

Hewlett Packard Enterprise

### Updated Node Power Management For HPE Cray EX255a and EX254n Blades

Brian Collum, HPE May 2024



CUG 2024 – © COPYRIGHT 2024 HEWLETT PACKARD ENTERPRISE

#### Abstract

Cray EX nodes have always supported a form of power capping that would allow customers to lower power usage of specific nodes as desired. With the introduction of the HPE Cray EX254n (NVIDIA Grace Hopper) and HPE Cray EX255a (AMD MI300a), this became critical as the overall rack power pushed beyond the maximum supported at some customers sites. With the HPE Cray EX254n in particular, the total TDP of the modules exceeds the maximum that can be delivered by the Cray EX infrastructure. This drove the decision to have to set a power limit on the Grace Hopper modules by default (a first for Cray EX). This presentation will walkthrough design goals of the blades, how power capping is implemented in the firmware, and how to configure the power limit in a running system. The presentation will also go through how to view the current limits configured via the node controllers Redfish API and in-band tools, where applicable, and how the in-band tools interact with the out-of-band configurations.

#### Agenda

#### • Redfish APIs

- Node Power Limit
- Accelerator Power Limit
- In-band APIs
  - PM Counters
  - nvidia-smi
  - rocm-smi
- HPE Cray EX255a
- HPE Cray EX254n



#### **Redfish – Node Power Limit**

- Configures the maximum input power the node to draw from the chassis
- Spikes may be seen above SetPoint value may be seen occasionally but on average will be below

```
curl https://<nC>/redfish/v1/Chassis/Node0/Controls/NodePowerLimit
  "ControlMode": "Automatic",
  "Sensor": {
    "DataSourceUri": "/redfish/v1/Chassis/Node0/Sensors/ChassisVoltageRegulator0InputPower",
    "Reading": 614
 },
  "SetPoint": 2500,
  "SettingRangeMax": 2852,
  "SettingRangeMin": 1600,
  "Status": {
    "Health": "OK"
```

Example from EX254n

#### **Redfish – Accelerator Power Limit**

- Configures the maximum power to be used by the given Accelerator
- Available on GPU Compute Blades
- Separate from the NodePowerLimit, but NodePowerLimit takes precedence

```
curl https://<nC>/redfish/v1/Chassis/Node0/Controls/Accelerator0PowerLimit
{
    ...
    "ControlMode": "Automatic",
    ...
    "SetPoint": 400,
    "SettingRangeMax": 900,
    "SettingRangeMin": 100,
    "Status": {
        "Health": "OK"
    }
}
```

Example from EX254n

#### **In-band Tools**

#### • PM Counters

- Read-only
- /sys/cray/pm\_counters/power\_cap
- /sys/cray/pm\_countes/accel[0-3]\_power\_cap
- nvidia-smi
  - <u>https://developer.nvidia.com/system-management-interface</u>
- rocm-smi
  - <u>https://rocm.docs.amd.com/en/latest/</u>



Nodes

APUs

• 4x AMD MI300a APUs



#### HPE Cray EX255a – Redfish NodePowerLimit

- Evenly distributed across the four APUs
- Power Limit is configured per APU via AMD's APML 'Set Package Power Limit'
- When disabled, the limit is removed from the APU

# $APUPowerLimit = \frac{((NodePowerLimit * EffeciencyFactor) - OtherNodePower)}{NumberOfAPUs}$

#### HPE Cray EX255a – View In-Band APU Limit

- Can read current limit with
- # rocm-smi --showmaxpower

================	======================================
=================	======================================
GPU[0]	: Max Graphics Package Power (W): 550.0
GPU[1]	: Max Graphics Package Power (W): 550.0
GPU[2]	: Max Graphics Package Power (W): 550.0
GPU[3]	: Max Graphics Package Power (W): 550.0
=======================================	
	======================================

- Does not reflect OOB limit
- Lower limit wins (OOB vs. In-band)!

#### HPE Cray EX255a – Set In-Band APU Limit

- Can set module limit to a value <u>lower</u> than OOB limit with
  - rocm-smi --setpoweroverdrive <power-limit>
- Should be able to read it back as well if successful
- # rocm-smi --setpoweroverdrive 277
- # rocm-smi --showmaxpower

=======================================	======================================
=======================================	======================================
GPU[0]	: Max Graphics Package Power (W): 277.0
GPU[1]	: Max Graphics Package Power (W): 277.0
GPU[2]	: Max Graphics Package Power (W): 277.0
GPU[3]	: Max Graphics Package Power (W): 277.0
	======================================

• Lower limit wins (OOB vs. In-band)!

Nodes Modules

# HPE Cray EX254n4x NVIDIA GH200Superchip Modules



#### HPE Cray EX254n – Redfish NodePowerLimit

- Evenly distributed across the four modules
- Power Limit is configured per module via NVIDIA SMBPBI's 'Set Total Module Power Limit'
- When disabled, the maximum supported value is configured on the modules
  - This is the only Cray EX blade that sets a limit when NodePowerLimit is disabled
  - i.e. effectively NodePowerLimit.SetPoint == 2852W (660.12W per module)
- Power is managed by the GPU, preference is given to the CPU and the GPU is dynamically capped to keep module power under the configured limit
- Takes precedence over Accelerator[0-3]PowerLimit

$$ModulePowerLimit = \left(\frac{\left((NodePowerLimit * EfficiencyFactor) - OtherNodePower\right)}{NumberOfModules}\right) * LimitAccuracy$$

#### HPE Cray EX254n – View In-Band Module Limit

• Can read current limit with

nvidia-smi -q

#### • Shows a LOT of data, but in there is this for each module:



• Lower limit wins (OOB vs. In-band)!

#### HPE Cray EX254n – Set In-Band Module Limit

Can set module limit to a value lower than OOB limit with

nvidia-smi --power-limit <power-limit> --scope 1

• After setting, able to read it back



• Lower limit wins (OOB vs. In-band)!

#### HPE Cray EX254n – Accelerator[0-3]PowerLimit

- Power Limit applied to individual Hopper GPU
- Power Limit is configured per module via NVIDIA SMBPBI's 'Set Total GPU Power Limit'
- When disabled, the maximum supported value is configured on the GPU
  - i.e. effectively Accelerator[0-3]PowerLimit.SetPoint == 900W
- Module Limit (from NodePowerLimit) takes precedence, so higher value here has little meaning

#### HPE Cray EX254n – View In-Band GPU Limit

• Can read current limit with

nvidia-smi -q

#### • Shows a LOT of data, but in there is this for each GPU:



- Note: Can be misleading if Module Power Limit is set lower.
- Lower limit wins!

#### HPE Cray EX254n – Set In-Band GPU Limit Can set module limit to a value lower than OOB limit with nvidia-smi --power-limit <power-limit> Should be able to read it back as well if successful GPU Power Readings : 92.29 W Power Draw : 500.00 W Current Power Limit **Requested In-Band** Requested Power Limit : 500.00 W Limit Default Power Limit : 900.00 W Min Power Limit. : 100.00 W Max Power Limit : 900.00 W

• Lower limit wins!

#### **Management Systems**

- HPCM and CSM both support reading/setting
- CLI and API support
- Utilize same Redfish APIs

#### • HPCM:

- mpower -n <nids> --get-limit
- mpower -n <nids> --set-limit <power>

#### • CSM:

- cray capmc get\_power\_cap create --nids <nids>

## Thank you

CUG 2024 – © COPYRIGHT 2024 HEWLETT PACKARD ENTERPRISE