

# Smart Infrastructure Management System: A Unified Proof-of-Concept Design

## Executive Summary

The exponential growth of AI/ML workloads in modern data centers has created extraordinary demands for sophisticated infrastructure monitoring and management. This proof-of-concept design presents a unified, intelligent management system specifically engineered for large-scale AI training facilities with hundreds of GPU servers running continuous ML workloads.

The proposed solution addresses the critical challenge of managing complex, interdependent infrastructure layers—from data center power and cooling systems to GPU-accelerated compute resources and containerized workloads. By integrating established DCIM protocols with modern cloud-native technologies, this system transforms reactive infrastructure management into a proactive, automated, and intelligent discipline.

Key capabilities include:

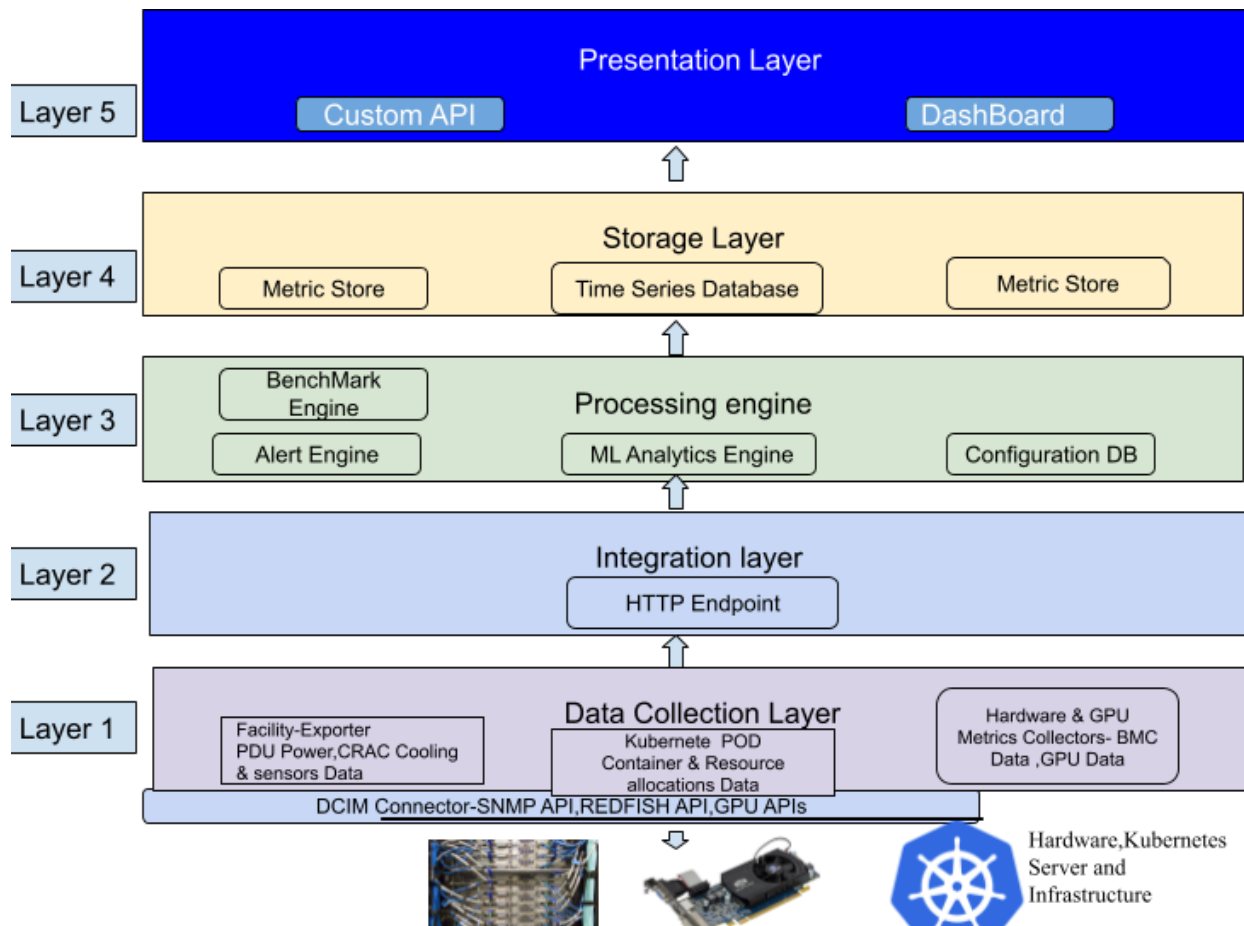
- End-to-end, workload-contextual monitoring—correlating AI/ML job performance with underlying hardware, facility telemetry, and environmental conditions.
- Continuous benchmarking and health validation—automatically detecting performance regressions and quantifying the impact of hardware/software changes.
- Proactive alerting and automated response—including predictive fault detection, root-cause analysis, and workload-aware remediation actions.

## Abbreviations and Acronyms

Abbreviation	Full Form	Definition
AI	Artificial Intelligence	Computer systems that can perform tasks typically requiring human intelligence
API	Application Programming Interface	Set of protocols and tools for building software applications

BMC	Baseboard Management Controller	Specialized microcontroller that manages the interface between system management software and platform hardware
CRAC	Computer Room Air Conditioning	Specialized air conditioning units designed for data center cooling
CRAH	Computer Room Air Handler	Air handling units that circulate conditioned air in data centers
DCGM	Data Center GPU Manager	NVIDIA's suite for managing and monitoring GPUs in data center environments
DCIM	Data Center Infrastructure Management	Software that monitors, manages, and controls data center operations
GPU	Graphics Processing Unit	Specialized electronic circuit designed to accelerate graphics rendering and parallel processing
IPMI	Intelligent Platform Management Interface	Standard for out-of-band management of computer systems
JSON	JavaScript Object Notation	Lightweight data-interchange format
ML	Machine Learning	Type of artificial intelligence that enables computers to learn without explicit programming
NVML	NVIDIA Management Library	C-based API for monitoring and managing NVIDIA GPU devices
PDU	Power Distribution Unit	A device that distributes electrical power to multiple loads
PoC	Proof of Concept	Demonstration that validates the feasibility of a concept or theory
SNMP	Simple Network Management Protocol	Internet standard protocol for collecting and organizing information about managed devices

# System Architecture Diagram



## Five-Layer Architecture Overview

**Layer 1: Data Collection Layer.** This foundational layer interfaces directly with physical infrastructure components, including facility systems, server hardware, and container platforms. Each collector uses protocol-specific native APIs (SNMP v3, Redfish, DCGM GPU APIs & Kubernetes API) to gather raw telemetry data and standardize it into a uniform metrics format for upstream consumption. GPU metrics are collected from Datacenter GPU Manager(DCGM), BMC with Redfish REST APIs, and DCIM from SNMP APIs

**Layer 2: The Integration Layer** serves as the central data aggregation hub, scraping standardized metrics from all collection endpoints using an HTTP

pull-based model. This layer performs protocol translation, data normalization, and correlation across disparate infrastructure domains to create a unified data stream.

**Layer 3: Processing Engine** implements intelligent analysis and decision-making capabilities through machine learning algorithms, alert management, automated response systems, and embedded performance benchmarking modules. In addition to transforming raw infrastructure data into actionable insights, predictive analytics, and automated remediation workflows, this layer continuously validates the performance of infrastructure and AI/ML workloads. Integrated benchmarking tools periodically execute standardized workloads, measuring throughput, latency, and resource efficiency, then correlate these results with real-time telemetry and workload assignments. This enables the system to detect performance regressions, quantify the impact of hardware or firmware changes, and support proactive optimization across the environment.

**Layer 4: Storage Layer** provides specialized data persistence through three complementary database systems optimized for different data types and access patterns. The architecture separates real-time metrics storage, historical time-series analytics, and structured configuration management to ensure optimal performance and scalability.

**Layer 5: Presentation Layer** provides unified interfaces (dashboards, APIs, alerts) for comprehensive infrastructure visibility and control. Users can visualize correlations, integrate with enterprise systems, and manage incidents while abstracting underlying system complexity.

## Kubernetes Correlation

Kubernetes integration for workload correlation is achieved by incorporating collectors and APIs within the Data Collection and Integration layers that directly interface with the Kubernetes control plane. These collectors use the Kubernetes REST API to continuously gather real-time information on running pods, workloads, and node resource allocations—including which GPU servers are running which AI/ML jobs. This data is streamed alongside hardware and facility telemetry into the Integration Layer, where it is correlated with system inventory and operational health metrics from servers, GPUs (via DCGM), and power/cooling systems (via DCIM collectors).

The Processing Engine then links specific workload performance, resource utilization, and potential bottlenecks to underlying infrastructure conditions, enabling root-cause analysis and automated responses—such as relocating jobs away from hardware with emerging faults or thermal issues. Ultimately, correlated insights are presented in

unified dashboards within the Presentation Layer, giving operators end-to-end visibility into how infrastructure health and configuration affect the performance and reliability of AI/ML workloads orchestrated by Kubernetes. This tight integration ensures that infrastructure monitoring is context-aware, workload-centric, and capable of proactive, automated management in a large-scale GPU-based environment system streams data on GPU server AI/ML job allocation, hardware, and facility telemetry into the Integration Layer. This data is then correlated with system inventory and operational health metrics from servers, GPUs (via DCGM), and power/cooling systems (via DCIM collectors). The Processing Engine links workload performance, resource utilization, and bottlenecks to infrastructure conditions. This enables root-cause analysis and automated responses, such as relocating jobs from faulty hardware or thermal issues. Ultimately, correlated insights are presented in unified dashboards within the Presentation Layer, offering operators end-to-end visibility into how infrastructure health and configuration impact the performance and reliability of AI/ML workloads orchestrated by Kubernetes. This tight integration ensures context-aware, workload-centric, and proactive automated management in large-scale GPU environments.

## Automated Benchmarking

In the system architecture, **Performance benchmarking tools** are primarily integrated in the Processing Engine, with their insights visualized in the Presentation Layer dashboards. Specialized benchmarking modules—deployed as additional microservices or agents—periodically execute standardized AI/ML workloads (using frameworks like TensorFlow and PyTorch) and synthetic load tests on GPU servers across the cluster. These benchmarks produce quantitative metrics on throughput, latency, resource utilization, and training efficiency, which are then ingested alongside real-time hardware telemetry via the Data Collection Layer.

The Integration Layer normalizes these benchmarking results and correlates them with relevant infrastructure and workload metadata (such as firmware versions, thermal states, or ongoing Kubernetes jobs). The Processing Engine uses this unified data to provide baseline performance validation, detect regressions following hardware or firmware changes, and enable historical trend analysis for capacity planning.

Automated benchmarking can also be triggered before and after system maintenance, firmware updates, or resource reallocations to quantify the direct impact of these changes on workload performance.

Operators access benchmarking results and historical performance reports through intuitive dashboards in the Presentation Layer, allowing for rapid identification of bottlenecks and data-driven optimization. By embedding these benchmarking tools within the end-to-end monitoring and management workflow, the system ensures that

infrastructure health, hardware/software changes, and AI/ML workload outcomes are tightly and continuously linked—enabling proactive and optimal management at scale.

## Data Flow -DCIM

PDU/CRAC → SNMP API → Facility Exporter (Layer 1: Data Collection) → Integration Layer (Layer 2) → Storage Layer (Layer 3) → Processing Engine (Layer 4) → Presentation Layer (Layer 5)

In Layer 1, Facility Exporters utilize SNMP Library APIs to query and collect this telemetry. They then standardize the collected data into Prometheus metrics format, making it accessible through HTTP endpoints.

Layer 2, the Integration Layer, scrapes this standardized data. It also gathers hardware and workload metrics from other collectors. This layer performs crucial data correlation and normalization before forwarding the refined information to Layer 3, the Storage Layer.

The Storage Layer persists these metrics in specialized databases optimized for time-series data, historical analytics, and configuration management. This robust storage facilitates efficient retrieval for further processing.

Layer 4, the Processing Engine, queries the stored data. It employs sophisticated analytics algorithms to detect patterns, generate intelligent alerts, and trigger automated responses. These responses can include workload migration or the creation of maintenance tickets.

Finally, Layer 5, the Presentation Layer, provides operators with critical tools for oversight. This includes interactive dashboards, programmatic APIs, and alert management interfaces. These tools display the correlated infrastructure health status and enable manual control over automated system

This seamless, upward data flow rapidly converts raw facility telemetry into actionable intelligence, enabling proactive infrastructure management within seconds. This swift transformation helps prevent costly AI/ML training job failures by promptly detecting issues such as cooling system degradation or power distribution anomalies.

## Data Flow -Hardware

### Server Hardware:

Server BMC → Redfish REST API → Hardware Collector (Layer 1: Data Collection)  
→  
Integration Layer (Layer 2) → Storage Layer (Layer 3) → Processing Engine (Layer 4) → Presentation Layer (Layer 5)

### GPU Hardware:

GPU Hardware → DCGM Library API → Hardware Collector (Layer 1: Data Collection) →  
Integration Layer (Layer 2) → Storage Layer (Layer 3) → Processing Engine (Layer 4) → Presentation Layer (Layer 5)

### Combined Hardware Flow:

Server BMC + GPU Hardware → Redfish API + DCGM API → Hardware Collector (Layer 1: Data Collection) → Integration Layer (Layer 2) → Storage Layer (Layer 3) → Processing Engine (Layer 4) → Presentation Layer (Layer 5)

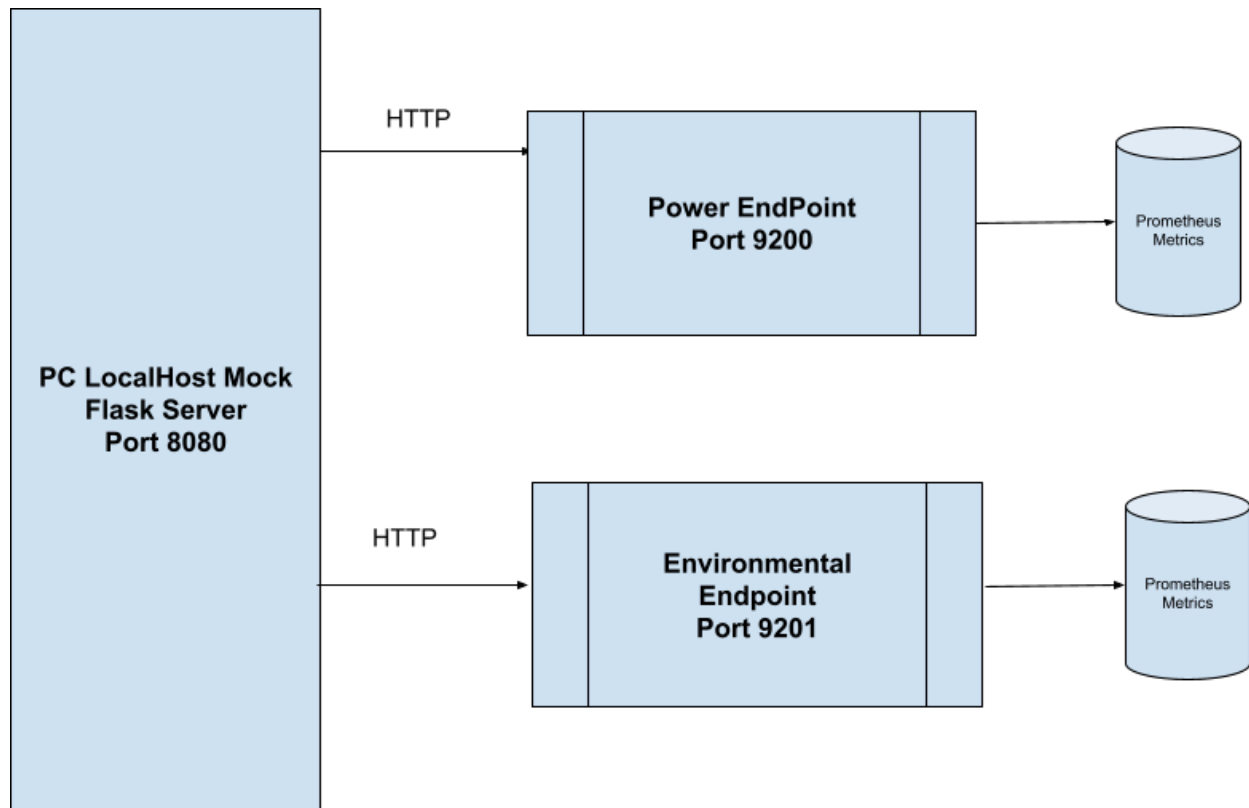
## DCIM Integration EndPoints

DCIM systems manage power, space and environmental data in data centers. Endpoints are the specific URLs/Interfaces where data requested. DCIM Integration endpoints, monitor systems in presentation layer and pulls data from DCIM management systems endpoint. SNMP agent embedded in all DCIM management systems.

Monitoring System ↔ [ENDPOINT] ↔ DCIM System

Implementation:

Implemented a Datacenter mock server in my PC and developed endpoint applications to collect info from mock server



The `mock_dcim_server.py` script creates a Flask-based mock DCIM server, exposing three endpoints: one for rack power consumption, one for environmental temperatures, and another for cooling unit status. Each endpoint, when accessed, generates and serves random data that mimics what you'd get from real hardware sensors, such as fluctuating power readings or status changes in cooling units. This allows development and testing of monitoring and automation integrations in a local-only, fully virtualized scenario.

The `power_endpoint.py` script acts as a simulated PowerIQ Prometheus exporter. It repeatedly queries the Flask server's `/api/v2/racks/power` endpoint, extracts power consumption values for each rack, and exposes these as Prometheus metrics on port 9200. This would let your monitoring stack graph rack power, alert on anomalies, or track long-term trends, all using fully synthetic data.

The `environmental_endpoint.py` script performs a similar function but focuses on environmental data. It collects temperature readings and cooling unit statuses via the Flask server's temperature and cooling endpoints, then exports these as Prometheus metrics on port 9201. This structure accurately mirrors production architectures, where exporters are coupled to real DCIM APIs, but here runs entirely on localhost using



generated data—streamlining development and reducing risk before moving code to production.

## Python POC Code:

(Note: Data collected from my dev machine (GPUModel) intelguy@AIRsearch)

(GPUModel) intelguy@AIRsearch:~/infra\$ cat environmental\_endpoint.py

```
import requests
import time
from prometheus_client import Gauge, start_http_server

class EnvironmentalEndpoint:
    def __init__(self):
        # Prometheus metrics
        self.temperature_gauge = Gauge('datacenter_temperature_celsius',
                                       'Temperature readings', ['location', 'sensor_type'])
        self.cooling_status_gauge = Gauge('cooling_unit_status',
                                          'Cooling unit status', ['unit_id'])

    def collect_environmental_data(self):
        """Collect environmental data from mock DCIM system"""
        try:
            # Get temperature data
            temp_response = requests.get('http://localhost:8080/nlyte/api/v1/sensors/temperature', timeout=5)
            cooling_response = requests.get('http://localhost:8080/nlyte/api/v1/cooling/units', timeout=5)

            if temp_response.status_code == 200:
                temp_data = temp_response.json()
                print(f"Collected temperature data: {len(temp_data['sensors'])} sensors")

                for sensor in temp_data['sensors']:
                    location = sensor['location']
                    sensor_type = sensor['type']
                    temperature = sensor['temperature_celsius']

                    self.temperature_gauge.labels(
                        location=location,
                        sensor_type=sensor_type
                    ).set(temperature)

                print(f"{{location}} ({{sensor_type}}): {{temperature:.1f}}°C")

            if cooling_response.status_code == 200:
                cooling_data = cooling_response.json()
                print(f"Collected cooling data: {len(cooling_data['cooling_units'])} units")

                for unit in cooling_data['cooling_units']:
                    unit_id = unit['id']
```

```

        status = 1 if unit['status'] == 'operational' else 0

        self.cooling_status_gauge.labels(unit_id=unit_id).set(status)
        status_text = "OK" if status else "FAULT"
        print(f" {unit_id}: {status_text}")

    return True

except Exception as e:
    print(f"Connection error: {e}")
    return False

def start_monitoring(self):
    """Start environmental monitoring"""
    start_http_server(9201)
    print(" Environmental Endpoint running on http://localhost:9201/metrics")

    while True:
        self.collect_environmental_data()
        print("Waiting 60 seconds...")
        time.sleep(60)

if __name__ == "__main__":
    endpoint = EnvironmentalEndpoint()
    endpoint.start_monitoring()

(GPUModel) intelguy@AIRsearch:~/infra$ cat power_endpoint.py
import requests
import json
import time
import threading
from prometheus_client import Gauge, start_http_server

class InfraPowerEndpoint:
    def __init__(self):
        # Prometheus metrics
        self.rack_power_gauge = Gauge('rack_power_consumption_watts',
                                      'Power consumption per rack', ['rack_id'])
        self.pdu_status_gauge = Gauge('pdu_status',
                                      'PDU operational status', ['rack_id'])

    def collect_power_data(self):
        """Collect power data from mock DCIM system"""
        try:
            # Call mock DCIM server
            response = requests.get('http://localhost:8080/api/v2/racks/power', timeout=5)

            if response.status_code == 200:
                data = response.json()
                print(f" Collected power data: {len(data['racks'])} racks")

```

```

        for rack in data['racks']:
            rack_id = rack['id']
            power = rack['power_consumption_watts']

            # Update Prometheus metrics
            self.rack_power_gauge.labels(rack_id=rack_id).set(power)
            self.pdu_status_gauge.labels(rack_id=rack_id).set(1) # Assume operational

            print(f" {rack_id}: {power}W")
            return True
        else:
            print(f" Error: HTTP {response.status_code}")
            return False

    except Exception as e:
        print(f" Connection error: {e}")
        return False

    def start_monitoring(self):
        """Start the power monitoring endpoint"""
        # Start Prometheus metrics server
        start_http_server(9200)
        print(" Power Endpoint running on http://localhost:9200/metrics")

        while True:
            self.collect_power_data()
            print(" Waiting 30 seconds...")
            time.sleep(30)

if __name__ == "__main__":
    endpoint = InfraPowerEndpoint()
    endpoint.start_monitoring()
(GPUModel) intelguy@AIResearch:~/infra$

```

## Test Results

```

(GPUModel) intelguy@AIResearch:~/infra$ python mock_dcim_server.py
Mock DCIM Server starting...
Power data: http://localhost:8080/api/v2/racks/power
Temperature data: http://localhost:8080/nlyte/api/v1/sensors/temperature
* Serving Flask app 'mock_dcim_server'
* Debug mode: on
WARNING: This is a development server. Do not use it in a production deployment. Use a production WSGI server instead.
* Running on http://127.0.0.1:8080
Press CTRL+C to quit
* Restarting with stat
Mock DCIM Server starting...
Power data: http://localhost:8080/api/v2/racks/power
Temperature data: http://localhost:8080/nlyte/api/v1/sensors/temperature
* Debugger is active!

```

\* Debugger PIN: 610-110-647

```
(GPUModel) intelguy@AIRsearch:~/infra$ python power_endpoint.py
Power Endpoint running on http://localhost:9200/metrics
Collected power data: 3 racks
rack-01: 7803W
rack-02: 7805W
rack-03: 8044W
Waiting 30 seconds...
```

```
(GPUModel) intelguy@AIRsearch:~/infra$ python environmental_endpoint.py
Environmental Endpoint running on http://localhost:9201/metrics
Collected temperature data: 3 sensors
rack-01 (ambient): 25.2°C
rack-02 (supply): 21.8°C
server-room (return): 26.8°C
    Collected cooling data: 3 units
crac-01: OK
crac-02: OK
crac-03: OK Waiting 60 seconds...
```

## Conclusion

This PoC has successfully delivered:

- ❖ DCIM Integration Endpoints: Two fully tested mock collectors ( rack power and environmental metrics) that export Prometheus-compatible telemetry
- ❖ Layered Architecture Validation: A comprehensive five-layer blueprint that confirms seamless data flow from physical sensors through processing to visualization
- ❖ Kubernetes Workload Correlation: REST-API collectors that ingest pod-level metadata and correlate it with GPU and DCIM metrics for end-to-end visibility
- ❖ Automated Performance Benchmarking: Microservices executing TensorFlow/PyTorch workloads, reporting throughput, latency, and resource utilization in real time, and integrating results into the unified telemetry stream for regression analysis
- ❖ Inventory & Firmware Management Foundation: Redfish/BMC-based collectors for server inventory, health, and firmware lifecycle monitoring embedded in the Data Collection layer

