# Unleash AI Faster

Automated Cisco AI Pod Setup with Intersight
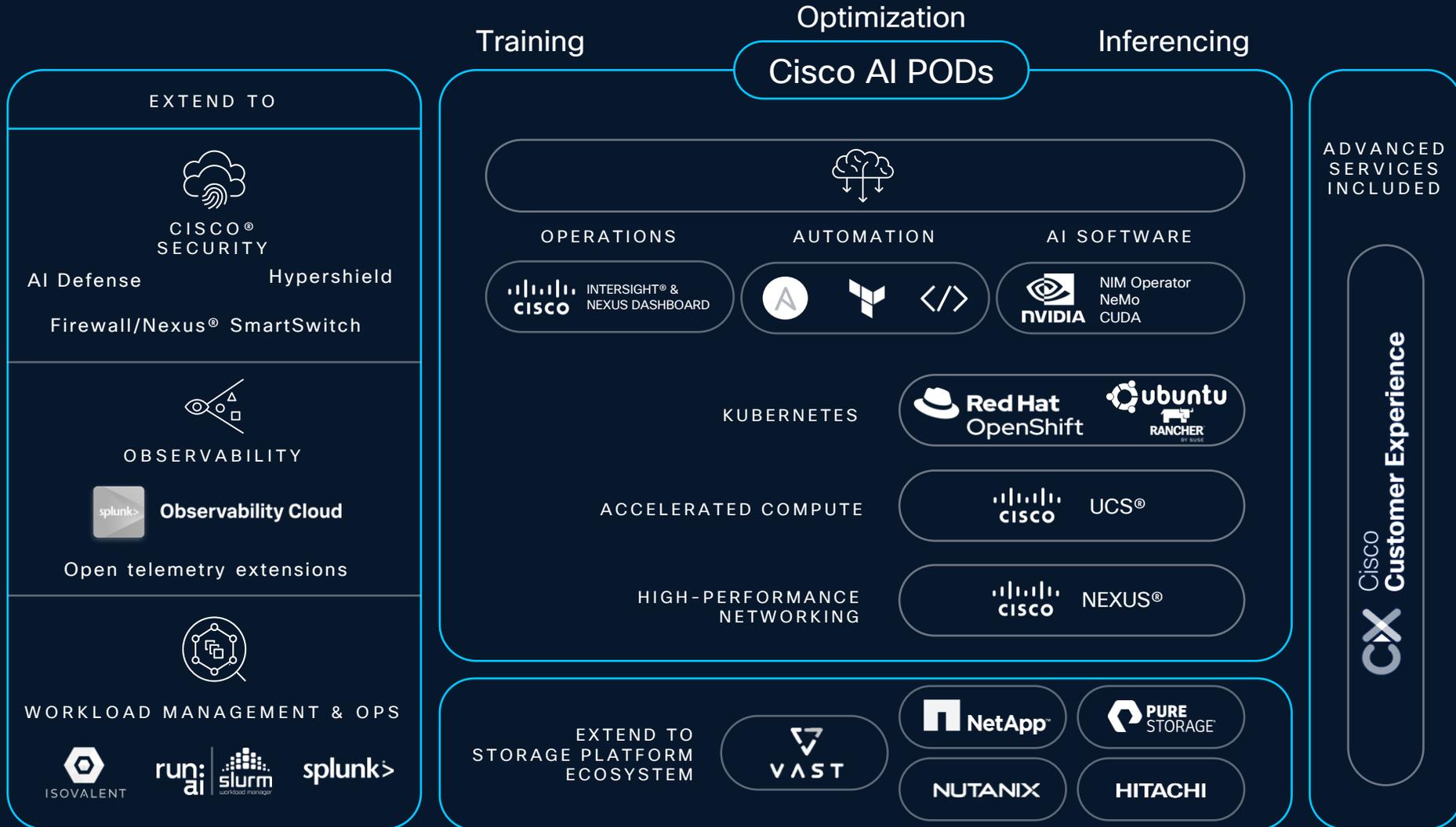
Marc Abu El Ghait
Customer Success Specialist – Compute & AI
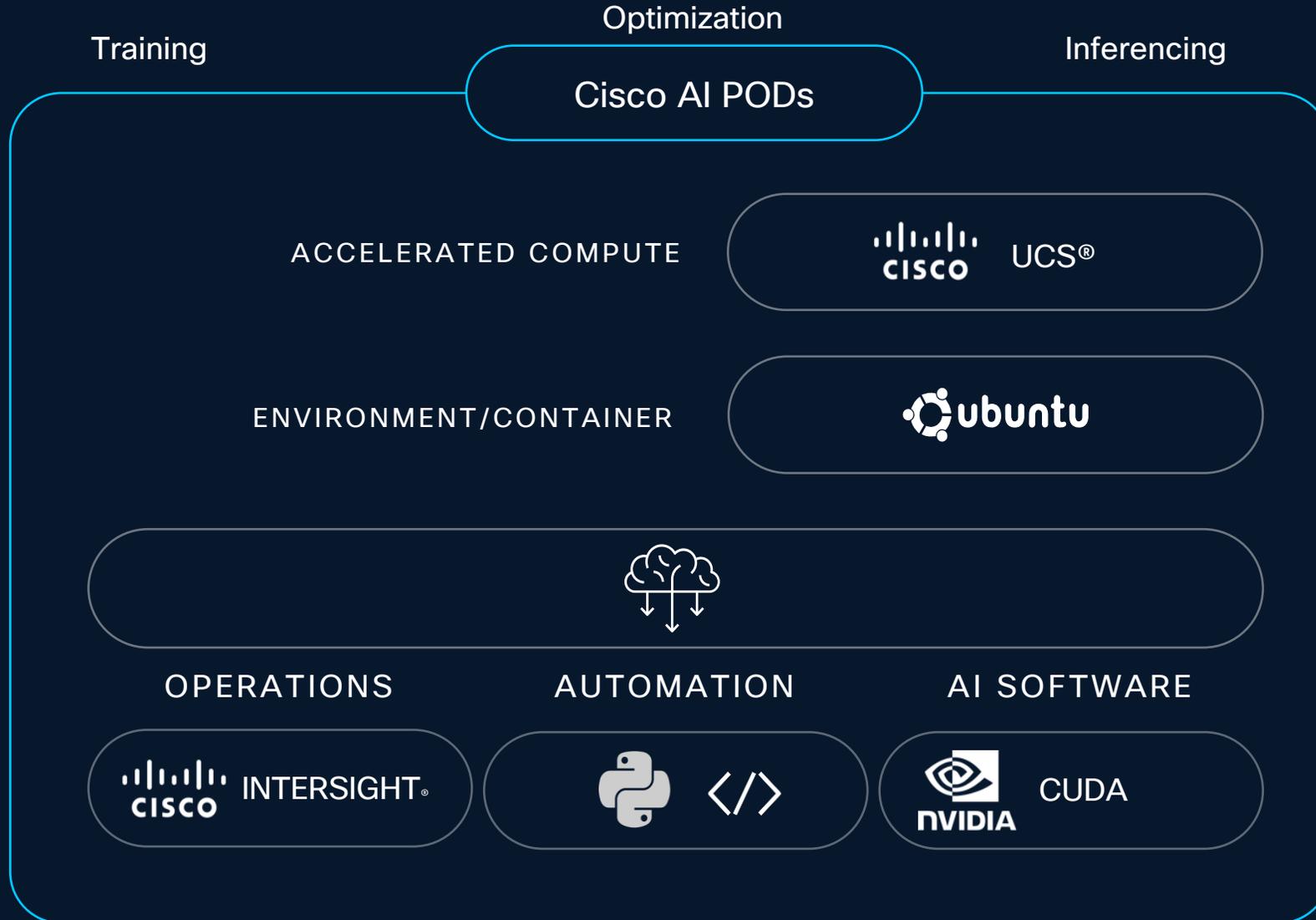
# Cisco AI PODs

## Introducing AI POD "Integrated Offerings"

Training    Optimization    Inferencing

**Cisco AI PODs**

**EXTEND TO**

CISCO® SECURITY

AI Defense    Hypershield

Firewall/Nexus® SmartSwitch

OBSERVABILITY

splunk> **Observability Cloud**

Open telemetry extensions

WORKLOAD MANAGEMENT & OPS

ISOVALENT    run:ai    slurm workload manager    splunk>

---

OPERATIONS    AUTOMATION    AI SOFTWARE

CISCO INTERSIGHT® & NEXUS DASHBOARD    A    </>    NVIDIA NIM Operator NeMo CUDA

KUBERNETES    RedHat OpenShift    ubuntu RANCHER BY SUSE

ACCELERATED COMPUTE    CISCO UCS®

HIGH-PERFORMANCE NETWORKING    CISCO NEXUS®

EXTEND TO STORAGE PLATFORM ECOSYSTEM    VAST    NetApp    PURE STORAGE    NUTANIX    HITACHI

---

ADVANCED SERVICES INCLUDED

Cisco **Customer Experience**    CX

DEVNET-2488

# Cisco AI PODs

Compute AI POD

Training

Optimization

Inferencing

Cisco AI PODs

ACCELERATED COMPUTE

CISCO UCS®

ENVIRONMENT/CONTAINER

ubuntu

OPERATIONS

AUTOMATION
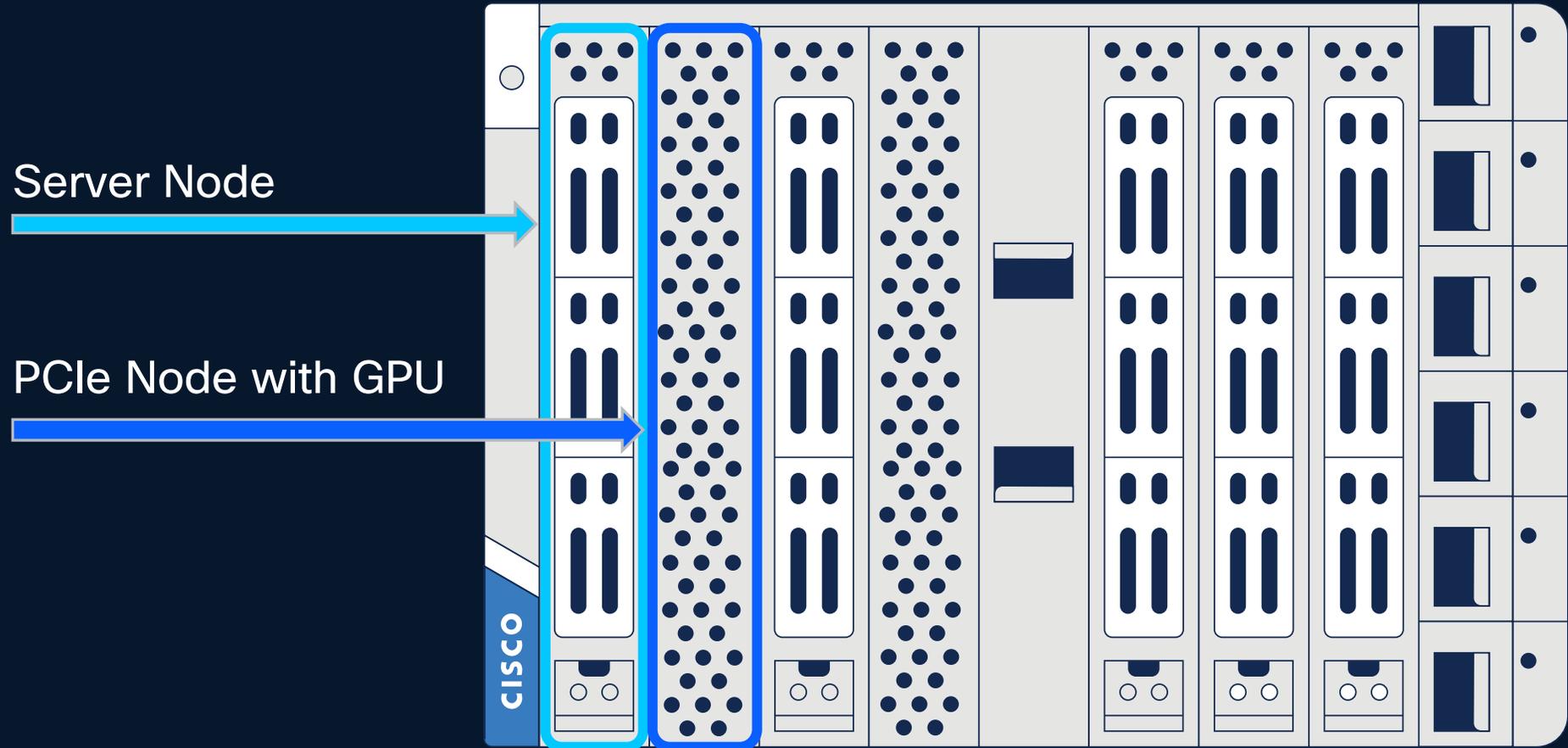
AI SOFTWARE

CISCO INTERSIGHT®

</>

NVIDIA CUDA

DEVNET-2488

# Cisco AI POD Compute Infrastructure

AI PODs can come in multiple types of UCS products :

• UCS C845A

• UCS C885A

• UCS C880A

• **UCS X-Series**

Server Node

PCIe Node with GPU

# What's needed to setup AI?

From a compute point of view

### Hardware:
- UCS Servers
- GPUs

### Software Solutions:
- **Intersight** for physical configuration
- Operating System
- Containers to host workloads
- GPU software stack

### Challenges:
- Respect Cisco Validated Designs
- Learn new technological dimensions
- Find proper AI use cases

# All those considerations will **delay the value realization**

# Intersight AI Bridge

A free community project

Complete **automated** build-up:

From **brand new AI Pod servers** to configured AI platform with **embedded use-cases**

- **Deploy** Intersight server configuration

- **Install** Operating System

- **Setup** OS and GPU configurations

- **Validate** with AI scenarios

mabuelgh / intersight-ai-bridge

## 📖 README

# Intersight AI Bridge [DEVNET] [published]

Intersight AI Bridge **simplifies and accelerates** the initial installation and usage of **AI workloads** such as Cisco AI Pods.

> **Note**: Starting from **Step 3**, these tools can also be used on any Linux system, even without Cisco UCS hardware.

This project provides scripts and configurations to:

1. Deploy a **Server Profile** on Cisco Intersight.
2. Install an **Operating System** through the Intersight OS Install feature (requires an *Advantage* license, otherwise can be done manually).
3. **Set up your environment for GPU-based AI workloads** with four possible use cases:
   - Chatbot with vLLM + OpenWebUI
   - Chatbot with Text Generation WebUI
   - Chatbot with vLLM + Retrieval-Augmented Generation (RAG)
   - Stresstest with vLLM

## Getting Started

# Intersight AI Bridge

Goals:

- Accelerate **time-to-value** for AI Pods

- Ensure consistent and **compliant deployment**

- Follow the **guidelines from Cisco Validated Designs**, NVIDIA and Ubuntu

- Show the automation capabilities of the products and solutions

Available:

- On GitHub, for free : https://github.com/mabuelgh/intersight-ai-bridge

DEVNET-2488

# Intersight AI Bridge

## Operational Workflow

**Deploy Server Config**

Create Intersight server configuration based on the **CVD** and deploy automatically using Python.

**Install OS**

Deploy and install the OS with **recommended releases.**
Use Intersight's OS Install feature to deploy **Ubuntu.**

**Setup system & NVIDIA software**

Use scripts to **install all the required system packages** to use NVIDIA hardware and toolkit on the OS.

**Validate AI scenario**

Intersight AI Bridge comes with embedded scripts to automatically deploy AI use-cases using Docker containers.

# Server Profile

UCS Server configuration in Intersight is simplified with full hardware abstraction.

Each server is configured with various policies detailling the configuration (BIOS, Boot, Disk, etc. and identity (IP, MAC, etc.) inside a **Server Profile.**

The **BIOS Policy** is customized to **handle AI workloads efficiently.**

BIOS Policy →

Boot Order Policy →

Power Policy →

IMC Access Policy →

vKVM Policy →

Local User Policy →

Storage Policy →

LAN Connectivity →

Intersight Server Profile

# Server Profile

UCS Server configuration in Intersight is simplified with full hardware abstraction.

Each server is configured with various policies detailling the configuration (BIOS, Boot, Disk, etc. and identity (IP, MAC, etc.) inside a **Server Profile.**

The **BIOS Policy** is customized to **handle AI workloads efficiently.**

Thanks to Intersight, every policy created can be **reuse for every server – less manual work – more compliance**

# Step 1 : Deploy Server Config

The first step is using Intersight API to **create, deploy and activate** a configuration on the server.

Automated process breakdown:

- Shape a configuration based on the <u>Cisco Validated Designs</u>

- Create the configuration (policies and profile) with Intersight

- Assign and Deploy the Server Profile to the server

➜ **The server is now configured**

Intersight Server Profile

DEVNET-2488

# Step 1 : Deep Dive

- JSON Template is created based on the **CVDs**

- Intersight Python SDK is used to deploy JSON config through **EasyUCS** (GitHub project)

- All the **configuration Policies** and **Server Profile** are created inside Intersight and deployed on desired Server



JSON Template of configuration

Policies created inside Intersight

Profile created and assigned in Intersight

DEVNET-2488

Demo

# Step 2 : OS Installation

The second step is also using Intersight API to install the OS on the server with the « **Install OS** » Intersight feature.
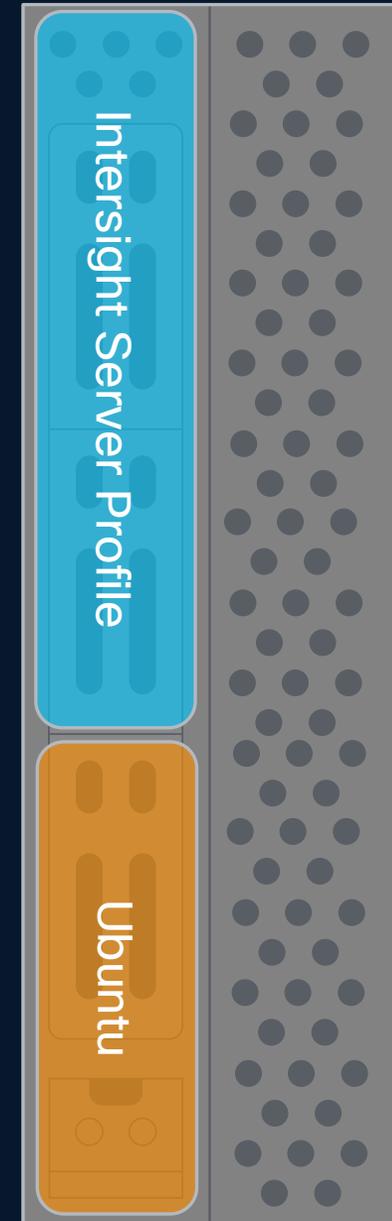
Manual process breakdown:

- Download the desired OS ISO file

- Download the recommended <u>SCU</u> (Server Configuration Utility) ISO file

Automated process breakdown:

- Create an OS configuration based on recommandations

- Install the OS through Intersight « Install OS » feature

➜ **The OS is installed on the server**

Intersight Server Profile

Ubuntu

# Step 2 : Deep Dive

- « **Install Operating System** » is a feature of Intersight to deploy automatically an OS through a user-friendly wizard

- Intersight **Python SDK** is used to do the process within a Python script instead of GUI

- Process can take some time and can be monitored using the UI



Install Operating System option in Intersight UI



Tracking Execution Flow in Intersight

Demo

# Step 3 : Setup system & NVIDIA software

The third step is about **deploying the packages, software and drivers** to use the GPUs.

Automated process breakdown:

- Install up to date packages

- Install NVIDIA drivers

- Install Docker and NVIDIA container toolkit

- Install NVIDIA CUDA toolkit

- Install Python & Hugging Face Hub

- Reboot the server to apply all the changes

➜ **The software stack is ready on the OS**

Intersight Server Profile

Ubuntu

**NVIDIA**
Software & Drivers

# Step 3 : Deep Dive

- Setup shell script will download and setup everything needed to use the GPUs, it includes:

  - **CUDA** (Compute Unified Device Architecture) is a toolkit from NVIDIA to allow developers to use GPUs for processing in addition to graphics

  - **Hugging Face Hub** is a Python package that provides an interface to Hugging Face and access thousands of pre-trained models

```
#-------------------------------------------------------#
# Updating the package list
#-------------------------------------------------------#

Scanning processes...
Scanning processor microcode...
Scanning linux images...

Running kernel seems to be up-to-date.

The processor microcode seems to be up-to-date.

No services need to be restarted.

No containers need to be restarted.

No user sessions are running outdated binaries.

No VM guests are running outdated hypervisor (qemu) binaries on this host.
Hit:1 http://security.ubuntu.com/ubuntu noble-security InRelease
Hit:2 http://archive.ubuntu.com/ubuntu noble InRelease
Hit:3 http://archive.ubuntu.com/ubuntu noble-updates InRelease
Hit:4 http://archive.ubuntu.com/ubuntu noble-backports InRelease
Reading package lists... Done
Reading package lists... Done
Building dependency tree... Done
Reading state information... Done
Calculating upgrade... Done
The following upgrades have been deferred due to phasing:
  dhcpcd-base
The following packages
  apparmor cloud-init                         1.2-packagekitglib-1.0
```

Setup everything automatically

```
## Docker is installed and is running successfully.


#-------------------------------------------------------#
# Running nvidia-smi to verify Nvidia drivers installation
#-------------------------------------------------------#

Fri Dec 12 15:31:52 2025
+-----------------------------------------------------------------------------+
| NVIDIA-SMI 580.95.05        Driver Version: 580.95.05      CUDA Version: 13.0 |
+-------------------------------+----------------------+----------------------+
| GPU  Name          Persistence-M | Bus-Id        Disp.A | Volatile Uncorr. ECC |
| Fan  Temp    Perf  Pwr:Usage/Cap |         Memory-Usage | GPU-Util  Compute M. |
|                                  |                      |               MIG M. |
|==============================+======================+======================|
|   0  NVIDIA L40S          Off | 00000000:3D:00.0 Off |                    0 |
| N/A   26C    P8     25W / 350W |      0MiB / 46068MiB |      0%      Default |
|                                  |                      |                  N/A |
+-------------------------------+----------------------+----------------------+
|   1  NVIDIA L40S          Off | 00000000:E1:00.0 Off |                    0 |
| N/A   27C    P8     24W / 350W |      0MiB / 46068MiB |      0%      Default |
|                                  |                      |                  N/A |
+-------------------------------+----------------------+----------------------+

+-----------------------------------------------------------------------------+
| Processes:                                                                  |
|  GPU   GI   CI          PID   Type   Process name            GPU Memory     |
|        ID   ID                                               Usage          |
|==============================================================================|
|  No running processes found                                                 |
+-----------------------------------------------------------------------------+

## Nvidia drivers in
```
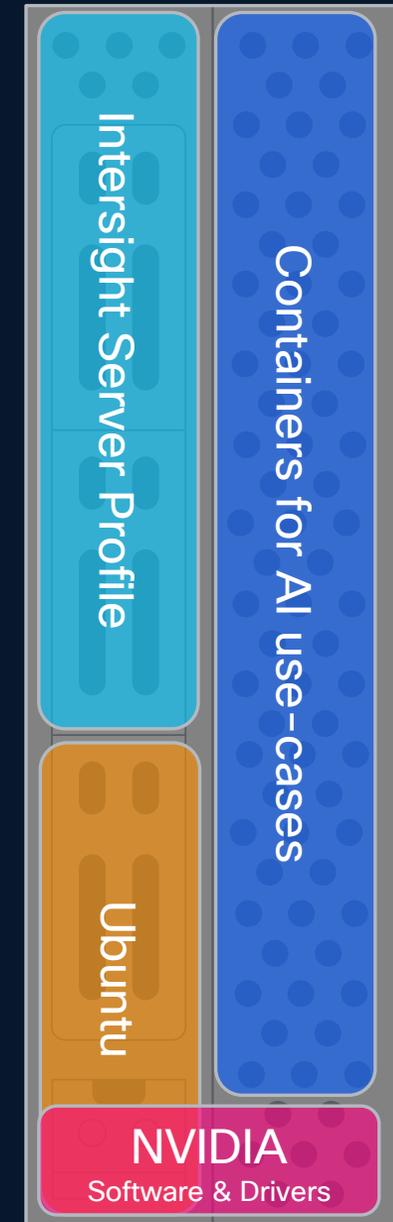
Script to check the setup

# Demo

# Step 4 : Validate with AI use-cases

The last step is about **choosing from a list of scenarios** the AI use-cases that best fit the usage needed.

Automated process breakdown:

- Download the recommended LLM (Large Language Model) on Hugging Face

- Use Docker compose on the pre-established scenarios

- Deploy container

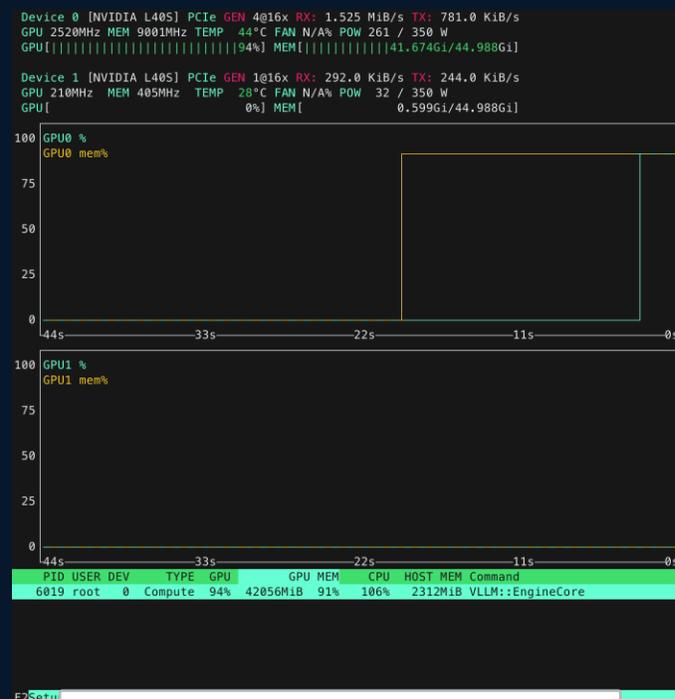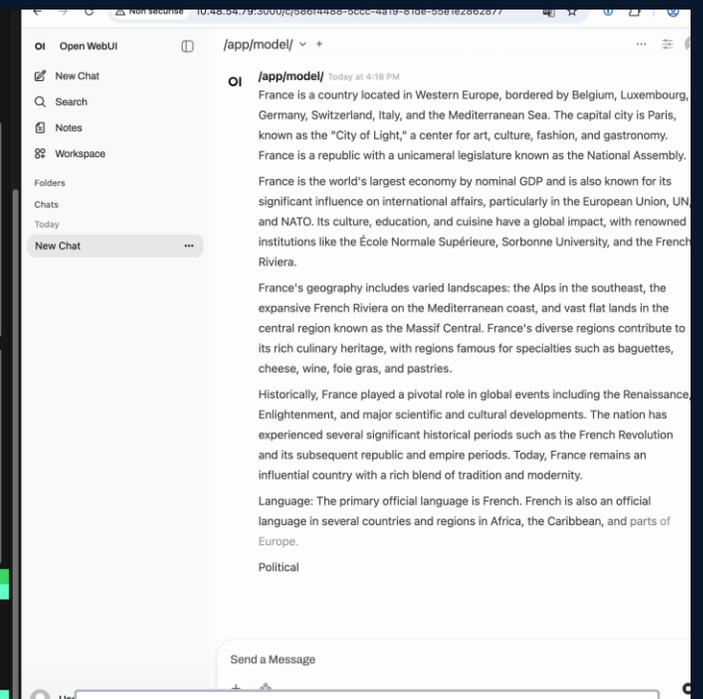➜ **The GPUs are now being used for AI use-cases**

Intersight Server Profile

Containers for AI use-cases

Ubuntu

NVIDIA
Software & Drivers

# Step 4 : Deep Dive

- **Pre-assembled scenarios using GPUs** for AI use-cases are available with shell scripting
  - Use **HuggingFace Hub** to download the LLM and Docker to **host and deploy** the workload
  - Deploy **vLLM**: a fast and easy-to-use library for LLM inference and serving with API capabilities
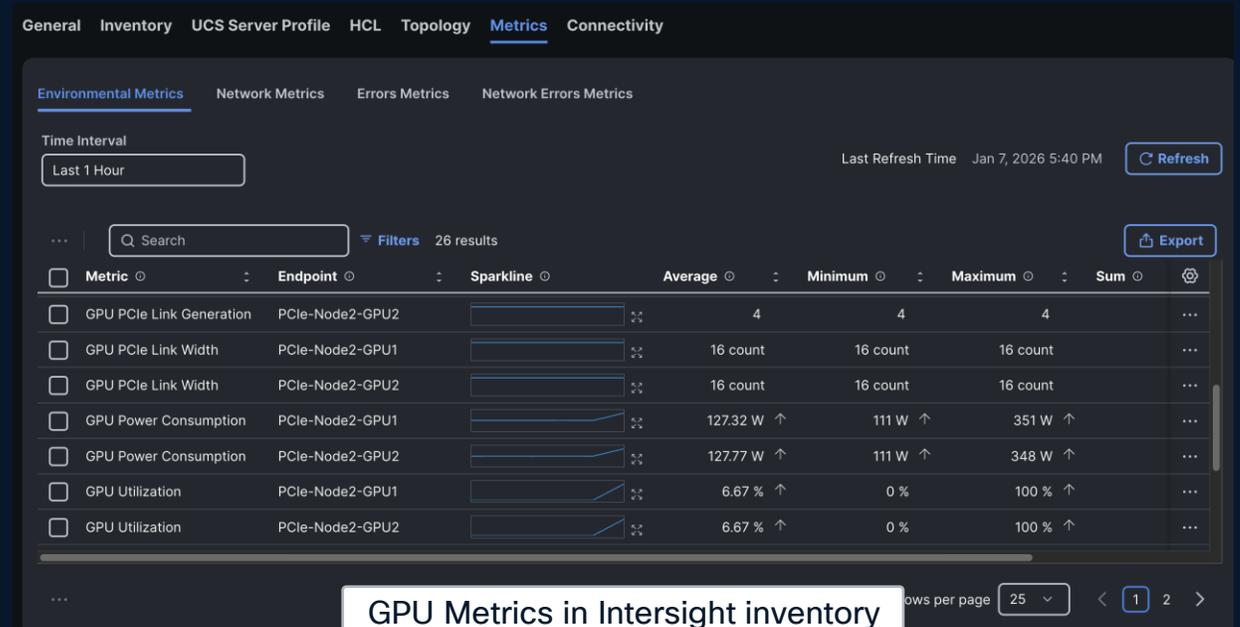


Launching scenario 2



Monitoring GPU Utilization



Chatbot using Web UI and LLM's API

Demo

# Intersight observability on GPU

Intersight can **display and explore the metrics** related to the **GPUs**:

- GPU Clockspeed
- GPU Memory Clockspeed
- GPU Power Consumption
- GPU Utilization
- And more!



GPU Metrics in Intersight inventory per server

# Intersight observability on GPU

Intersight can **display and explore the metrics** related to the **GPUs**:

- GPU Clockspeed
- GPU Memory Clockspeed
- GPU Power Consumption
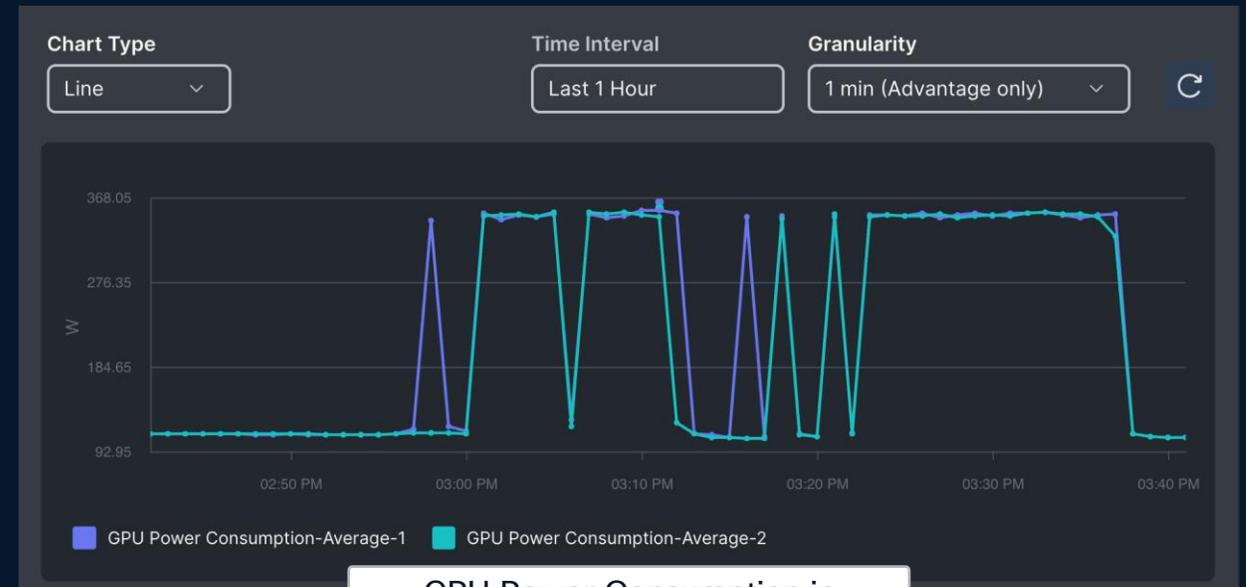- GPU Utilization
- And more!

without Operating System agent



| Chart Type | Time Interval | Granularity |
|---|---|---|
| Line | Last 1 Hour | 1 min (Advantage only) |

GPU Power Consumption-Average-1  GPU Power Consumption-Average-2

GPU Power Consumption in Intersight Metrics Explorer

CISCO

# Key Takeaways - Call to Action

Every step of an AI deployment can be **automated** and **accelerated**:

- Intersight configurations and OS deployments done with **API**

- All the OS requirements achieved **automatically**

- Get GPU observability **easily with Intersight**

**Intersight AI Bridge** is a **DevNet** community project:

- **Free**, available on <u>GitHub</u> and <u>DevNet</u> website

- Each step can work **independently**

- Very **modular**: Easily customisable

- **Create and share** new AI scenario with the community

➜ **From weeks of learning, work and experimentation to 1.5 hours of automated deployment**
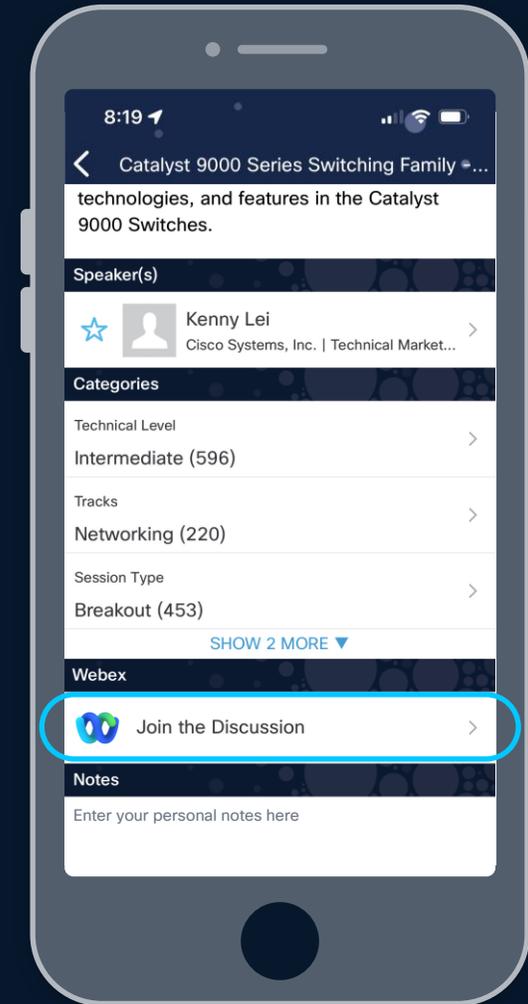
# Webex App

**Questions?**

Use Webex app to chat with the speaker after the session

**How**

(1) Find this session in the Cisco Events 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 27, 2026.**

https://ciscolive.ciscoevents.com/
ciscolivebot/#DEVNET-2488

DEVNET-2488

# Complete your session surveys

**Complete your surveys** in the Cisco Events app.

**Complete** a minimum of 4 session surveys and the overall event survey to receive a unique Cisco Live t-shirt.

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

DEVNET-2488

# Continue your education

**Visit** the Cisco Showcase for related demos

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

**Visit** the Technical Solutions Clinics to discuss your technical questions

**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

**Contact me at**: mabuelgh@cisco.com

Thank you

CISCO