

| GPU1 | $A$ | x | $B0$ | | $=$ | $C0$ | | → | $C$ |
| GPU2 | $A$ | x | | $B1$ | $=$ | | $C1$ | → | $C$ |
| GPU3 | $A$ | x | | $B2$ | $=$ | | $C2$ | → | $C$ |
| GPU4 | $A$ | x | | $B3$ | $=$ | | $C3$ | → | $C$ |

Groups of parameters (weights, gradients) within a layer are split across GPUs.

Uses *All-Gather* for all GPUs to get aggregated results

cisco Connect

# Digital resilience for AI-ready data centers

**Digital resilience**

Security    Assurance    Observability

## AI-ready data centers

| Build the Model I Training | Optimize the Model I Fine-tuning and RAG | Use the Model I Inferencing | Model the World I Digital Twin |



AI PODs Data Center
(inferencing suite)

Cisco UCS C-Series Accelerated
(GPUs)

Cisco Nexus
HyperFabric

Cisco Nexus
HyperFabric + AI Cluster

*Private cloud managed AI Infrastructure*

*Cloud managed AI Infrastructure*

PROTECT ASSETS | OPTIMISZE PERFORMANCE | MAINTAIN TRUST

## Security

## Observability

### Data Protection
**Protecting data from breaches and misuse**

Competitive Advantage

Regulatory Compliance

### Model Integrity
**Ensure models remain untampered and resistant to attacks**

Business Continuity

Customer Retention

### AI Operations
**Optimize operations to drive digital resilience**

Rapid Detection and Resolution

Controlled AI Costs

### Digital Experience
**Realtime insights assure the user AI experience**

Accelerated adoption

Improved Brand Loyalty

# Cisco Secure AI

**Changing the fundamentals of data center security**

| Visibility | Zone firewalling | Segmentation | Exploit protection |
| --- | --- | --- | --- |

Cisco Security Cloud Control

Secure Firewall

Secure Workload

Isovalent

Hypershield

AI Defense

*Security infused into the network through a fabric of enforcement points*

# Observability for AI
## Unified Experience for Infrastructure Operations
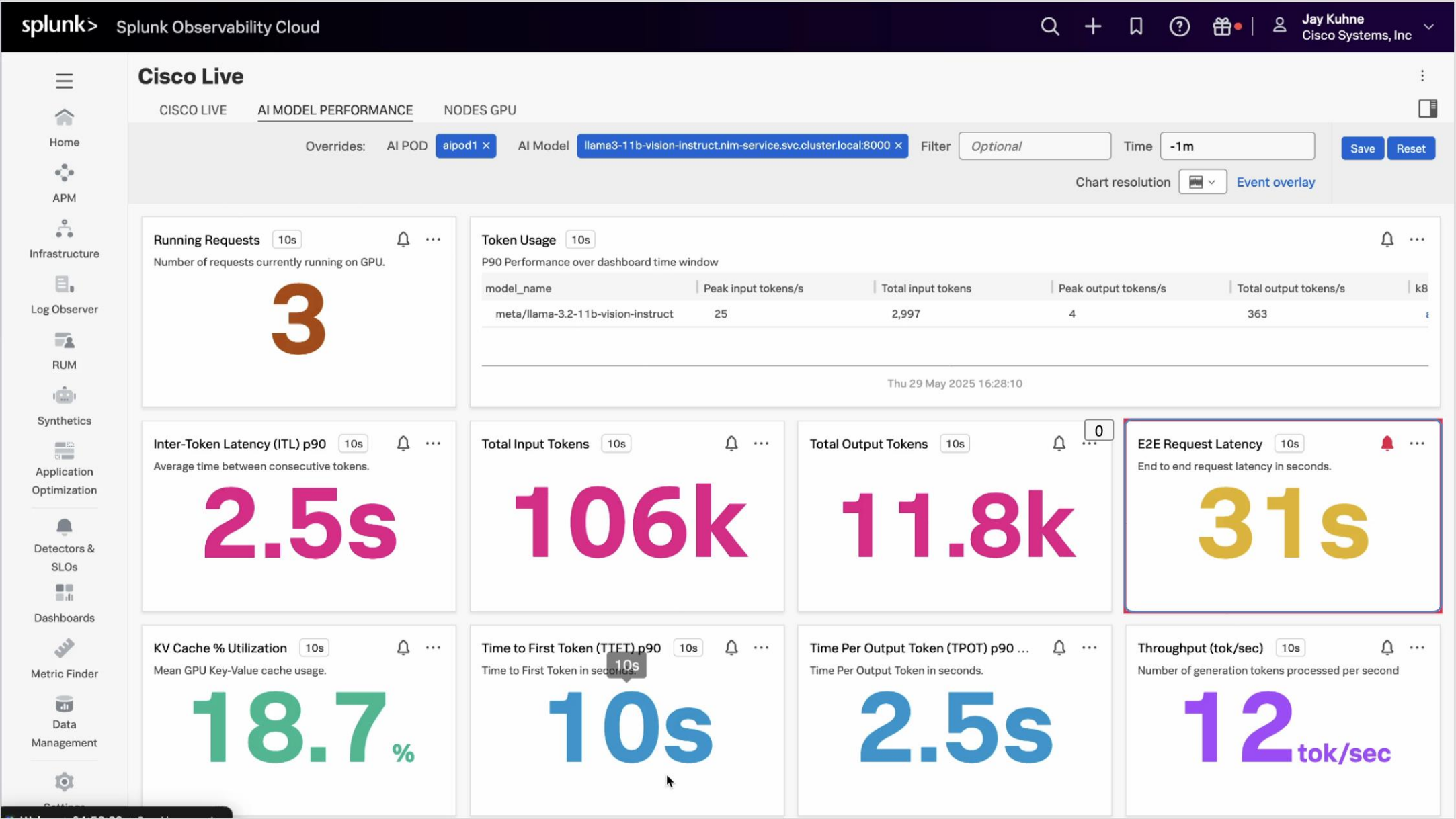


**Cisco Intersight**

**Cisco Nexus Dashboard**

*See   Connect   Secure   Automate*

*Simplified   Accelerated   Insights*

**Powerful Integrations with industry leading platforms**

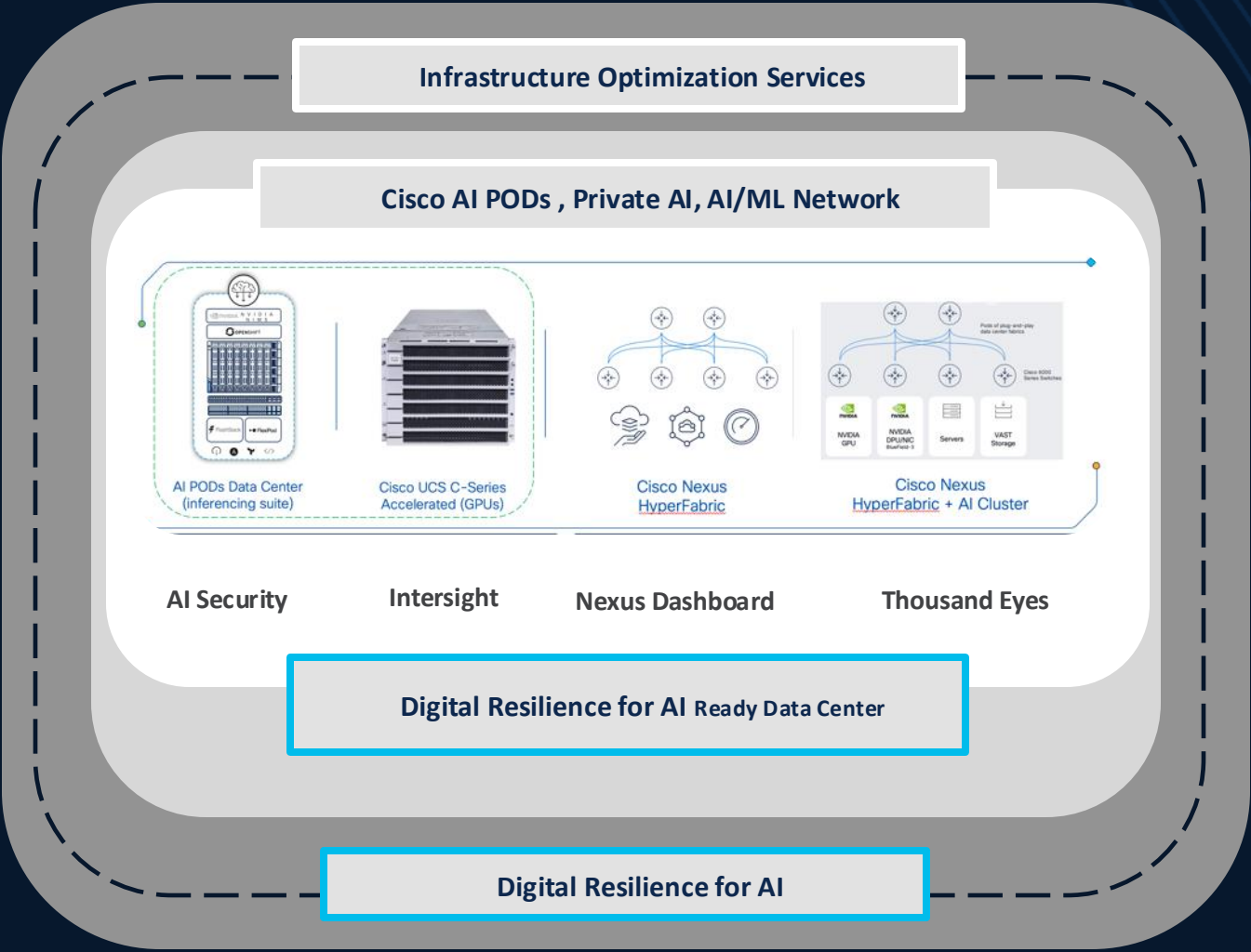# Splunk Observability for AI

# AI Ready Data Center Services Portfolio

## AI Adoption Journey/Framework



**STRATEGY**

- AI Use-Case Identification
- AI Solution Design
- AI Solution Validation
- AI Solution Proof of Concept
- AI Solution Build and Deploy
- AI Solution Adopt and Scale
- AI Solution Maintain

AI Governance

**+**

## CX Services Portfolio for AI Ready DCs

*Design , Deploy & Observ*

**Infrastructure Optimization Services**

**Cisco AI PODs , Private AI, AI/ML Network**



AI PODs Data Center (inferencing suite) | Cisco UCS C-Series Accelerated (GPUs) | Cisco Nexus HyperFabric | Cisco Nexus HyperFabric + AI Cluster

**AI Security**   **Intersight**   **Nexus Dashboard**   **Thousand Eyes**

**Digital Resilience for AI** Ready Data Center

**Digital Resilience for AI**

Product   |   Implementation Services   |   Premium Services

# AI Ready Enterprises

AI Platform

Infra Readiness

Plan & Design

Implement & Optimize

## CX Approach

Automate & Operate

Observability for AI

Hand Holding & Guidance

Life Cycle Services

New AI Operating Model

CX Capability Beyond AI Platform

Cloud + AI Infrastructure

*Ripple Effect - Beyond DC – Ex: DCI, SDWAN , Campus, Security , Ops, D2 Services ...*

# AI Ready Enterprises

# Enterprise AI Use cases

## Knowledgebase copilots
AI assistants

## Content & code generation
Text | Images | Video | Code

## Security
Physical | Virtual

## Language translation
Multilingual real-time communication

## Virtual agent & chatbots
Specialized domain specific chatbots

## Detection, prediction & reporting
Analytics | Forecasts | Anomalies | Insights

CISCO

# AI Ready Data Center – *CX Services*

## Customer Asks

| Early Alignment | Rapid Prototype | Security for AI | E2E Services | Storage & Data |
|---|---|---|---|---|

| AI Governance | As a Service | Adoptable Network | Expansion | Constraint Handling |
|---|---|---|---|---|

| Resource Management | Customer Ref | Cloud vs OnPrem | Convergence | Operational Challenges |
|---|---|---|---|---|

Cloud + AI Infrastructure

# AI Is Rising — The Real Wave Is Still Ahead

## Agentic workflows will power the next wave of demand

| LAYER | DRIVER | IMPACT |
|---|---|---|
| **Adoption** X **Capability** X **Interaction** | Commoditization 💸 | ↑ Lower cost & demand growth |
| | Better LLMs 🧠 | ↑ Addressable tasks |
| | Multi-modal 📊 | ↑ Payload size |
| | Reasoning 🔢 | ↑ Tokens & Compute |
| | AI Agents 🤖 | ↑ Transactions per request |
| | Multi-agent 🌐 | ↑ Requests per session |

**1448%**
Increase in tokens processed by AI models in the last 12 months[1]

**75%+**
Of inference-driven data creation and processing at the edge by 2030[2]

**10x**
Bandwidth required upstream compared to downstream

**36x**
Increase in AI traffic as early as 2023-2024[1]

*1: Openrouter.ai/rankings 2: KPMG*

# The Rise of Agentic AI



**Business value** (y-axis)

Timeline across top: **< 2022**, **2022-2023**, **2024**, **2025+**

- Traditional AI (< 2022)
- Single Modal Generative AI (2022-2023)
- Multimodal Generative AI (2024)
- Agentic AI (2025+)

### Agents driving traffic multiplication

**Network Traffic in the Digital Era**

Internet

Requested and processed at "human speed"

**Network Traffic in the Agentic Era**

Agent → Internet

Requested and processed at "agent speed", 24x7
Potential strain on surrounding infrastructure

$n \times rtt$
$n \times BW$

"33% of enterprise software applications will include
agentic AI by 2028, up from less than 1% in 2024"

— Gartner

# AI is re-shaping DCs



AI infrastructure
for enterprise

Enterprise data center
AI + traditional

AI infrastructure for hyperscalers

# Enterprise AI Convergence

**AI PODs**

**Convergence LLM Fabric**

**Legacy DC**

| NVIDIA Software |
| Platform Software |
| Cisco Networking & Optics |
| Cisco Compute GPU & CPU |
| Partner Storage |

Every AI request is unique – no content is cached

Multi-LLM Challenges

AI requests generate a high volume of tokens

Distribution of Agents and LLMs

AI drives a disproportional increase in upstream traffic

| Enterprise Apps |
| Platform Software |
| Cisco Networking & Optics |
| Cisco Compute   CPU |
| Partner Storage |

**AI Workloads**

**Non-AI Workloads**

# Future of AI !

# Industry view
## No analyst firms have an overall AI report

- No industry standard of the AI market

- AI made up of several market segments

- No magic quadrant for AI

- AI is Reported in all segments

- AI effects all market segments

- Market research suggest Enterprise AI will be consumed as a service
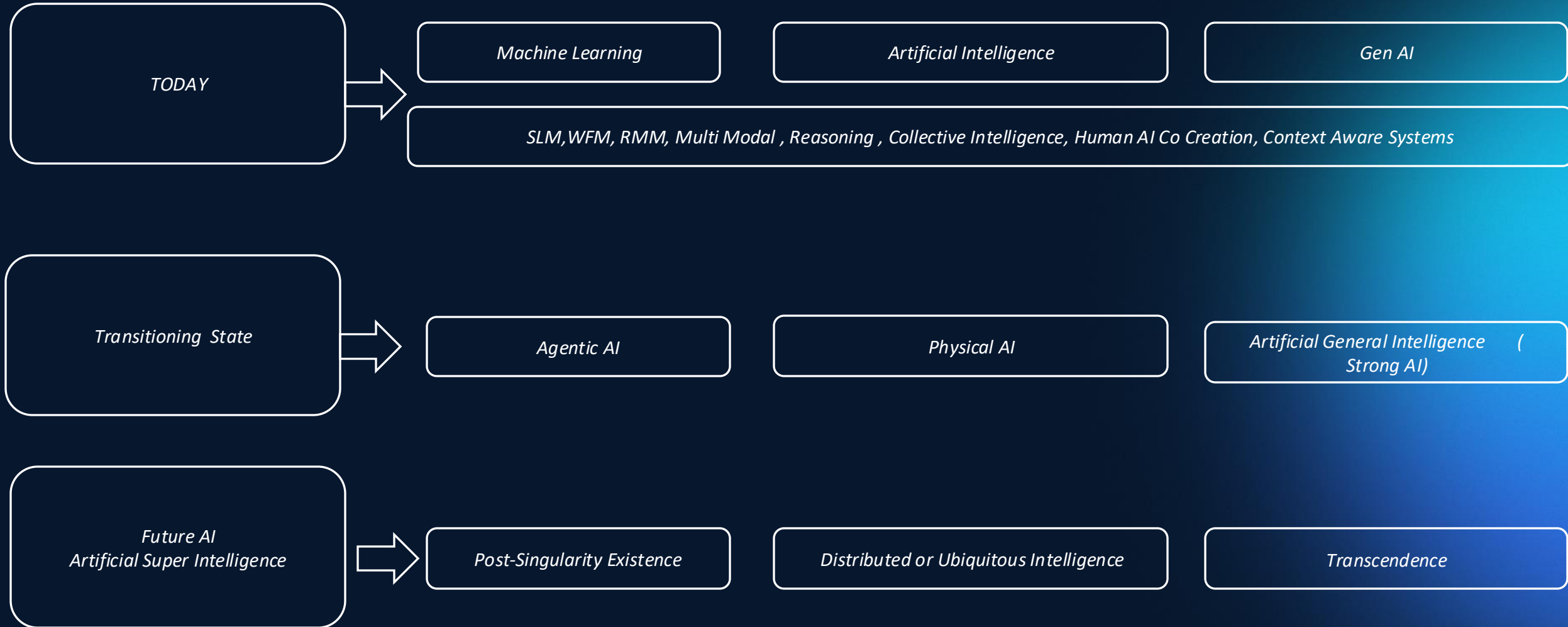


Largest publicly announced AI training runs by company — EPOCH AI

# The Future  !

| TODAY | → | Machine Learning | Artificial Intelligence | Gen AI |
|---|---|---|---|---|

SLM,WFM, RMM, Multi Modal , Reasoning , Collective Intelligence, Human AI Co Creation, Context Aware Systems

| Transitioning State | → | Agentic AI | Physical AI | Artificial General Intelligence    ( Strong AI) |
|---|---|---|---|---|

| Future AI Artificial Super Intelligence | → | Post-Singularity Existence | Distributed or Ubiquitous Intelligence | Transcendence |
|---|---|---|---|---|

CISCO

# The Roadmap



AI Chat

Attended AI Agents

Unattended AI Agents

| | 2022 | 2023 | 2024 | 2024/2025 | 2026/2027 | | 2030+? |
|---|---|---|---|---|---|---|---|
| | Standalone AI Chat | RAG Powered AI Chat | Reasoning AI Chat | Standalone AI Agent | Multi AI Agent System | Built for purpose | General purpose |
| Value | Low (Micro efficiencies) | Low/Medium | Medium | Medium/High | High | Very High | Revolutionary |
| Need for low latency | Low | Medium | High | High | Very High | Very High | Unknown |
| Volume / tokens | Medium | Low/Medium | High | High | Very High | Extremely High | Unknown |
| Risk exposure | Medium | High | High | Very High | Very High | Extremely High | Unknown |

Definitions: LLM = Large Language Model, LMM = Large Multimodal Model, UI = User Interface, RAG = Retrieval Augmented Generation, API = Application Programming Interface,

**AI will make our world of *8B* people feel like one with the capacity of *80B***

*"80B Agents + 8B"*

*Cisco CX along with our Partners unifies networking, compute, security, and observability to deliver AI-ready Accelerated data centers.*

CiscoConnect