# Hewlett Packard Enterprise

# AIOPS EMPOWERED FAILURE PREDICTION IN SYSTEM MANAGEMENT SOFTWARE TOOLS

Deepak Nanjundaiah
Subrahmanya Joshi
Sergey Serebryakov

May 09, 2024

**CUG2024**
diverse universe
PERTH AUSTRALIA

# AGENDA

- Why AI and ML for HPC systems?

- AIOps Machine Learning and Failure Prediction Framework

- AIOps Dashboards

- Current Status and Future Work

# WHY AI / ML FOR HPC DATA CENTERS?

## Large number of metrics (thousands)

*Real-time data streams at NREL datacenter: $O(10^3)$ facility data points per minute and $O(10^6)$ Eagle data points per minute.*

- Do not know where to look.
- Data is coming real time at high rates
- Threshold-based methods (setpoints) are used, but produce many false positives and thus not scalable
- Some anomalies can only be identified in high-dimensional spaces (multiple metrics).
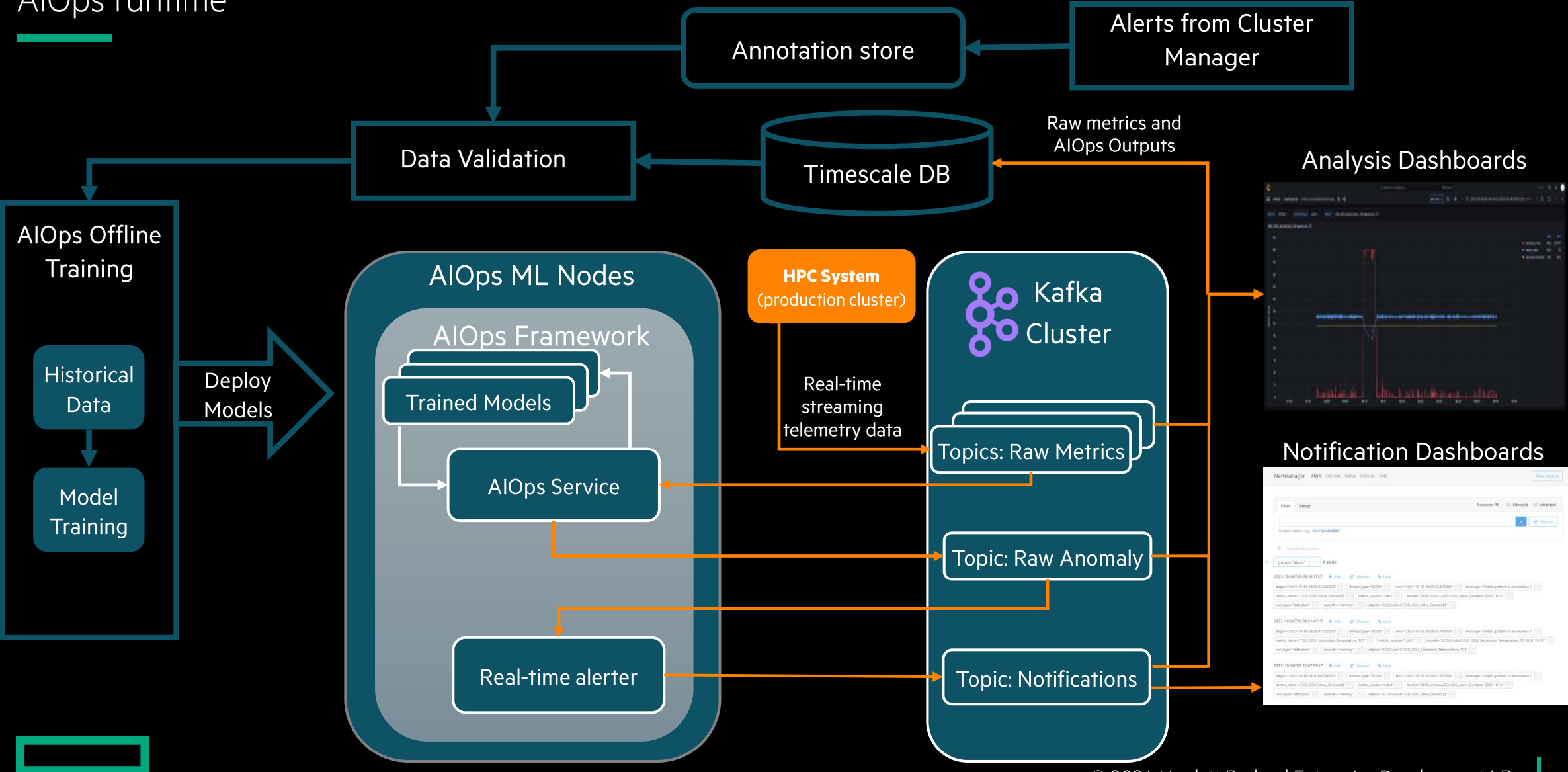- Setpoints and Dashboards are not always sufficient to identify anomalies

## Broad range of problems (not only anomaly detection)

- Anomaly Detection (single metric and multi-metric models)
- Forecasting
- Preventative maintenance
- Reducing carbon footprint by optimizing datacenter ops (Digital Twins)

# AIOPS MACHINE LEARNING FRAMEWORK

AIOps runtime

# AIOPS MACHINE LEARNING FRAMEWORK

From statistical to machine learning models for time series data

**Uni-variate models**

**Multi-variate models**

enable

- Anomaly Detection
- Failure Prediction
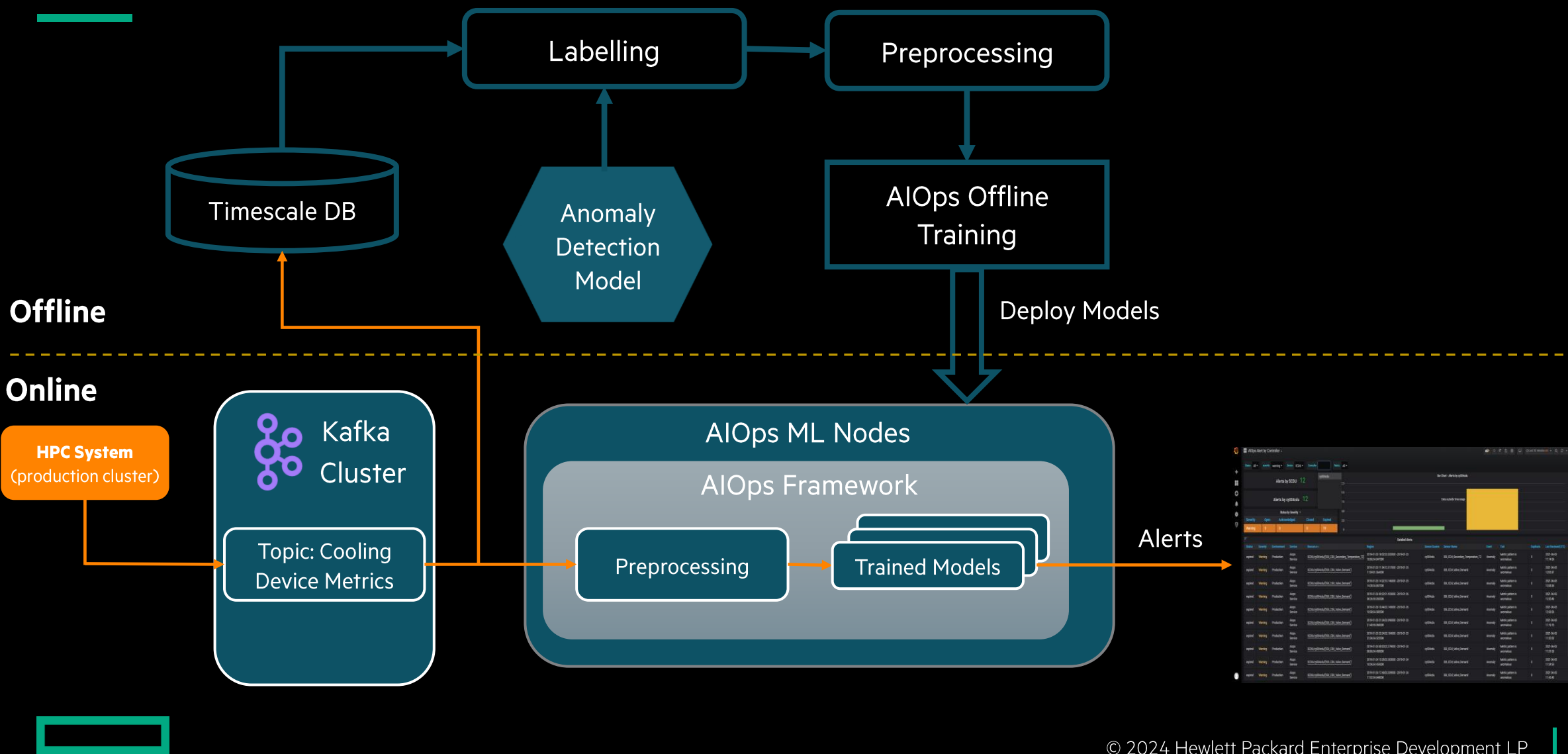- Forecasting

# AIOPS MACHINE LEARNING FRAMEWORK

## AIOps Failure Prediction Framework: Overview

- Some failures are easily predictable, while others are extremely hard to predict.

- 2 types of failures. Some failures occur suddenly, while others are preceded by some specific set of events.

- These occurrences could be anomalies or, at times, simply a normal sequence of events.

- We can train a model to understand these events, and during inference, could potentially predict them even before they happen.

- To accomplish this, we require historical data with known and accurately labeled failures.
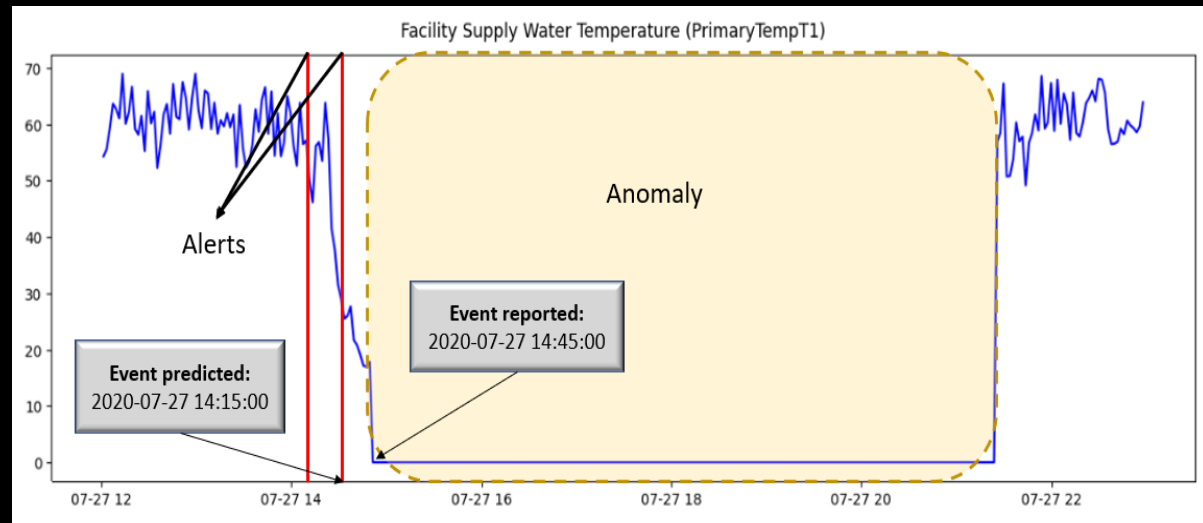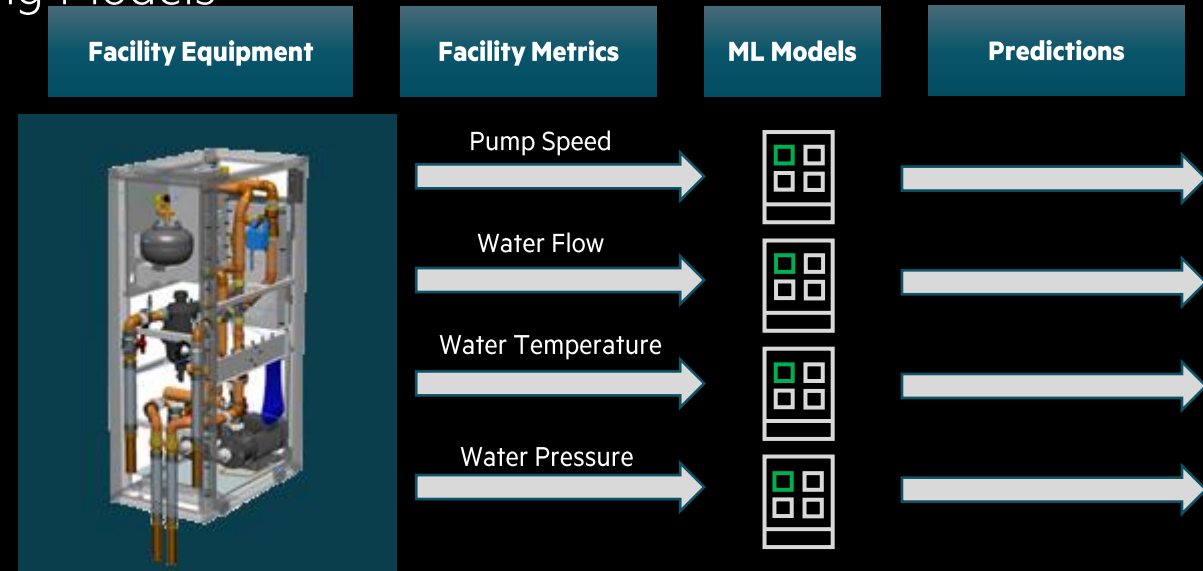
# AIOPS MACHINE LEARNING FRAMEWORK

AIOps Failure Prediction Framework: Online Inference and Offline training architecture

# AIOPS MACHINE LEARNING FRAMEWORK

## AIOps Failure Prediction Framework: Machine Learning Models

- We use LSTM and fully connected neural networks for failure prediction.

- LSTMs very effectively capture and utilize long term dependencies in sequential data.

- Anomaly detection models, which are part of AIOps itself label the data.

- Models are uni-variate, with each being trained and utilized for inference on a single metric.

- Models employ a semi-supervised approach, wherein labeled data is automatically generated from existing detection models.

# AIOPS DASHBOARDS

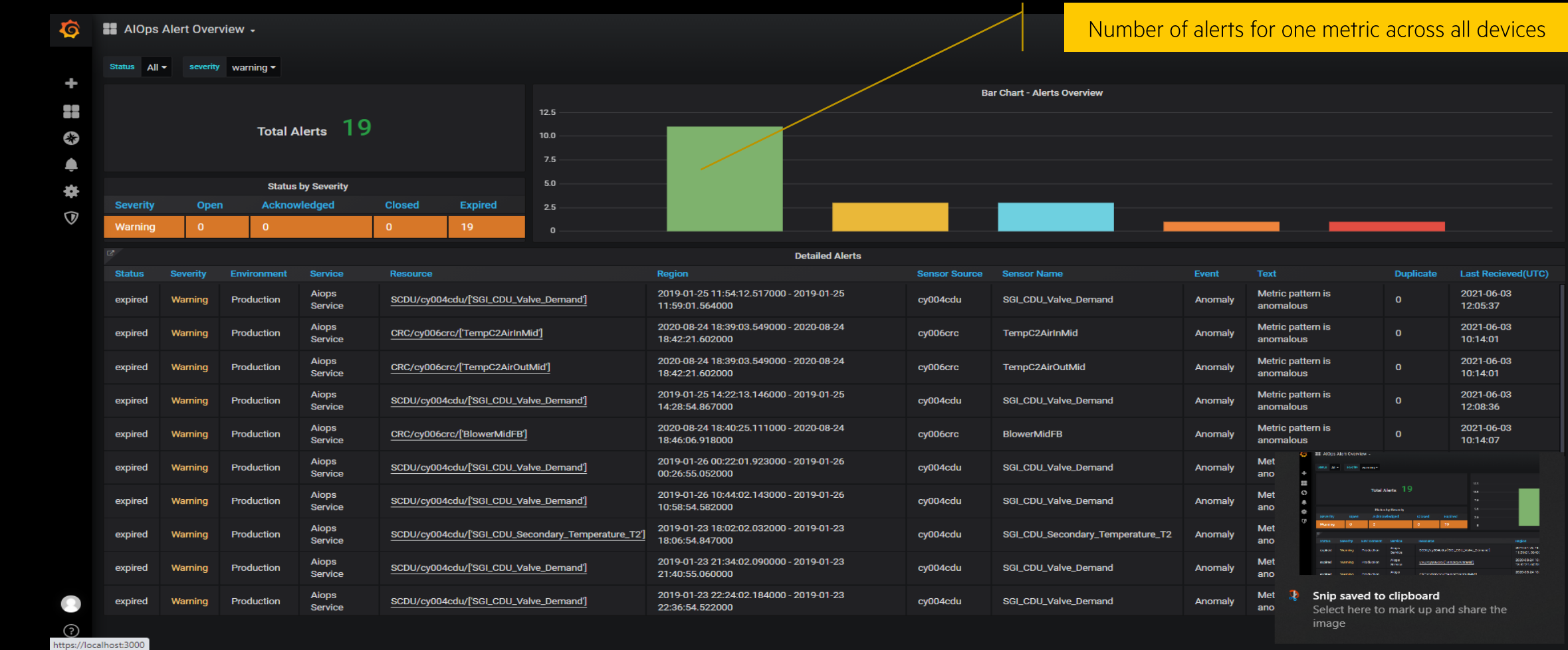Univariate anomaly detection dashboard for Coolant Distribution Units (CDUs)

# AIOPS DASHBOARDS

AIOps alert overview



Number of alerts for one metric across all devices

# AIOPS DASHBOARDS

## Connection with HPCM PCIM visualization



Click metric source to open PCIM interface for this device. This metric will be shown inside a pink rectangle.

| Status | Severity | Environment | Service | Resource ▲ | Region |
|--------|----------|-------------|---------|------------|--------|
| expired | Warning | Production | Aiops Service | SCDU/cy004cdu/['SGI_CDU_Secondary_Temperature_T2'] PCIM Interface | 2019-01-23 02:02.032000 - 2019-01-23 |
| expired | Warning | Production | Aiops Service | SCDU/cy004cdu/['SGI_CDU_Valve_Demand'] | 2019-01-25 11:54:12.517000 - 2019-01-25 11:59:01.564000 |
| expired | Warning | Production | Aiops Service | SCDU/cy004cdu/['SGI_CDU_Valve_Demand'] | 2019-01-25 14:22:13.146000 - 2019-01-25 14:28:54.867000 |

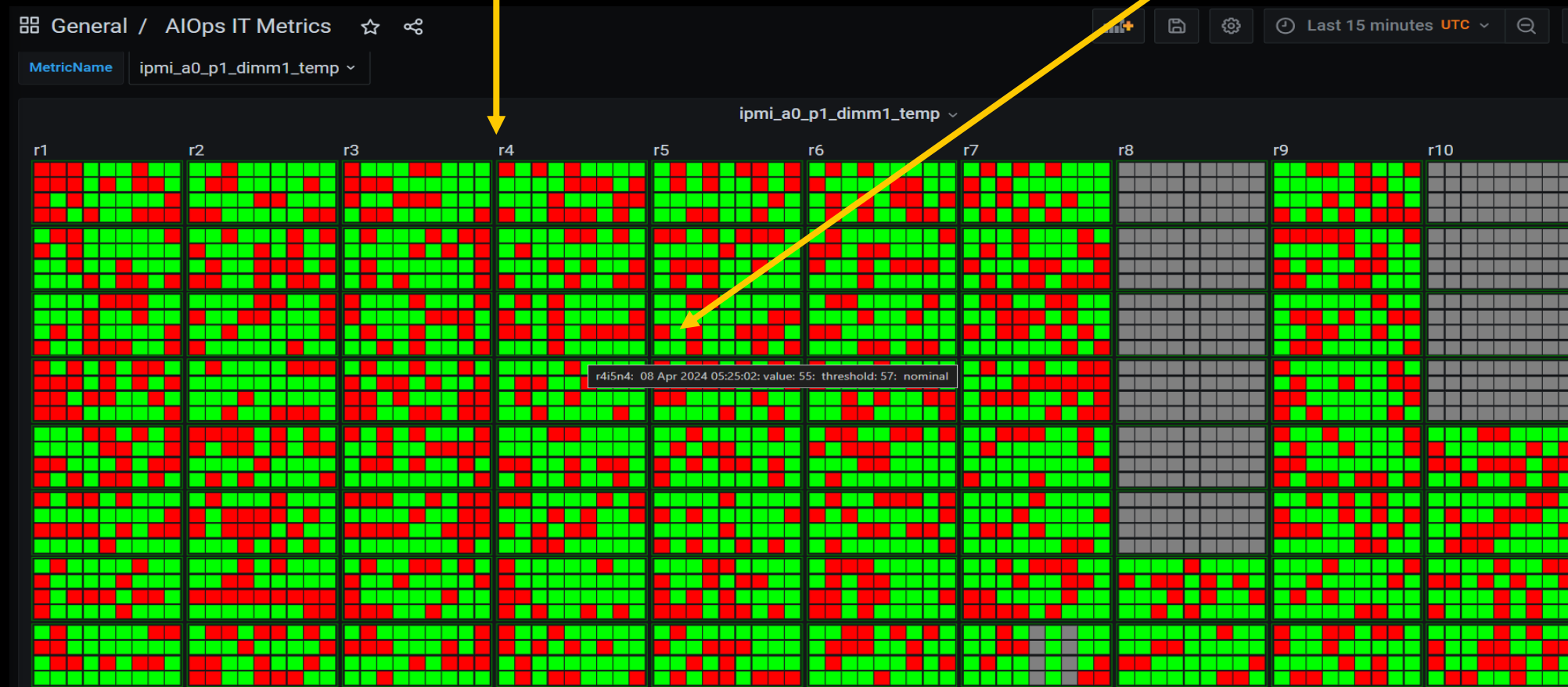PCIM – Power and Cooling Infrastructure Manager

# AIOPS DASHBOARDS
Visualizing IT metrics

compute nodes are displayed in the rack orientation

**Green is nominal , red represent anomalous**

# AIOPS DASHBOARDS
Visualizing Slingshot metrics

Switches are displayed in the group orientation

**Green is nominal , red represent anomalous**



**Click on each switch ,display detail dashboard**

# AIOPS DASHBOARDS

Visualizing Slingshot metrics

Switches are displayed in the group orientation

**Green is nominal , red represent anomalous**



**Click on each switch ,display detail dashboard**

14

# AIOPS DASHBOARDS
Visualizing Slingshot metrics



details of all ports inside the switch.

Green is nominal , red represent anomalous

# STATUS AND FUTURE WORK

**CURRENT STATUS**

- AIOps has been deployed in NREL's production environment since June 2020.
- AIOps has detected 50% more hardware-related anomalies (using historical data).
- AIOps has shown that 40% more high priority incidents that turned into critical events could have been prevented by demonstrated early anomaly detection.
- AIOps supports MLOps, Automation of Model training/creation and reload, model re-training and performance monitoring.
- AIOps supports anomaly detection for
  - Facility metrics CDU (cooling distribution unit), CRC (cooling rack controller), ARCS (Adaptive Rack Cooling System)
  - IT telemetry metrics , with new Grafana visualization panel
  - Slingshot telemetry metrics , with new Grafana visualization panel
- AIOps supports metric forecast and failure prediction for CDU metrics


- AIOps is integrated with HPCM 1.11 and CSM 1.4
- AIOps Integration with GreenLake : A proof of concept (POC) has already been conducted in collaboration with the Green Lake COM team for the purpose of anomaly detection and the reduction of carbon emissions by optimizing the data center ops.


**FUTURE WORK – 2024**

- Multi-variate failure prediction
- Platform independent solution that can be easily leveraged across HPE infrastructure portfolio (GreenLake).
- Analyzing and predicting job failures.
- Log analytics : Anomaly Detection and Diagnosis from System Logs through Deep Learning.

# THANK YOU

**Deepak Nanjundaiah, HPC System Software,** nanjundaiah@hpe.com
**Sergey Serebryakov, HPE Labs,** sergey.serebryakov@hpe.com
**Subrahmanya Joshi,** HPC System Software, subrahmanya-vinayak-joshi@hpe.com
**Amarnath Chilumukuru, HPC System Software,** amarnath.c@hpe.com
**Piyali Hazra, HPC System Software,** piyali.hazra@hpe.com