# GPU Profiles for Hyperscale Use Cases

**(On-behalf of OCP GPU Management Interface Work Group)**

Sivakumar Sathappan,  AMD

Hari Ramachandran,  Microsoft
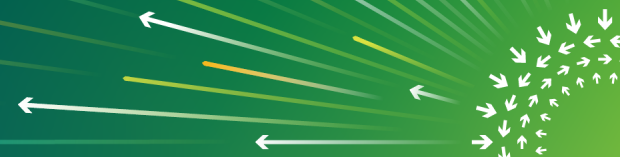
Linda Wu,  NVIDIA

# OCP GPU Management Interface

- The OCP GPU Management Interface Working Group represents a remarkable coalition of industry leaders including **AMD, Google, Meta, Microsoft, and NVIDIA**. This working group exemplifies the spirit of collaboration and innovation, as each member has contributed equally to the development of a standardized Redfish Interoperability Profile and a comprehensive roadmap.

- The presentation would showcase
  1. Redfish Interoperability Profiles helps to test the compliance for Hyperscale's
  2. GPU Management Interface contribute to Redfish Interoperability Validator
  3. Adapting Redfish Interoperability Profiles and how challenges were mitigated

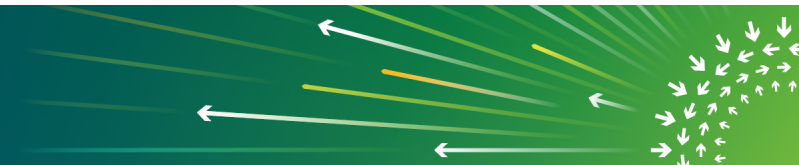# Motivation for the OCP GPU Management Working Group

**Suppliers**

- Converge customer management requirements to reduce engineering costs and speed time-to-market

**Hyperscalers**

- Rapidly adopt and deploy new GPU designs to our fleets

- Drive consistency across vendors

- Increase consistency between generations to streamline management at scale
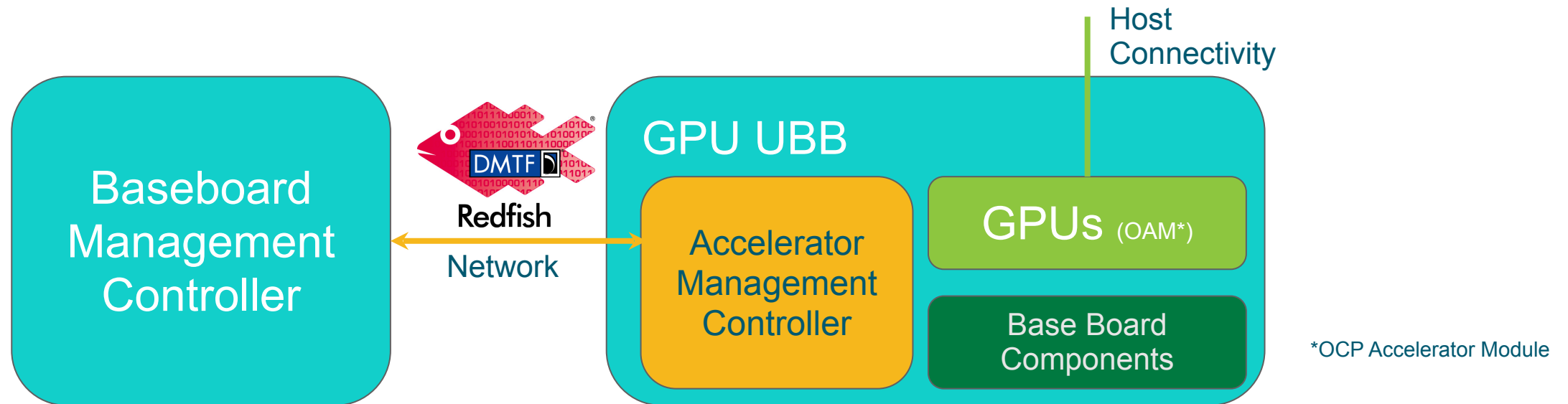
**Solution**

- Define profiles built on industry standards

# Anatomy of a GPU Universal Base Board

- <u>OCP HGX</u> defines the hardware for a UBB for GPUs

- UBB can be described as tray which consists of a baseboard and a set of Accelerator Modules (OAM) with the necessary high-speed interconnect and other I/O devices

- As per OCP HGX, GPU UBBs contain an Accelerator Management Controller (AMC)

- The UBB AMC is managed directly via Redfish interface to an enclosure management controller



*OCP Accelerator Module

# Standardizing Hyperscalers Requirements for Accelerators

## Please attend the OCP presentation

**Standardizing Hyperscaler Requirements for Accelerators**
SJCC - Concourse Level - 210CG

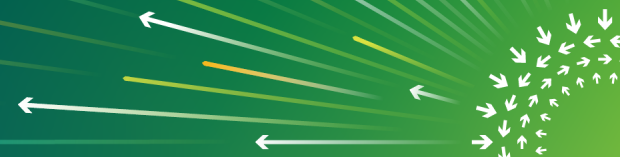By **Krishna Sugumaran, Rama Bhimanadhuni, Anil Agrawal, Sujoy Sen, Vishal Jain**

Wed, October 16, 9:00am - 9:30am | SJCC - Concourse Level - 210CG

## Standardizing Hyperscaler Requirements for Accelerators

- Hardware Management

With the advent and explosion of AI/ML, GPUs are playing an increasingly important role in hyperscale data centers. Rapid innovation in GPU technology and the evolving demands of workloads necessitate swift adoption of new GPU designs by hyperscalers. The current lack of standardization in GPU management, coupled with variations in hyperscalers' specific requirements, make the implementation and onboarding of new GPUs time-consuming.

In this presentation, AMD, Google, Meta, Microsoft, and NVIDIA will discuss how DMTF industry standards, such as MCTP, SPDM, PLDM, and Redfish are being used to manage GPU devices deployed by hyperscalers. We will provide an update on our work to define standards-based profiles for GPU management, aimed at accelerating the implementation and adoption of new designs. Additionally, we will provide an update on how this OCP work has improved GPU management over the last 12 months.

# GPU Profile Use Cases

**Standard Requirements**
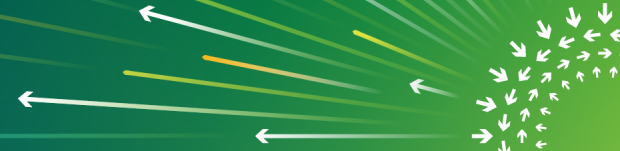
- Standardize (*Required, Recommended and Optional*) Redfish Resources and Properties
  - Provided by GPU platform vendors
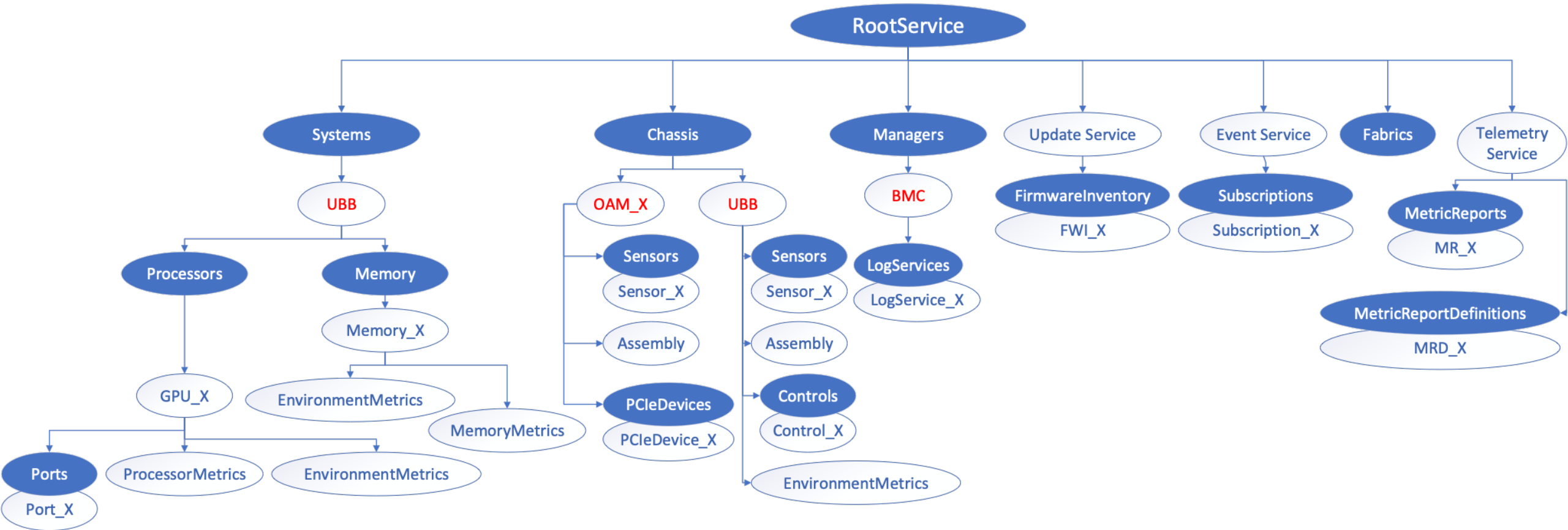  - Required by Hyperscalers

**Validation**

- Vendors can utilize the Redfish Interoperability Validator to provide compliance to the RIP

**Compliance**

- Hyperscalers can utilize the Redfish Interoperability Validator to test compliance against RIP

# UBB Management Redfish Model

# GPU Management Redfish Interop Profile

- Define the <u>baseline</u> set of Redfish Resources and Properties that are required by Hyperscalers that shall be provided across all GPU vendors

- GPU Management Interoperability Profile Resources

| Assembly | LogEntry | MetricReport | PortMetrics | TelemetryService |
|---|---|---|---|---|
| Certificate | LogService | MetricReportDefinition | Processor | ThermalMetrics |
| Chassis | Manager | PCIeDevice | Sensor | ThermalSubsystem |
| ComputerSystem | Memory | PCIeFunction | ServiceRoot | UpdateService |
| EnvironmentMetrics | MemoryMetrics | Port | SoftwareInventory | |

<u>OCP_UBB_BaselineManagement_v1_0_0.json</u>

# Redfish Interop Validator Tool

- Tool is critical for validation and compliance

- Tool is part of a broader Initiative CTAM (Compliance Tool for Accelerator Manageability)

CTAM - Compliance Tool for Accelerator Manageability
SJCC - Concourse Level - 210CG
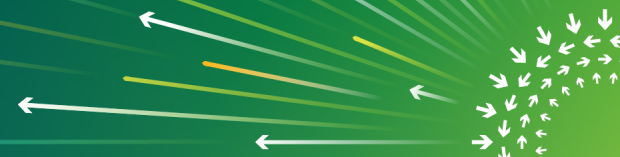
By **Rajat Madhusudan, Venkat Ramesh, Afsana Chowdhury**

Wed, October 16, 4:05pm - 4:25pm | SJCC - Concourse Level - 210CG

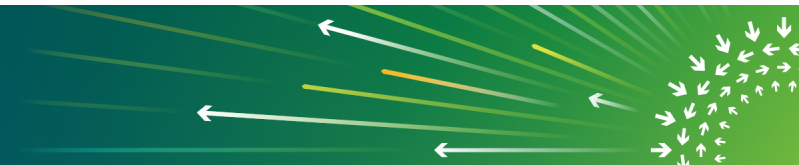## CTAM - Compliance Tool for Accelerator Manageability

● **Test & Validation**

The emergence of AI/ML has thrust AI accelerators, such as GPUs, into a pivotal role within hyperscale data center operations. In the past year, the OCP community has put forth considerable effort to unify AI accelerator management interfaces and standardize manageability workflows for hyperscalers, including firmware updates and RAS flows. To aid the industry in complying with these specifications, the Compliance Tool for Accelerator Management (CTAM) was introduced.

This presentation will navigate us through CTAM's evolution, shedding light on future developments. We will examine the underlying philosophy, architecture, and processes designed to align CTAM with the related OCP specifications. A recorded demo will provide a glimpse into the community's collaborative achievements, and the feedback gathered during this session will fuel further enhancements. CTAM is poised to become a fundamental component in the operational toolkit for GPU suppliers and hyperscalers alike.

# Adoption Challenges

- Collating all the properties and deciding the scope under "Mandatory", "Recommended and "Supported"

- Defining "Use Cases" without adding to the RIP complexity.

Thank You