# Unlock the Future: Leveraging Generative AIOPs for Enterprise Innovation and Performance

Sanjit Aiyappa – Cisco, Director Solutions Engineering
Joel Jose – Cisco, Senior Solutions Engineer
AIHUB - 1006

Cisco Live!

# Listening to the market

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

GenAI is expected to be the **#1 driver** of **infrastructure investments** over the next 18 months

**$300B** global spending on AI by 2026

**97%** of companies say the **urgency to deploy AI-powered technologies** has increased

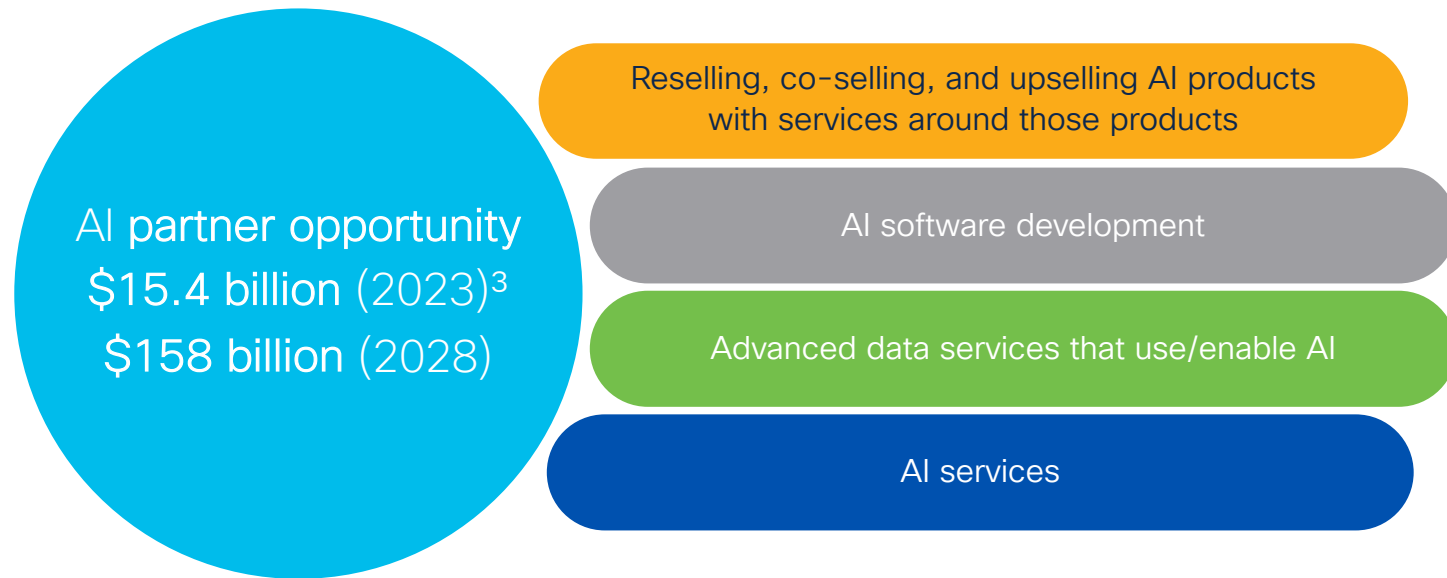**$15.7T**
Potential contribution to global economy by 2030

GenAI has potential to generate **$2.6 to $4.4 trillion** in **value across industries**[2]

[1]IDC Artificial Intelligence Systems Spending Guide
[2] McKinsey & Company
[3]Canalys estiates & forecasts, August 2023
Gartner, May 2024

# Capitalizing on AI for channel partners

AI partner opportunity
$15.4 billion (2023)[3]
$158 billion (2028)

Reselling, co-selling, and upselling AI products with services around those products

AI software development

Advanced data services that use/enable AI

AI services

# Heeding the voices of our partners

## Partner AI readiness index

**Biggest challenges supporting customers to deploy AI**

Lack of experience in new areas of tech deployment **62%**

Lack of knowledge of systems and processes **56%**

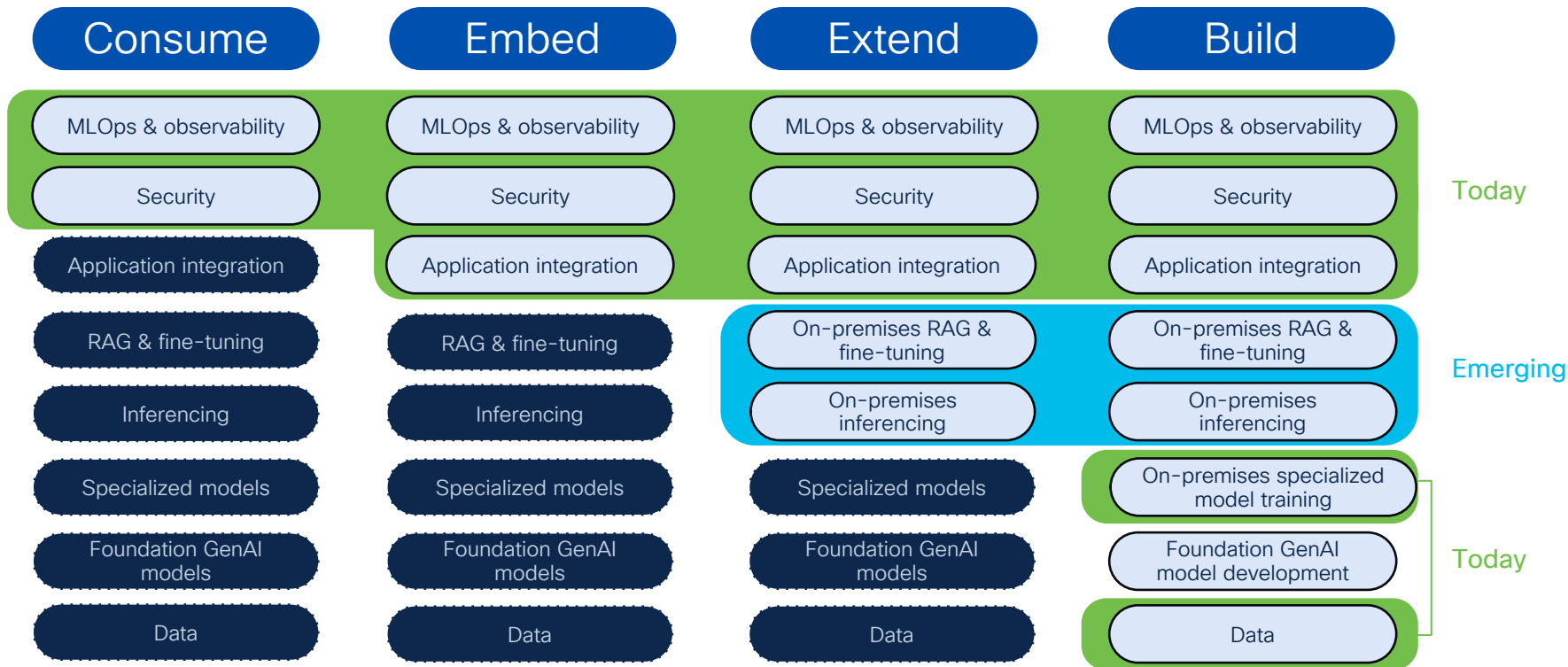Lack of knowledge of available technologies **53%**

27% of partners believe

**76-100%**

of **revenue** will come from **AI technologies** in next 4-5 years

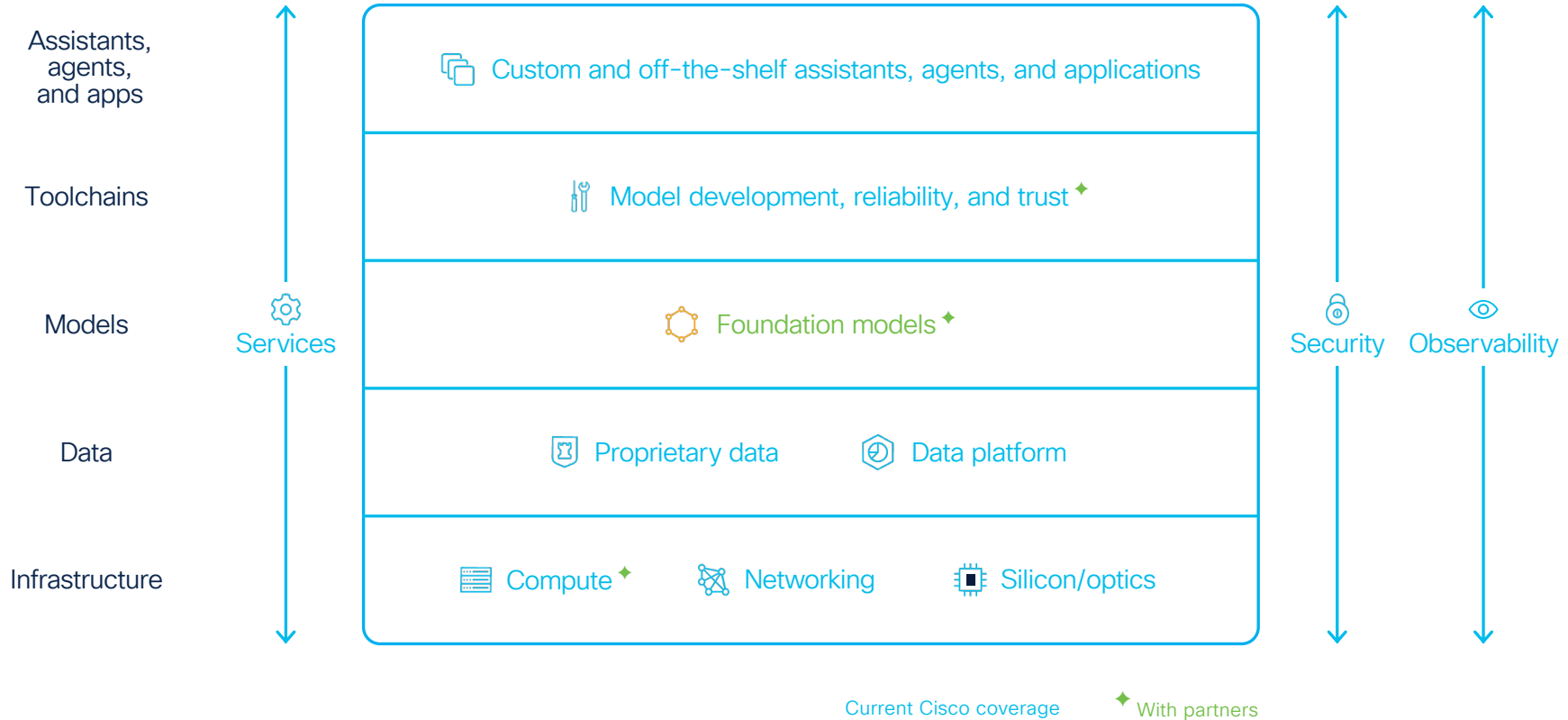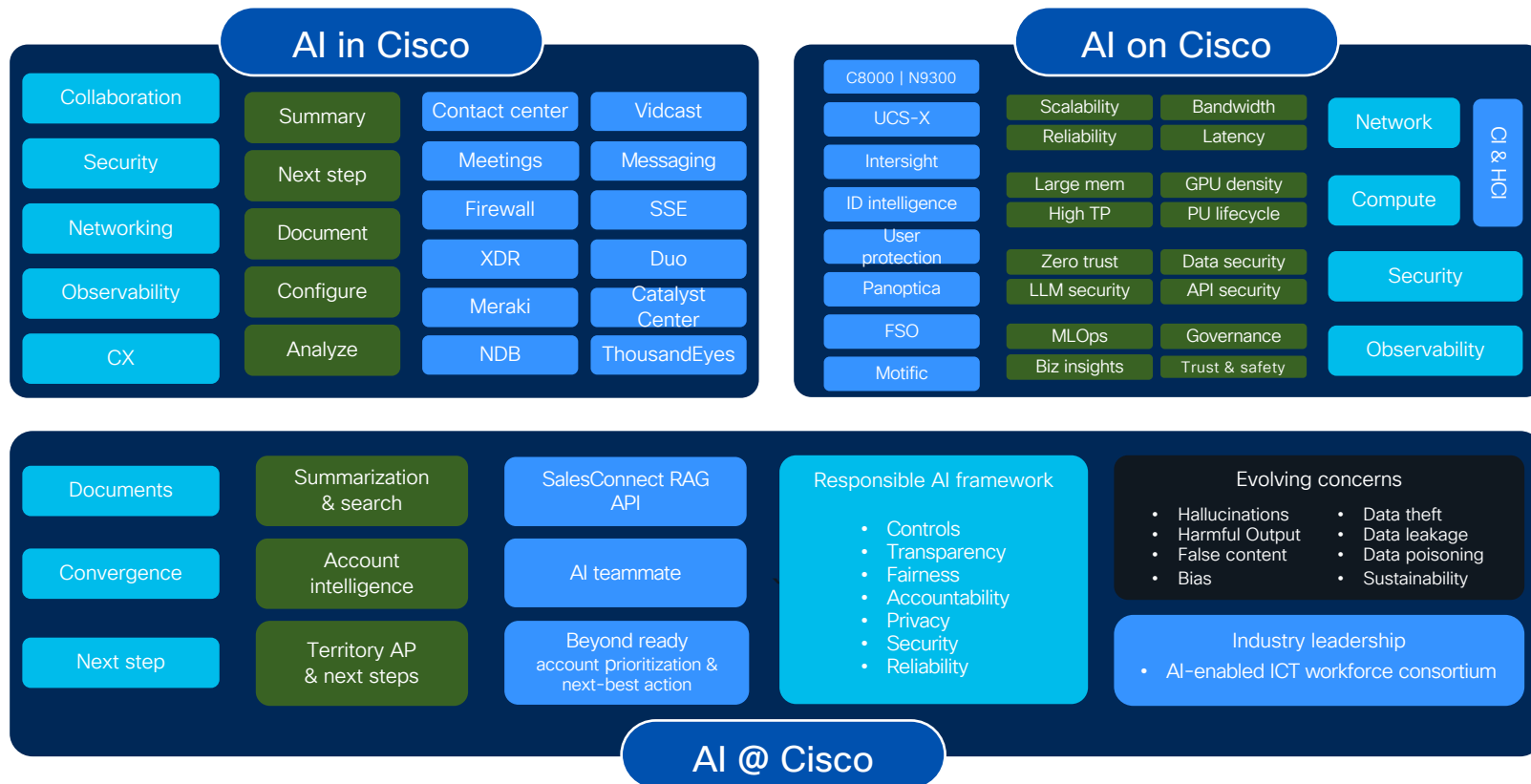**Biggest revenue drivers ALL aligned to Cisco**

Infrastructure **33%**

Cybersecurity **19%**

Customer experience **10%**

# Where is our opportunity?

Gartner AI consumption models and customer journey

| Consume | Embed | Extend | Build |
|---|---|---|---|
| MLOps & observability | MLOps & observability | MLOps & observability | MLOps & observability |
| Security | Security | Security | Security |
| Application integration | Application integration | Application integration | Application integration |
| RAG & fine-tuning | RAG & fine-tuning | On-premises RAG & fine-tuning | On-premises RAG & fine-tuning |
| Inferencing | Inferencing | On-premises inferencing | On-premises inferencing |
| Specialized models | Specialized models | Specialized models | On-premises specialized model training |
| Foundation GenAI models | Foundation GenAI models | Foundation GenAI models | Foundation GenAI model development |
| Data | Data | Data | Data |

Today

Emerging

Today

# Cisco supports customers across the AI stack



Assistants, agents, and apps — Custom and off-the-shelf assistants, agents, and applications

Toolchains — Model development, reliability, and trust

Models — Foundation models

Data — Proprietary data · Data platform

Infrastructure — Compute · Networking · Silicon/optics

Services

Security

Observability

Current Cisco coverage · With partners

# Map of AI in, on, and @ Cisco

## AI in Cisco

| Collaboration | Summary | Contact center | Vidcast |
| Security | Next step | Meetings | Messaging |
| Networking | Document | Firewall | SSE |
| Observability | Configure | XDR | Duo |
| | | Meraki | Catalyst Center |
| CX | Analyze | NDB | ThousandEyes |

## AI on Cisco

C8000 | N9300
UCS-X
Intersight
ID intelligence
User protection
Panoptica
FSO
Motific

| Scalability | Bandwidth | Network |
| Reliability | Latency | |
| Large mem | GPU density | Compute |
| High TP | PU lifecycle | |
| Zero trust | Data security | Security |
| LLM security | API security | |
| MLOps | Governance | Observability |
| Biz insights | Trust & safety | |

CI & HCI

## AI @ Cisco

| Documents | Summarization & search | SalesConnect RAG API |
| Convergence | Account intelligence | AI teammate |
| Next step | Territory AP & next steps | Beyond ready account prioritization & next-best action |

**Responsible AI framework**
- Controls
- Transparency
- Fairness
- Accountability
- Privacy
- Security
- Reliability

**Evolving concerns**
- Hallucinations
- Harmful Output
- False content
- Bias
- Data theft
- Data leakage
- Data poisoning
- Sustainability

**Industry leadership**
- AI-enabled ICT workforce consortium

CISCO *Live!*

# Why AIOPs

**15–30%**
cost savings

**30%**
increased productivity

**≠70%**
increase in FTTR

**5–10%**
revenue growth

**-25%**
compliance fines

# AIOPs outcomes

## Top use cases

- Anomaly detection
- Root cause analysis
- Predictive analytics
- Change impact analysis
- Capacity planning
- Security monitoring and compliance
- Incident management
- Performance monitoring and optimization
- Change management

## Benefits

- Operational efficiency
- Cost reduction
- Faster problem resolution
- Proactive monitoring and predictive analytics
- Enhanced scalability
- Optimized resource utilization
- Reduce alert noise and false positive
- Increase service reliability
- Data-driven insights for decision makers

# AI-driven portfolio today

## AI with Cisco for MSP

| AI ASSISTANT | | | | |
| --- | --- | --- | --- | --- |
| **Cloud and AI** | **Networking Cloud** | **Collaboration** | **Security** | **Observability** |

**AI-Driven Portfolio**

| Cloud and AI | Networking Cloud | Collaboration | Security | Observability |
| --- | --- | --- | --- | --- |
| • Enabled by Silicon One Scheduled Fabric Ethernet Solution<br>• Digital Experience Monitoring & predictive analytics<br>• AI PODs<br>• Nexus with RoCE<br>• Hyperfabric AI/ML<br>• Hyper Shield | • Predictive Path Recommendation (SD-WAN)<br>• Bandwidth Forecast (SD-WAN)<br>• Machine Reasoning Engine (MRE)<br>• AI Network Analytics<br>• AI-Driven Baselining<br>• AI-Enhanced RRM<br>• Anomaly Detection (SD-WAN)<br>• Client Analytics and RCA | Cloud Contact Center<br> • AI Powered Analytics<br> • AI Call Summaries<br> • AI Scripted Agents<br> • Automatic CSAT Scores<br> • Agent Wellness<br>Webex Calling<br> • Catch Up and Summarization<br>Webex Meetings<br> • AI Meeting Summary<br> • AI Codec for Voice and Video | • Extended/Managed Detection & Response<br>• AI Assistant Experience (inc. Firewall, XDR, & cross-portfolio)<br>• SPLUNK (SOC)<br>• Autonomous Actions (incl. incident response, recs, & automation) | • Statistical Modeling with Baselining<br>• Anomaly Detection with dynamic baselining<br>• Intelligent Automation<br>• SPLUNK – Oly & ITSI, AI assistant for SPL<br>• Prompt interface with AI-powered workflows<br>• AI assistant for summarization |

**Partner led usecases**

| Cloud and AI | Networking Cloud | Collaboration | Security | Observability |
| --- | --- | --- | --- | --- |
| • Network Management & Deployment<br>• Service Assurance<br>• Anomaly Detection<br>• Change Automation<br>• Predictive Internet Insights<br>• Customized AI Models for Compliance Monitoring<br>• Machine Learning-Enhanced Predictive Analytics for Customer Satisfaction | • Anomaly Detection and Alert Prioritization<br>• Predictive Incident Management<br>• Incident Root Cause Analysis (RCA) with Domain-Specific Insights<br>• Automated, Contextualized Remediation Actions<br>• Dynamic Incident Correlation and Noise Reduction | • NLP Chat bot for ops workflow (knowledge Base Integration)<br>• Custom ML bots integration with CC.<br>• AI model to detect voice/video issues across network path<br>• Proactive notification and auto meeting room booking for employee behalf.<br>• Intelligent workflows to see usage consumption. of meeting rooms<br>• Predictive insights on network usage for Voice and video calls. | • Proactively hunting - searching for hidden threats to prevent the attack from happening<br>• Understanding and Mitigating Advanced Threats<br>• Executing the incident response plan<br>• Define automation and playbooks | • Intelligent root cause analysis with event prioritization<br>• Realtime business health monitoring with data tagging<br>• Data security posture management<br>• Adaptive thresholding for optimized SLAs |

# Types of AI

**Generative AI**

Synthesize signal to improve user productivity and outcomes

**Foundational AI**

Make sense of the signal in vast amounts of data

# Rightsizing GenAI

**1. Understand Workload Needs**
    Training vs. Inference
    Model Complexity
    Batch Size/Dataset

**2. Profile & Benchmark**
    Tools: NVIDIA Nsight, TensorBoard, PyTorch Profiler.
    Key Metrics: GPU/CPU utilization, memory, I/O bottlenecks.

**3. Optimize Infrastructure**
    Compute: Use GPUs (e.g., A100 for training, T4 for inference) or accelerators (TPUs, FPGAs).
    On-Prem vs. Cloud

**4. Scale Efficiently**
    Horizontal
    Vertical
    Elastic

**5. Enhance Resource Usage**
    Quantization/Pruning: Optimize model efficiency.
    Data Pipelines: Caching, sharding for faster I/O.
    Orchestration: Kubernetes

**6. Monitor Continuously**

# Advanced RAG Implementations

- **Hybrid Retrieval (Dense + Sparse)**: Combines the strengths of dense (e.g., vector-based) and sparse (e.g., BM25) methods for more robust and accurate document retrieval.

- **Iterative RAG (Re-RAG)**: Refines outputs by iteratively retrieving additional context and updating the response for higher accuracy and reduced ambiguity.

- **Memory-Augmented RAG**: Integrates short-term and long-term memory mechanisms to maintain context and personalize interactions over time.

- **Multimodal RAG**: Extends RAG to handle multiple modalities, such as text, images, and videos, enabling applications like visual question answering and multimodal search.

- **RAG with Feedback Loops**: Incorporates human-in-the-loop or automated feedback to iteratively improve retrieval and generation quality, using reinforcement learning or fine-tuning.

# Powered by Splunk AI

## Product overview



**SECURITY**

Enterprise security with enterprise security content updates (ESCU)

User behavior analytics

**AI assistant**

**OBSERVABILITY**

IT service intelligence

Application performance monitoring

Infrastructure monitoring

**AI assistant**

*Included embedded AI/ML capabilities*

**Assistive intelligence experiences**

**AI assistant for SPL**

**Anomaly detection**

**Customizable ML**

**Machine learning toolkit**

**Data science and deep learning**

**Python for scientific computing**

**THE SPLUNK PLATFORM**

**Splunk Cloud Platform**

**Splunk Enterprise**

*Free assistive and customizable apps & tools*

# Splunk Observability portfolio

Application Performance
Monitoring (APM)

Digital Experience
Monitoring

Infrastructure
Monitoring

AIOps

Splunk Observability

Log Analysis

Business Risk
Observability

Incident
Response

Network
Monitoring

AI / ML / LLM
Observability

Real-Time Insights

AI Powered

Enterprise Grade

Open Telemetry Native

Extensible

Cross MELT

Business Context

# Cisco Networking + Splunk: AI & ML Based

- AI to process the data and find the issues (AI, ML)
- AI to assist in understanding and interacting with the data (Gen AI Assistants)
- AI & ML throughout the incident response

One Catalyst App on Splunk

**AI Assistants**

**IT Services Intelligence and o11yCloud**

**AppD**

**splunk>**
a **CISCO** company

| Infrastructure Metrics, Events, Logs, and Traces | Catalyst Center | Catalyst SD-WAN | Cisco Meraki | Thousand Eyes | 3rd Party |
|---|---|---|---|---|---|

Essential for correlation of events across disparate data types and sources

Radio Resource Mgmt
ThousandEyes Agent integrations

WAN Insights with Predictive Path analysis, Bandwidth Forecasting
Sites, Clients, Circuits, Apps dashboard
Traffic flow and App distribution patterns

Radio Resource Mgmt,
WiFi 3D Analyzer
WiFi 6/6E planning tools
Endpoint Analytics
PoE Analytics
Cisco ecosystem connectivity
Site & Event Analytics

WAN Assurance
Active Application Monitoring
Internet Outage

cisco *Live!*

# Catalyst Common Splunk Application – Outcomes



**Benefits**
- Holistic view of Network, Security and Identity events
- Long time data retention for compliance needs
- Charts combine data from multiple sources for ease of visualization (ex Top Threats combine threats from WAN and ISE)

**Outcomes**
- Improve troubleshooting by tracking device and endpoint issues
- Reduce Mean Time to Resolution (MTTR) by quickly identifying critical events
- Single repository for NGFW events (Malware Protection, Intrusion Prevention, URL-Filtering) and Security Advisories
- Secure grip over the network by tracking Access Control List (ACL) entries across WAN and LAN

# Catalyst Common Splunk Application – Use Cases

## Consolidated Visibility

- Consume data from different solutions (DNA/SD-WAN/ISE) using standard data sets
- Common dashboard that visualizes data from all solutions

- Outcomes : Realtime Monitoring, Historical Insight, Network and Client Insight, Security Insight, Compliance & Security advisory

**MVP Focus**

## Analytics Dashboard

- Correlate information from the various solutions and establish baselines

- Outcome: Detect and report on anomalies based on deviation from the baselines

## Playbook Driven Response

- Playbooks that look for certain event triggers and generate API calls back to the appropriate domain (SD-WAN/DNA/ISE)

- Outcome: Automated Event Response

## Splunk Ecosystem Partner

- Trigger Notifications to 3rd Party Platforms
- Splunk Already Integrates with (1,000+) existing apps

- Outcome: Integration with third party event management systems

Catalyst SD-WAN

Catalyst Center

Cisco ISE

Common Splunk Dashboard

# Demo

- MSP BOT with troubleshooting using GenAI

- Advanced Cost-effective RAG for Enterprise

- Interactive API documentation using GenAI

- Alert reduction and troubleshooting using GenAI

# Cisco supports to build outcomes



AIOPS hackathon workshops for partners
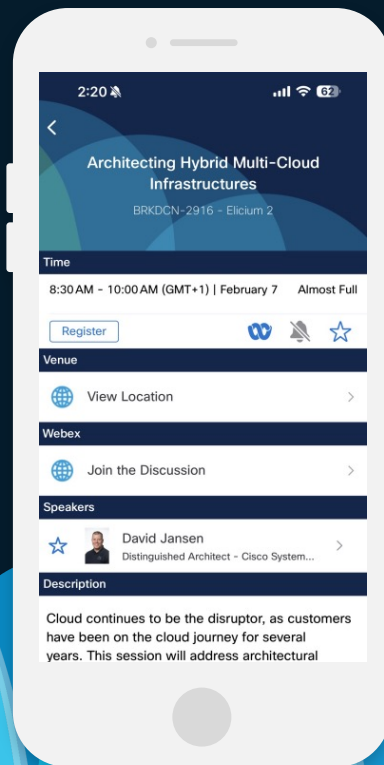
Cisco Black Belt for AI

# Webex App

## Questions?
Use the Webex app to chat with the speaker after the session

## How

1. Find this session in the Cisco Events mobile app

2. Click "Join the Discussion"

3. Install the Webex app or go directly to the Webex space

4. Enter messages/questions in the Webex space

Webex spaces will be moderated
by the speaker until February 28, 2025.

# Fill Out Your Session Surveys

Participants who fill out a minimum of 4 session surveys and the overall event survey will get a unique Cisco Live t-shirt.

(from 11:30 on Thursday, while supplies last)

All surveys can be taken in the Cisco Events mobile app or by logging in to the Session Catalog and clicking the 'Participant Dashboard'

**Content Catalog**

# Continue your education

- Visit the Cisco Showcase for related demos

- Book your one-on-one Meet the Engineer meeting

- Attend the interactive education with DevNet, Capture the Flag, and Walk-in Labs

- Visit the On-Demand Library for more sessions at ciscolive.com/on-demand. Sessions from this event will be available from March 3.

Thank you

GO BEYOND