

# Adastra energy

**B. Cirou, Journée « Energie, infrastructure et calcul, quel équilibre ? »**

**Montpellier - 27 mars 2025**

# Layout

**Adastra Architecture**

**Metrics**

**Data**

**Work in progress**

# ADASTRA Architecture

## HPE Cray EX4000

- 544 Genoa nodes

2 x 9654 @ 2.4GHz

- 356 MI250X nodes

8 x GFX90A @ 1.5GHz

3<sup>rd</sup> Green 500 (2023)

- 28 MI300A nodes

4 x GFX942 @ 1.8GHz

3<sup>rd</sup> Green 500 (2024)

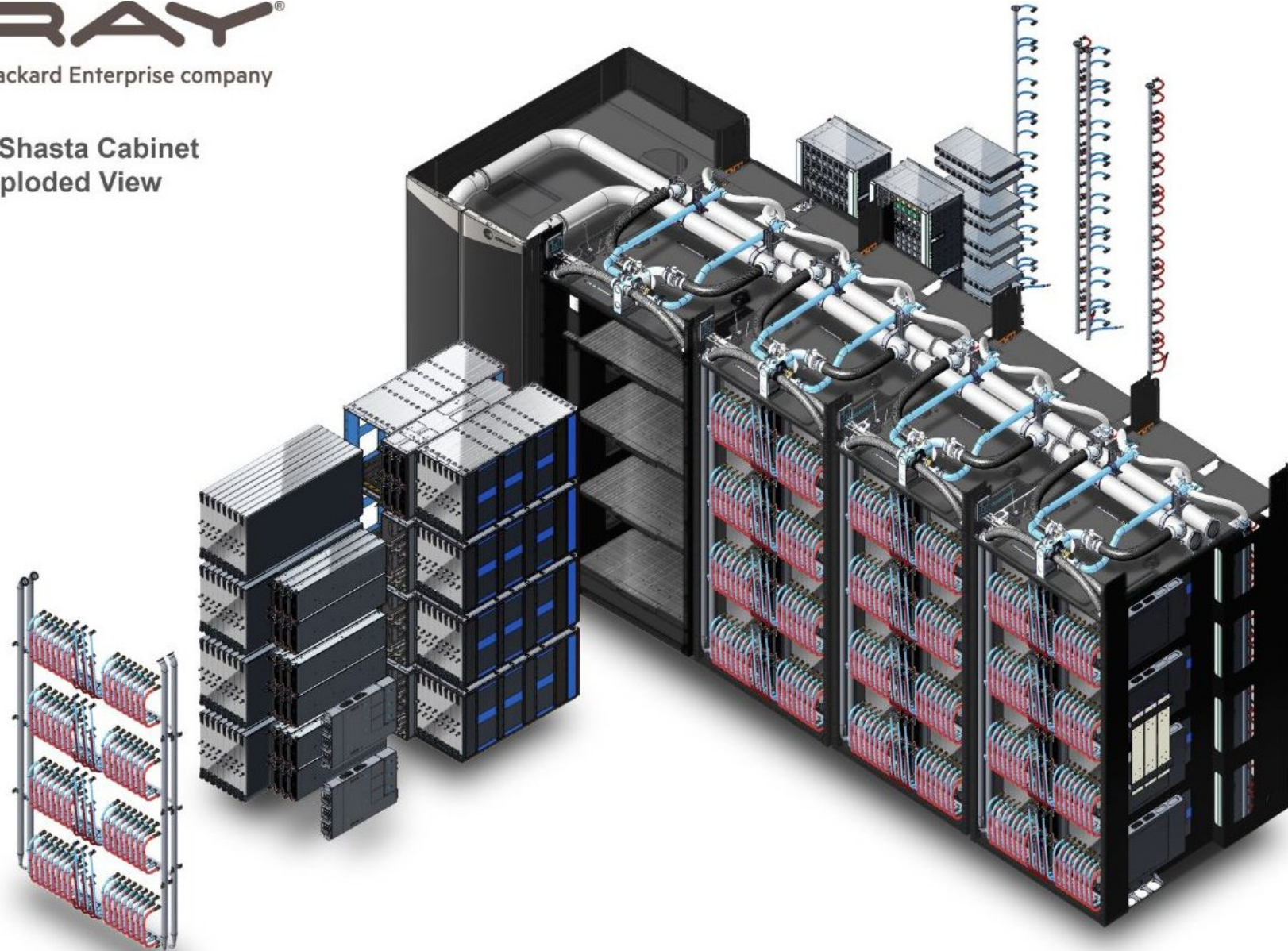


# ADASTRA Architecture

**CRAY**<sup>®</sup>

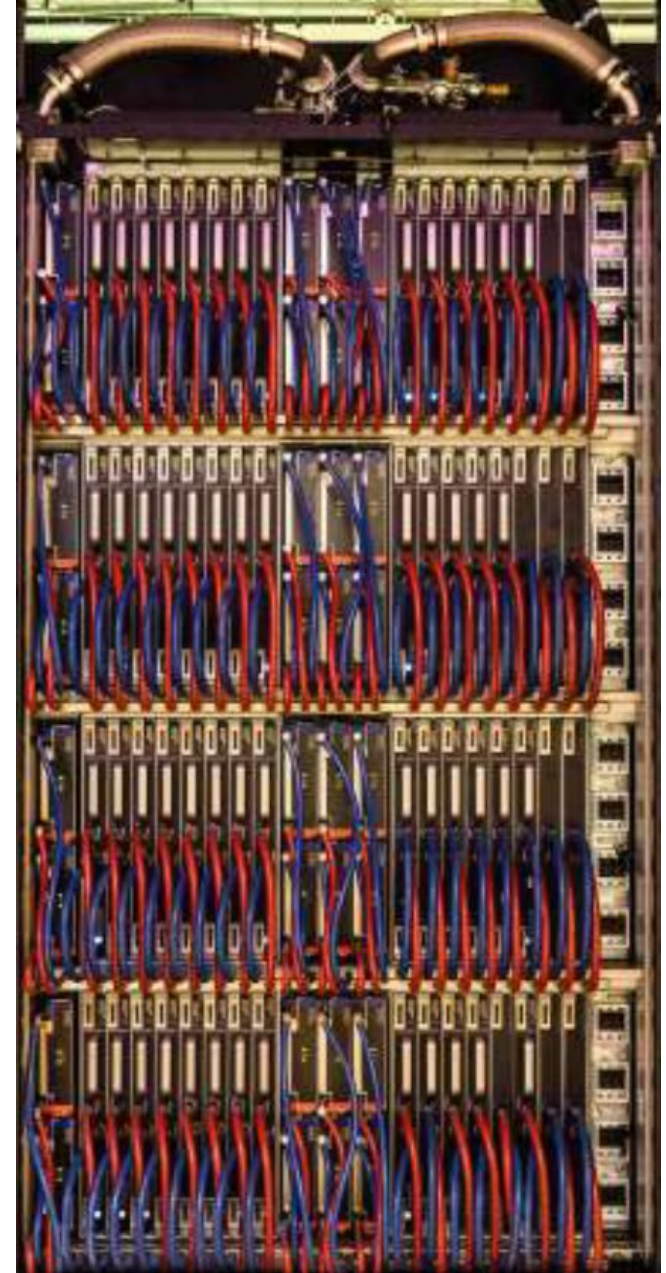
a Hewlett Packard Enterprise company

Cray Shasta Cabinet  
Exploded View



## Cooling

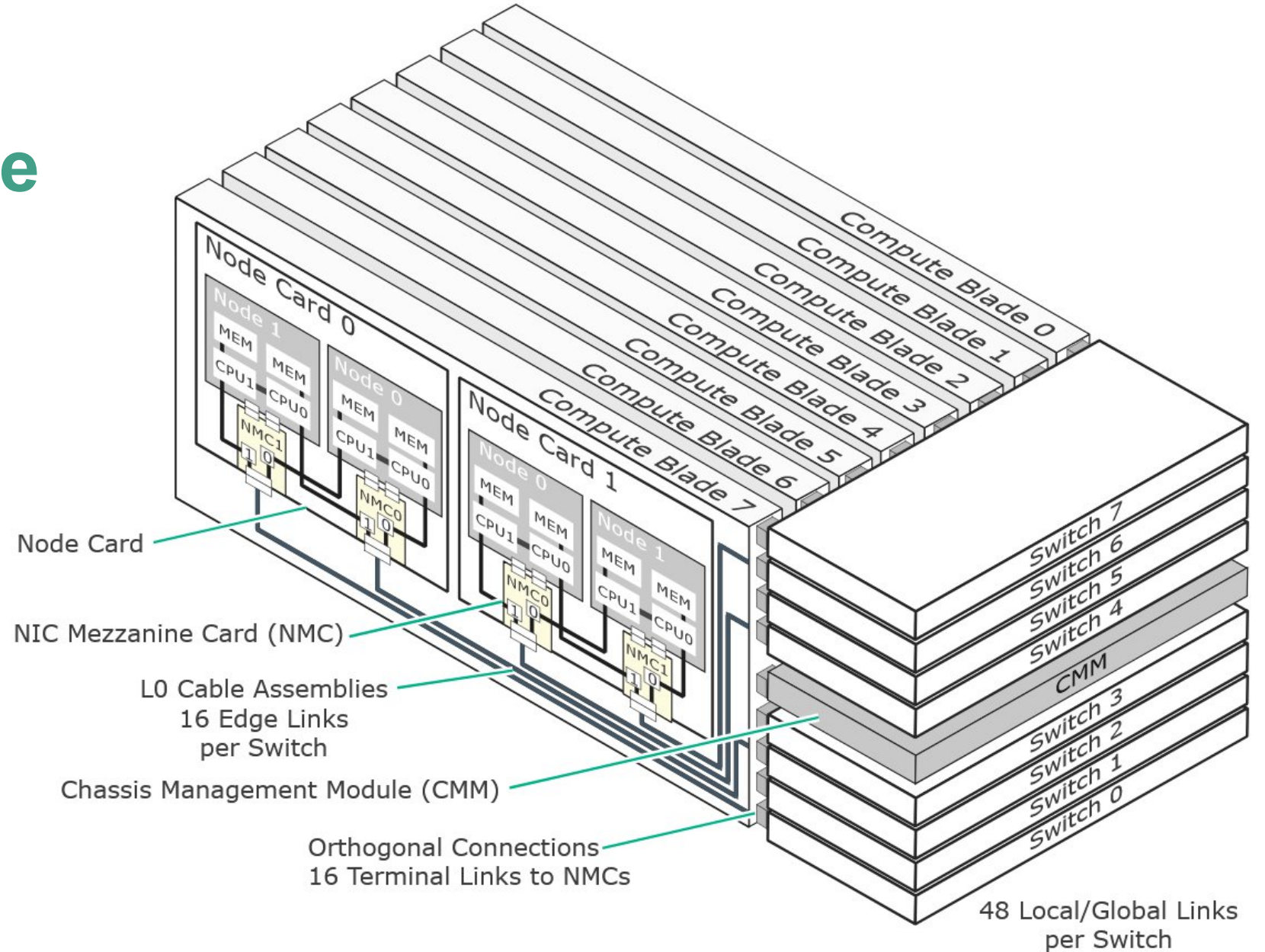
- **97% watercooled**
- **1493 kW peak** (MI250X + GENOA blades)
- **CDU Inlet temp: 31°C**
- **CDU Outlet temp: 42°C**
- **Flow: 90 m<sup>3</sup>/h**
- **Pressure drop: 160 kPa**



# ADASTRA Architecture

## Genoa enclosure

### Direct Liquid Cooling





# METRICS

## Monitoring: Cray HPCM (Kibana)





## Admin tools : Cray HPCM

### Get Node Energy

```
ncn-w001# cray capmc get_node_energy create --nids NID_LIST \ --start-time '2020-03-04 12:00:00' --end-time '2020-03-04 12:10:00' --format json
```

### Get Node Energy Stats

```
ncn-w001# cray capmc get_node_energy_stats create --nids NID_LIST \ --start-time '2020-03-04 12:00:00' --end-time '2020-03-04 12:10:00' --format json
```

### Get Node Energy Counter

```
ncn-w001# cray capmc get_node_energy_counter create --nids NID_LIST \ --time '2020-03-04 12:00:00' --format json
```

### Get Node Power Control and Limit Settings

```
ncn-w001# cray capmc get_power_cap create --nids NID_LIST \ --format json
```

### Get System Power

```
ncn-w001# cray capmc get_system_power create \ --start-time '2020-03-04 12:00:00' --window-len 30 --format json
```

```
ncn-w001# cray capmc get_system_power_details create \ --start-time '2020-03-04 12:00:00' --window-len 30 --format json
```

### Get Power Capping Capabilities

```
ncn-w001# cray capmc get_power_cap_capabilities create --nids NID_LIST \ --format json
```

### Set Node Power Limit

```
ncn-w001# cray capmc set_power_cap create --nids NID_LIST \ --node 225 --format json
```

### Remove Node Power Limit (Set to Default)

```
ncn-w001# cray capmc set_power_cap create --nids NID_LIST \ --node 0 --format json
```

## Node metrics: Cray PM Counters

`/sys/cray/pm_counters`

- **power**: Point-in-time power (Watts).
- **energy**: Accumulated energy, in joules.
- **cpu\_power**: Point-in-time power (Watts) used by the CPU domain.
- **cpu\_energy**: The total energy (Joules) used by the CPU domain.
- **cpu\_temp**: Temperature reading (Celsius) of the CPU domain.
- **memory\_power**: Point-in-time power (Watts) used by the memory domain.
- **memory\_energy**: The total energy (Joules) used by the memory domain.
- **accel\_energy**: Accumulated accelerator energy (Joules).
- **accel\_power**: Accelerator point-in-time power (Watts).
- **raw\_scan\_hz**: The power management scanning rate for all data in pm\_counters.  
[...]

## Job scheduler: SLURM

### SLURM configuration

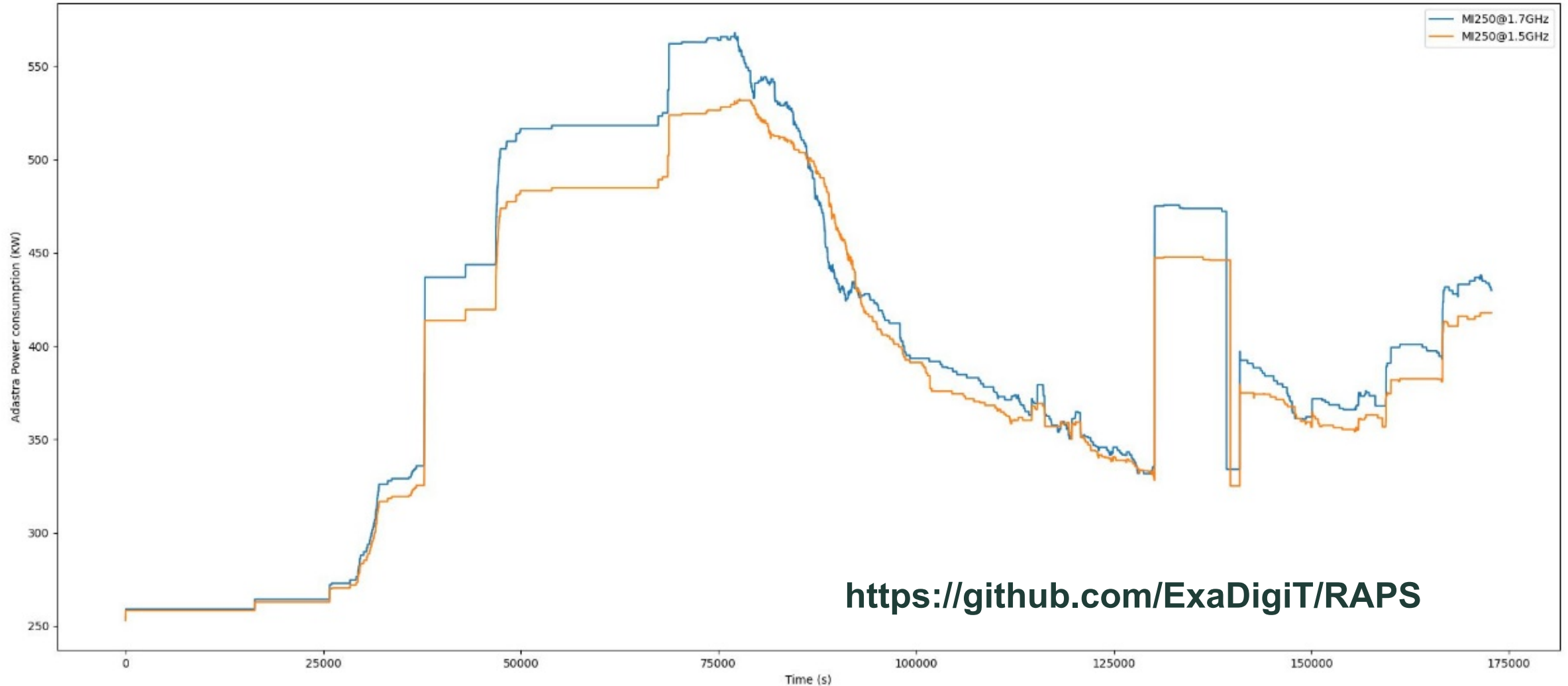
*AcctGatherEnergyType=acct\_gather\_energy/pm\_counters*  
*JobAcctGatherFrequency=task=30,energy=30,filesystem=30*

### Frequency and power capping

- Defaults on ADAstra MI250 nodes
  - Turbo ON on Trento CPU
  - Frequency capping to 1500 MHz on MI250x GPU
- Special SLURM plugin for dedicated job tuning (SPANK)
  - **Power capping management**
    - MI250X : 100 - 560W
    - MI300A : 100 - 550W
    - #SBATCH --gpu-power-cap=430
  - **Frequency capping management**
    - MI250X : 500 - 1700 MHz
    - MI300A : 500 - 1800 MHz
    - #SBATCH --gpu-srange=500-1400

# METRICS

## PM\_counter based (RAPS simulation) Theoretical pure MI250X racks



<https://github.com/ExaDigiT/RAPS>

## Job scheduler: SLURM

- Consumption of current running job :

```
sstat -j <jobid> --format=JobID,ConsumedEnergy -P JobID | ConsumedEnergy  
<jobid>.<current step> | 3.20K
```

- Consumption of ended jobs :

```
sacct -j <jobid> --format=JobID,ConsumedEnergy -P JobID | ConsumedEnergy  
<jobid>.<current step> | 3.20K
```

- Job report :

*CINES Job Report:*

-----

*o Estimated energy consumption: 254076 Joules  
(representing ~ 27% of the maximum nodes utilization)*

## **CINES Data Cluster (ELK, VictoriaMetrics, ...)**

- **Data Lake**
- **Data Warehouse**
- **Data Mart**
  
- **Jobs Data Warehouse**
  - Times
  - Node list
  - Job Energy
  - GPU energy
  - Maxrss
  - Modules loaded
  - Pm counters
  - ...

https://zenodo.org/records/14007065

**zenodo** Search records... 

**Planned intervention:** On Tuesday March 18th 06:30 UTC Zenodo will be unav

Published October 29, 2024 | Version 1.0

## Adastra jobs MI250 15days

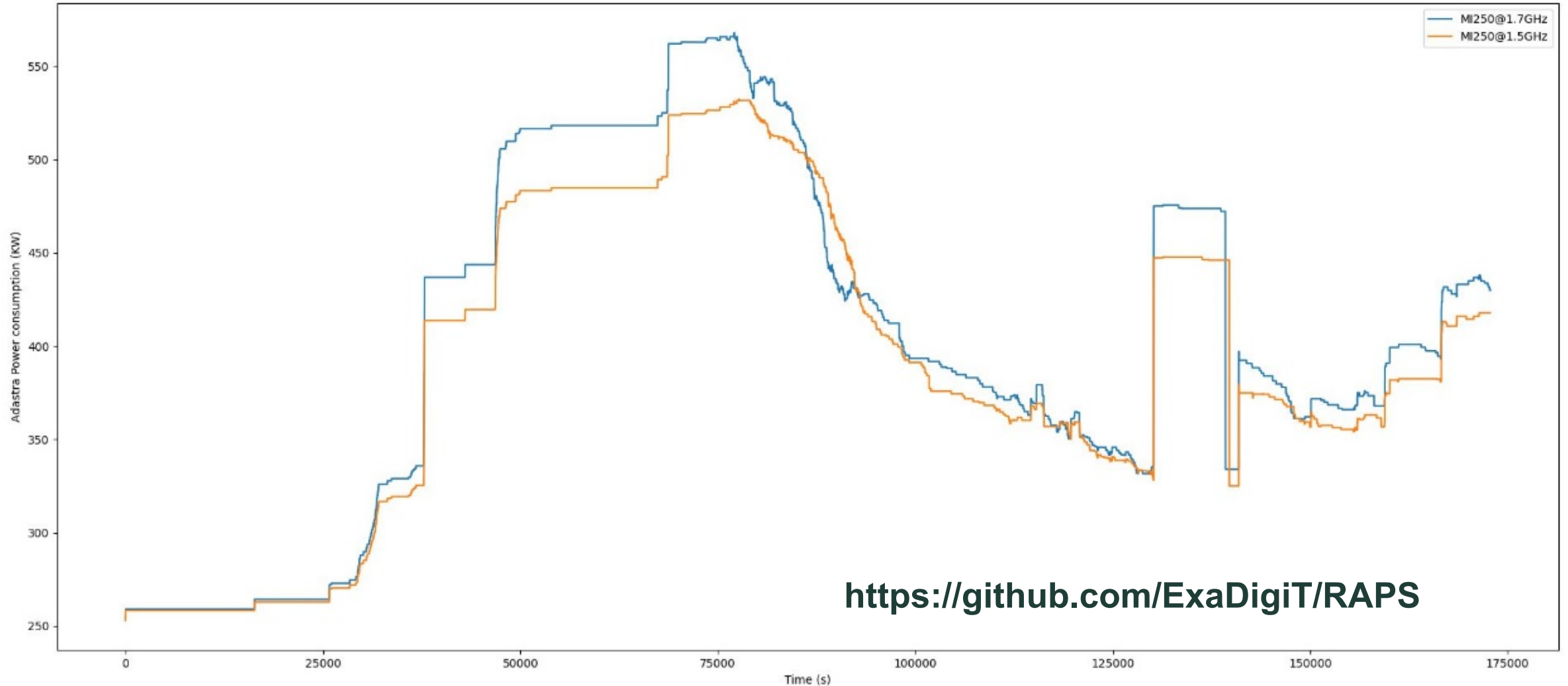
Cirou (Producer)<sup>1</sup> 

Adastra jobs MI250 15days

### Files

Name
<a href="#">AdastaJobsMI250_15days.parquet</a> md5:0e875eaacc1c9f31426715744b3da754e 

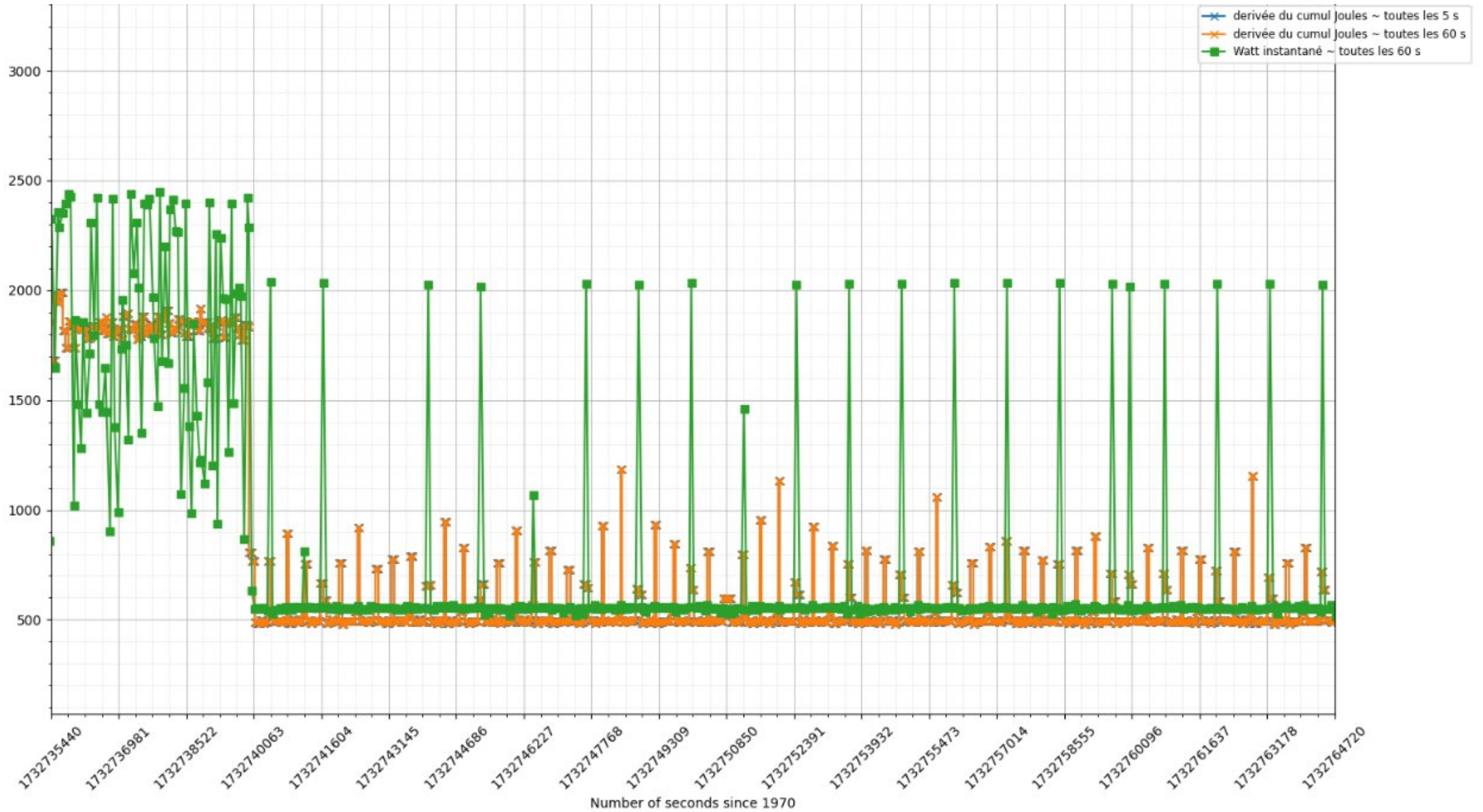
## PM\_counter based (RAPS simulation) Theoretical pure MI250X racks



<https://github.com/ExaDigiT/RAPS>



# Data



# Work in progress

## Digital Twin project:



### Participating Organizations



### Industry Partners



**Hewlett Packard  
Enterprise**



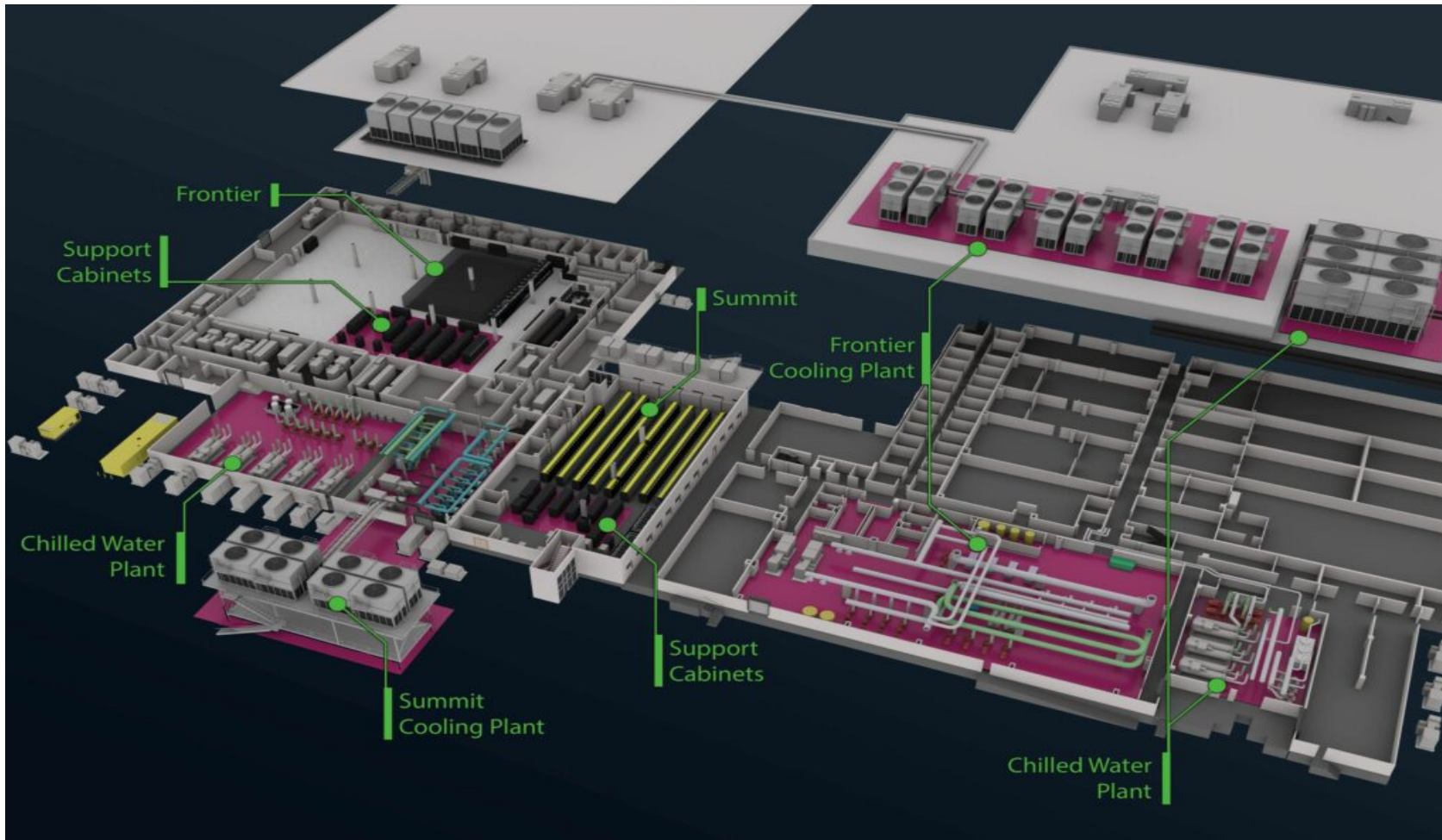
**NVIDIA**

## RAPS: Resource Allocator and Power Simulator

RAPS module can replay workloads from telemetry, reschedule them, or simulate synthetic workloads on the supercomputer to analyze the resulting energy consumption; further details are provided in Section

# Work in progress

## Modelica-based thermo-fluids cooling mode OpenModelica / Dymola®

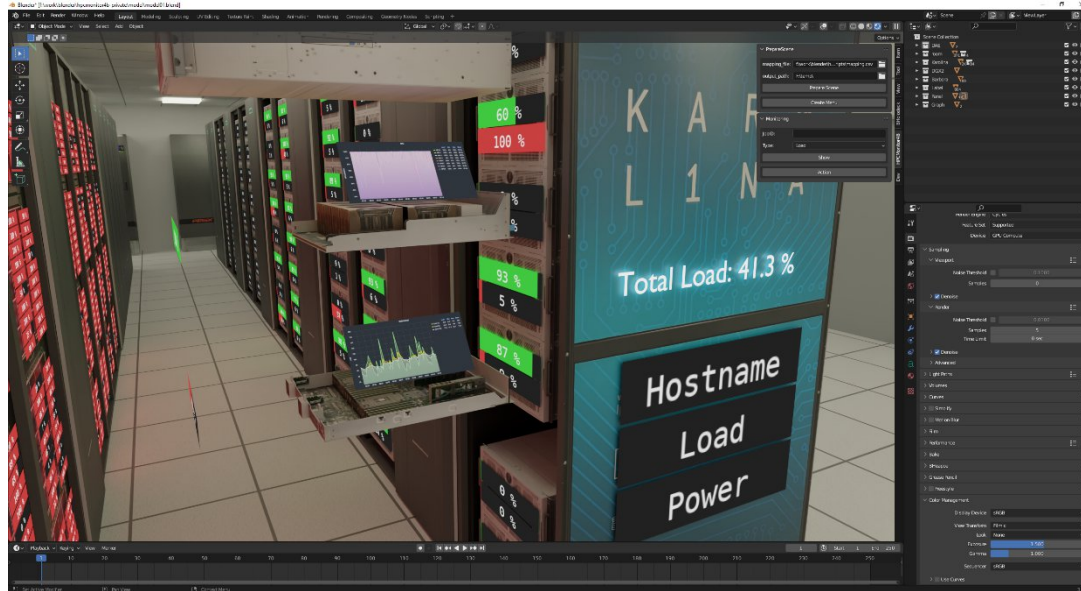


# Work in progress

## Visualization

**HPCMonitor4B**

<https://code.it4i.cz/blender/hpcmonitor4b>



**UnrealEngine**



<https://www.youtube.com/watch?v=SLWVqEonVgw>

