# BM@N DAQ Data Center

BM@N Experiment at the NICA Facility
14th Collaboration Meeting
JINR, Dubna, May 13 – 15, 2025

ILIA SLEPNEV, JINR

# DAQ Data Center
## Outline

**DAQ Data Center**

- Mission & Design Principles
  - Purpose of the DAQ Data Center
  - Guiding Architecture Principles
- Data Flow & Network Fabric
  - End-to-End Data Path
- BM@N DAQ Network Topology
  - DAQ Network Performance & Reliability
  - High Availability & Storage Resilience
  - Readout Link Redundancy Constraints
- Compute, Storage & Virtualization
  - Distributed Storage Cluster (CephFS Layer)
  - Virtual & Bare-Metal Compute Tiers

**LHEP Computing Resources**

- LHEP data centers in operation

**Monitoring**

- Incident Resolution
- Node-RED Automation
- Grafana: DDC dashboard

**Extra**

- Infrastructure Management
- Monitoring Architecture
- Log Message Processing and Analysis

# DAQ Data Center

## Mission & Design Principles

**BM@N**

**Purpose of the DAQ Data Center**

- Central hub for experiment data reception and archiving
- Decouples micro-second readout from second-scale processing
- Guarantees continuous acquisition during long physics runs
- Operates autonomously, no external IT dependencies
- Hosts online monitoring for quality-of-data checks
- Designed to evolve without disruptive rewiring or downtime

**Guiding Architecture Principles**

- All critical paths redundant, no single point of failure
- Software-defined everything: storage, network, virtualization
- Hardware chosen for low-latency performance
- Horizontal scalability—add nodes, no redesign required
- Observability first: fine-grained metrics, logs, alerts
- Security via dedicated VLANs, JINR SSO and firewall zoning
- Documentation-driven operations; infra declared in Git
- Emphasizes open-source, vendor-neutral technologies

# DAQ Data Center
## Data Flow & Network Fabric



**End-to-End Data Path**

- Detector front-ends push UDP streams via custom MStream protocol
- Hardware IP core in FPGA provides dat transfer and control
- First-Level Processor buffers, validates, formats packets
- Event builders assemble multi-detector fragments asynchronously
- Asynchronous layers prevent back-pressure on readout electronics
- FLP quality checks tag corrupted or partial events
- Event files streamed to high-availability Ceph storage
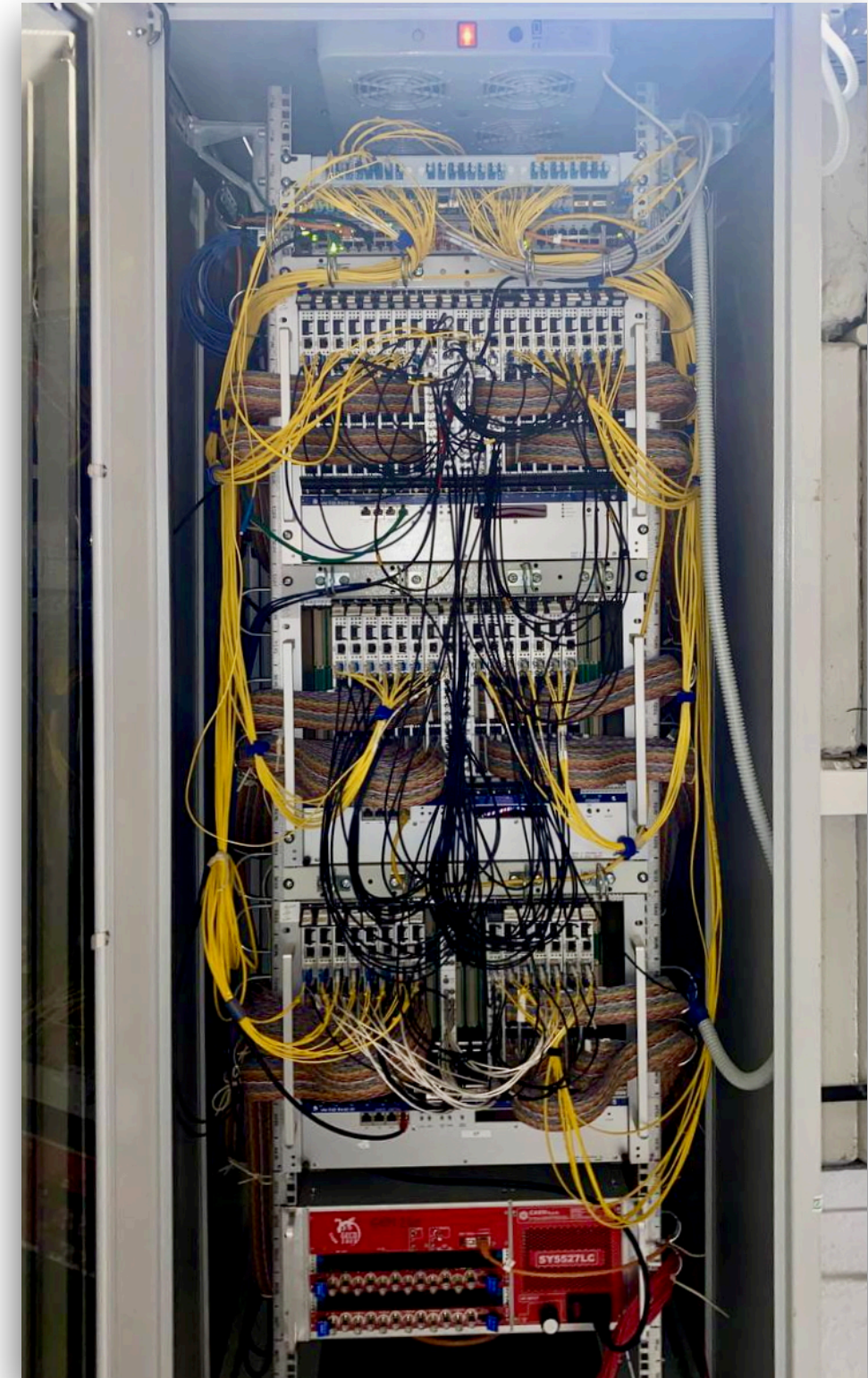- Complete files are synchronized to Offline farm immediately



**TOF Readout Board**

- VXS 6U 160 cm form factor
- Front-Panel: detector I/O
- Backplane: sync, readout 1Gb/s



**21-slot, VXS Chassis**

- Power, Cooling, Remote control
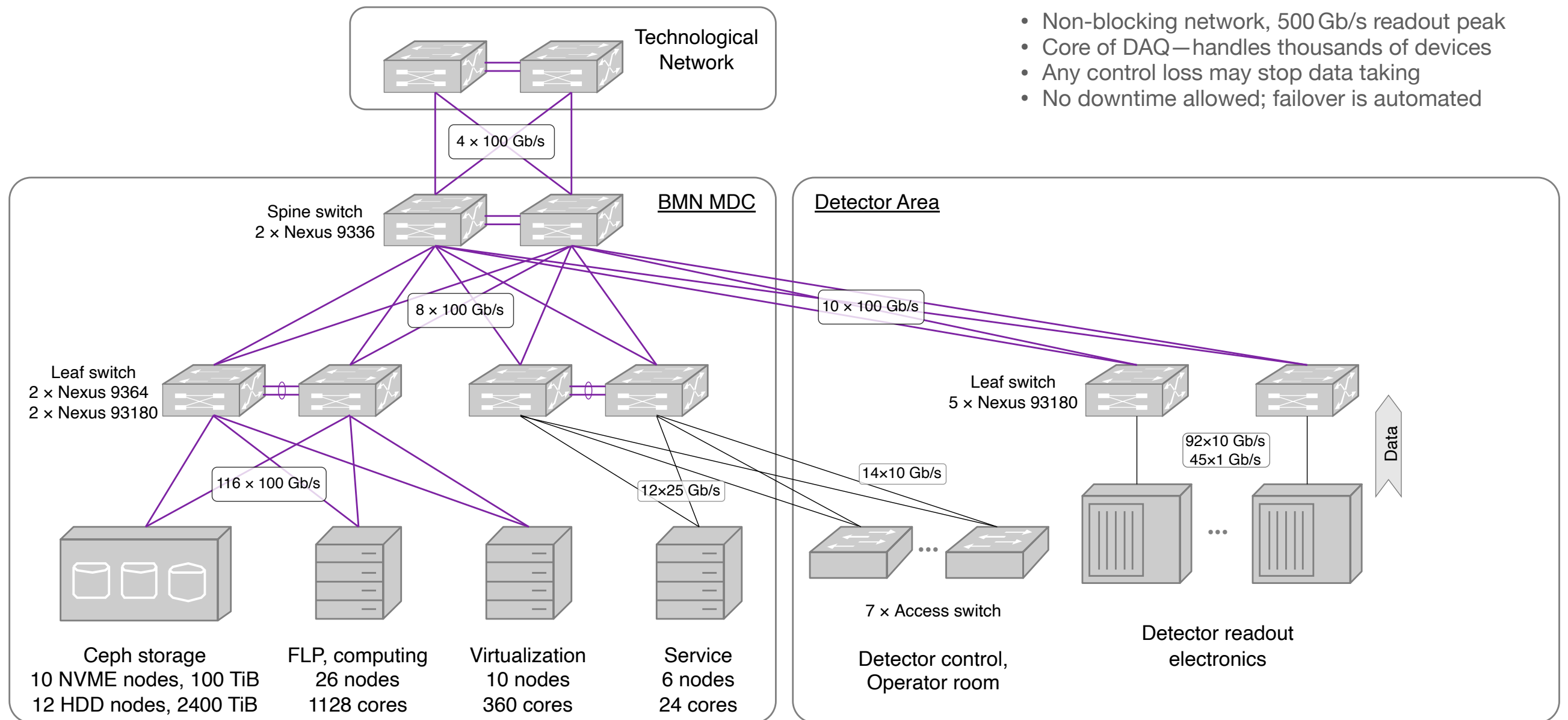- Backplane: Dual-Star, 10 Gb/s per slot

**Readout Electronics Rack: GEM**

- 3 × VME64x chassis with ADC64 boards
- Fiber-optical front panel readout links
- Ribbon cable detector I/O
- Top-of-Rack Cisco Nexus 93180 switch

# DAQ Data Center
## BM@N DAQ Network Topology



**DAQ Network Performance & Reliability**

- Non-blocking network, 500 Gb/s readout peak
- Core of DAQ—handles thousands of devices
- Any control loss may stop data taking
- No downtime allowed; failover is automated

Technological Network

4 × 100 Gb/s

Spine switch
2 × Nexus 9336

BMN MDC

Detector Area

8 × 100 Gb/s

10 × 100 Gb/s

Leaf switch
2 × Nexus 9364
2 × Nexus 93180

Leaf switch
5 × Nexus 93180

92×10 Gb/s
45×1 Gb/s

Data

116 × 100 Gb/s

12×25 Gb/s

14×10 Gb/s

7 × Access switch

Detector readout
electronics

Ceph storage
10 NVME nodes, 100 TiB
12 HDD nodes, 2400 TiB

FLP, computing
26 nodes
1128 cores

Virtualization
10 nodes
360 cores

Service
6 nodes
24 cores

Detector control,
Operator room

**High Availability & Storage Resilience**

- Ceph ensures 24/7 access with auto-recovery on failure
- HDD pools use erasure coding, NVMe SSD pools triple-replication
- Storage remains online with no data loss or downtime
- External users can access data at all times

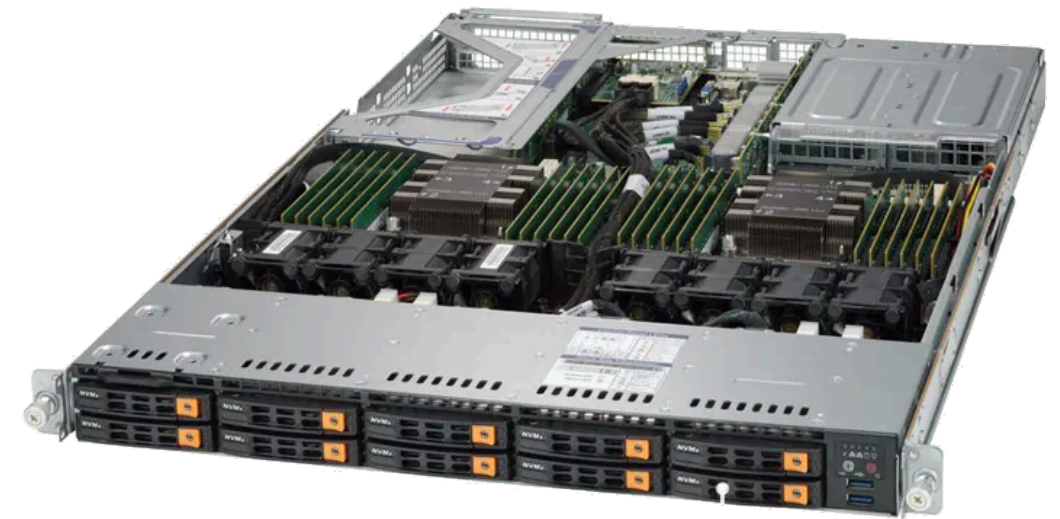**Readout Link Redundancy Constraints**

- No redundancy in readout links due to custom hardware
- Failover is complex, costly, and resource-limited
- Manpower and design constraints limit improvements
- Focus is on robustness and failure monitoring

# DAQ Data Center
## Compute, Storage & Virtualization
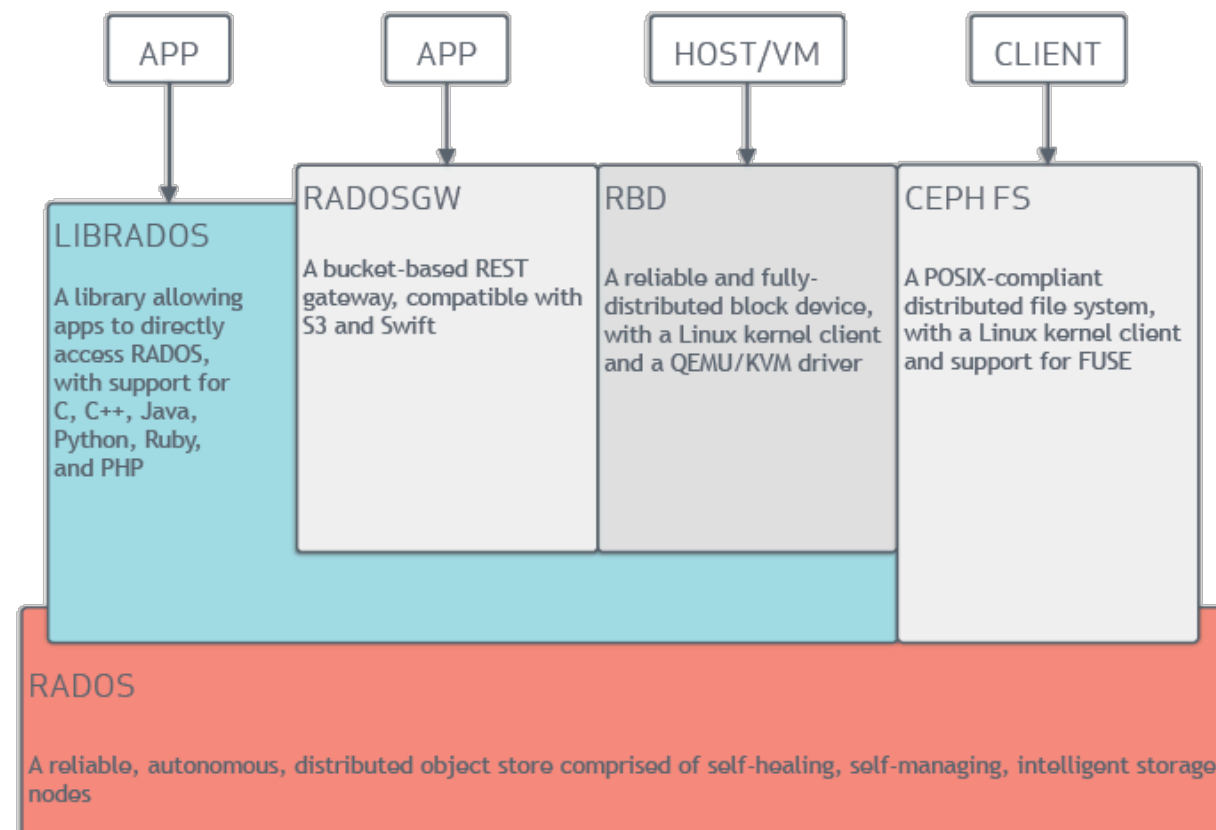
**Distributed Storage Cluster (CephFS Layer)**

- NVMe replicated pools for low-latency RBD workloads
- HDD pools with erasure coding for cost-effective capacity
- CRUSH algorithm enforces fault-domain-aware placement
- Self-healing scrubs verify checksums without operator action
- POSIX CephFS exports RAW data to FLP and event builders
- Grafana dashboards track PG-health, latency, rebuild rate
- Capacity expansion performed online by adding OSD nodes
- Data-at-rest encrypted; keys managed by vault-backed KMS
- Network-isolated background replication to off-site cluster
- Retention policy aligns with collaboration data mandate



Ceph OSD node: 10 × 3.8 TB NVMe SSD

| APP | APP | HOST/VM | CLIENT |
|---|---|---|---|

| LIBRADOS | RADOSGW | RBD | CEPH FS |
|---|---|---|---|
| A library allowing apps to directly access RADOS, with support for C, C++, Java, Python, Ruby, and PHP | A bucket-based REST gateway, compatible with S3 and Swift | A reliable and fully-distributed block device, with a Linux kernel client and a QEMU/KVM driver | A POSIX-compliant distributed file system, with a Linux kernel client and support for FUSE |

**RADOS**

A reliable, autonomous, distributed object store comprised of self-healing, self-managing, intelligent storage nodes



Ceph OSD node: 24 × 18 TB HDD

# DAQ Data Center
## Compute, Storage & Virtualization

**Virtual & Bare-Metal Compute Tiers**

- KVM / LXC cluster hosts control, monitoring, DB, web portals
- Bare-metal nodes reserved for FLP and heavy event builders
- Live-migration keeps services online during maintenance
- Snapshot roll-back guards against faulty software updates
- Resources re-balanced when DAQ idle to boost offline jobs
- Continuous incremental backups: hourly, daily, weekly, yearly
- SSO-protected self-service portal for VM lifecycle requests
- Performance tuned: CPU-governor, NUMA pinning, NIC interrupts



Dual-node compute server



**PROXMOX** Virtual Environment 8.3.5 — Search — Documentation — Create VM — Create CT — islepnev@JINR
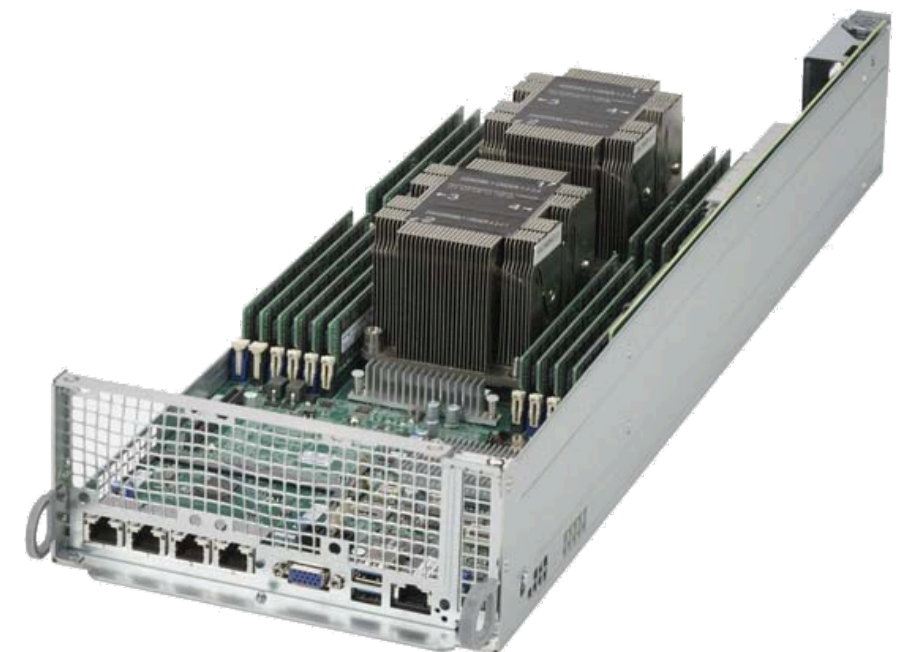
Folder View

- Datacenter (C5)
  - LXC Container
  - Nodes
  - Resource Pool
    - BMN-APP
    - BMN-SC
    - Batch
    - DAQ
    - Local
    - MDC
    - STS
    - TestMPD
    - common
  - Virtual Machine
  - SDN
  - Storage
    - backup-bk1 (c5n01)
    - bmn-daq (c5n01)
    - bmn-sc (c5n01)
    - iso (c5n01)
    - local-zfs (c5n01)
    - mdc (c5n01)
    - pbs1 (c5n01)
    - rbd-batch (c5n01)
    - rbd-sts (c5n01)
    - teleport (c5n01)

Virtual Machine — Help — Search:

| Type ↑ | Description | Disk usage… | Memory us… | CPU us |
|---|---|---|---|---|
| qemu | 501 (bmn-gem-test) | | | |
| qemu | 502 (bmn-csc) | 0.0 % | 35.5 % | 5.3% o |
| qemu | 503 (bmn-fhcal-1) | | | |
| qemu | 504 (bmn-fhcal) | 0.0 % | 33.2 % | 3.5% o |
| qemu | 505 (vmc) | | | |
| qemu | 506 (bmn-ecal) | | | |
| qemu | 508 (bmn-tof700) | 0.0 % | 13.1 % | 10.7% |
| qemu | 509 (bmn-fsd) | 0.0 % | 14.1 % | 12.7% |
| qemu | 512 (bmn-gem) | 0.0 % | 18.3 % | 12.7% |
| qemu | 515 (bmn-daq) | 0.0 % | 7.4 % | 7.7% o |
| qemu | 519 (bmn-msc-1) | 0.0 % | 17.7 % | 6.2% o |
| qemu | 520 (bmn-fsd-win) | | | |
| qemu | 521 (bmn-t0) | 0.0 % | 8.4 % | 0.6% o |
| qemu | 522 (bmn-gem-2) | | | |
| qemu | 524 (bmn-tof400) | 0.0 % | 10.2 % | 11.2% |
| qemu | 525 (bmn-ts) | 0.0 % | 15.2 % | 0.4% o |
| qemu | 526 (bmn-tof400-evb) | 0.0 % | 5.5 % | 0.1% o |
| qemu | 529 (bmn-ceph-fs1) | | | |
| qemu | 532 (bmn-radius) | | | |

**Compute node**
RAM: 384 GB
CPU: Dual Xeon Gold 6154, 6342

# LHEP Computing Resources
## Data Centers in Operation

### DDC — DAQ Data Center

- Data taking, online processing and monitoring
- Primary RAW data storage
- CephFS for immediate data access
- Secondary role: batch jobs outside DAQ periods

### NCX — Offline Cluster

- Batch and interactive jobs at large scale
- Experimental and simulation data storage: EOS
- Shared by multiple experiments, 200+ users

| | DDC (BMN) | DDC (MPD) | NCX |
|---|---|---|---|
| **Location** | Building 215, BMN area | MPD Hall | Building 216, room 115 |
| **Operating since / Last upgrade** | 2019 / 2023 | 2021 / 2024 | 2019 / 2023 |
| **CPU architecture** | 3.0 GHz Skylake<br>2.8 GHz Ice Lake | 3.1 GHz Cascade Lake<br>2.8 GHz Ice Lake | 2.6 GHz Broadwell-EP<br>2.0 GHz Skylake<br>2.5 GHz Cascade Lake |
| **CPU cores — Total** | 1488 | 1760 | 4200 |
| **CPU cores — Batch** | 700 [1] | 1000 [1] | 3000 [2] |
| **RAM, GB per CPU core** | 6–7.5 | 7.5 | 9.6–16 |
| **Node uplink, Gb/s** | 2×100 | 2×100 | 100 |
| **Local node storage (/tmp)** | 32 GB SSD | 32 GB SSD | 1 TB HDD |
| **Shared storage (workspace, experimental and simulation data)** | 100 TB NVMe (CephFS)<br>2.5 PB HDD (CephFS) | | 124 TB NVMe (NFS+ZFS)<br>11 PB HDD (EOS) |

[1] Additional CPU cores available when DAQ is inactive
[2] May be temporarily unavailable due to maintenance

### CephFS

- Distributed POSIX-compliant file system on Ceph
- Supports high-throughput, parallel data access
- Used for RAW data from DAQ in DDC clusters
- Mounted natively on compute and DAQ nodes

### EOS

- CERN-developed distributed file system
- Optimized for large-scale data access via XRootD
- Used in NCX cluster for simulation and analysis
- Suited for shared access by many experiments
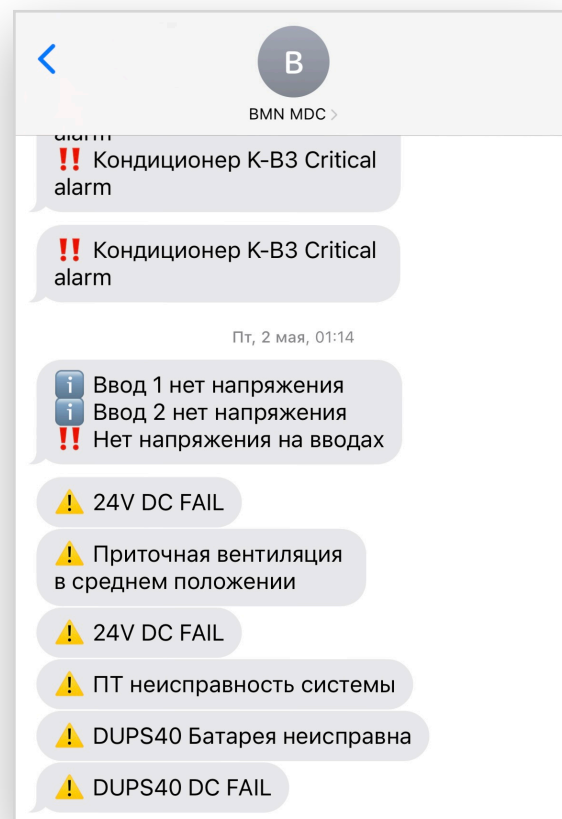
# Computing Resources
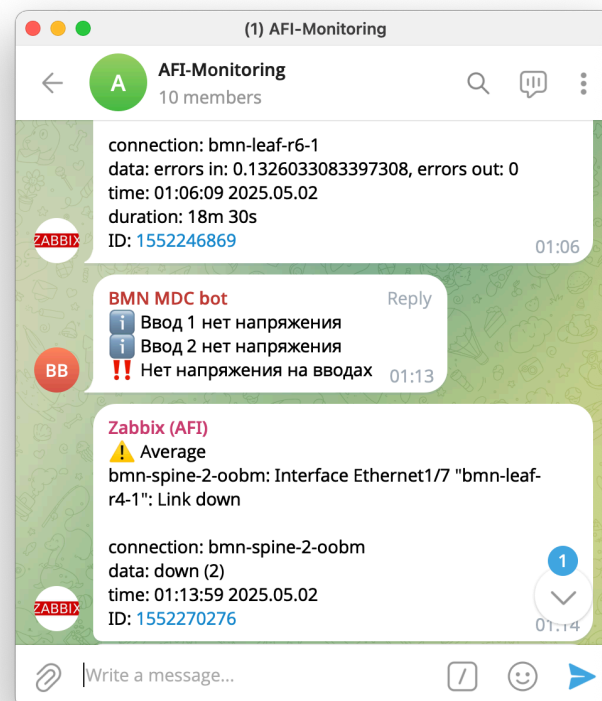## Data Centers in Operation



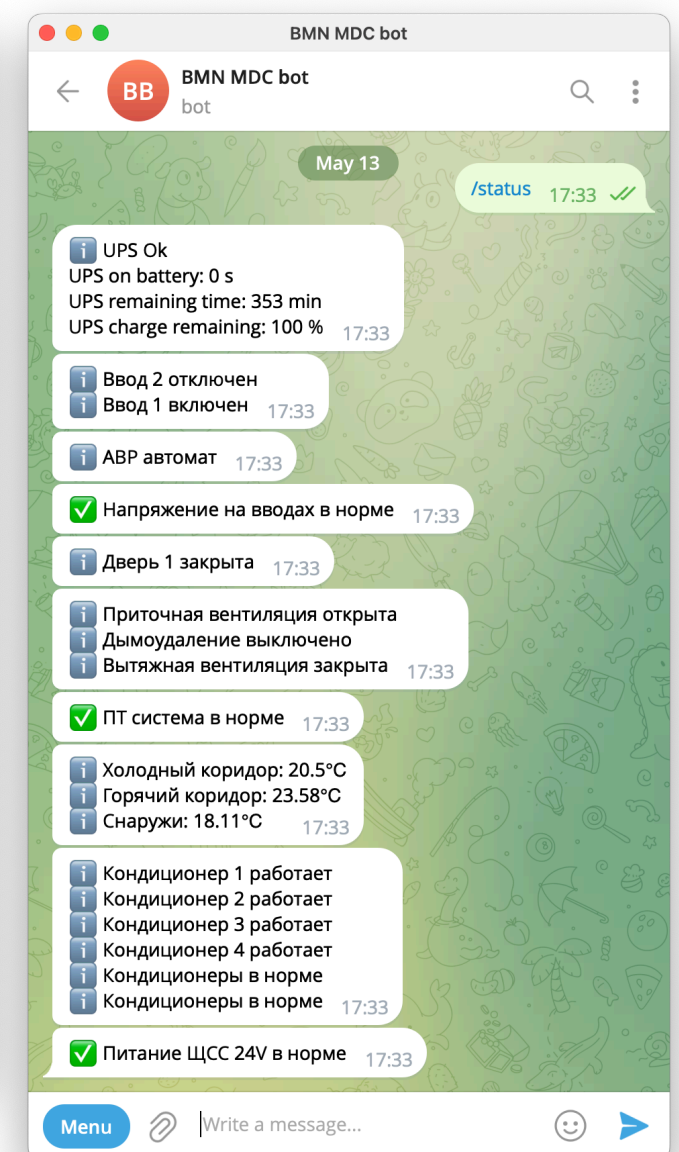BM@N DAQ Data Center
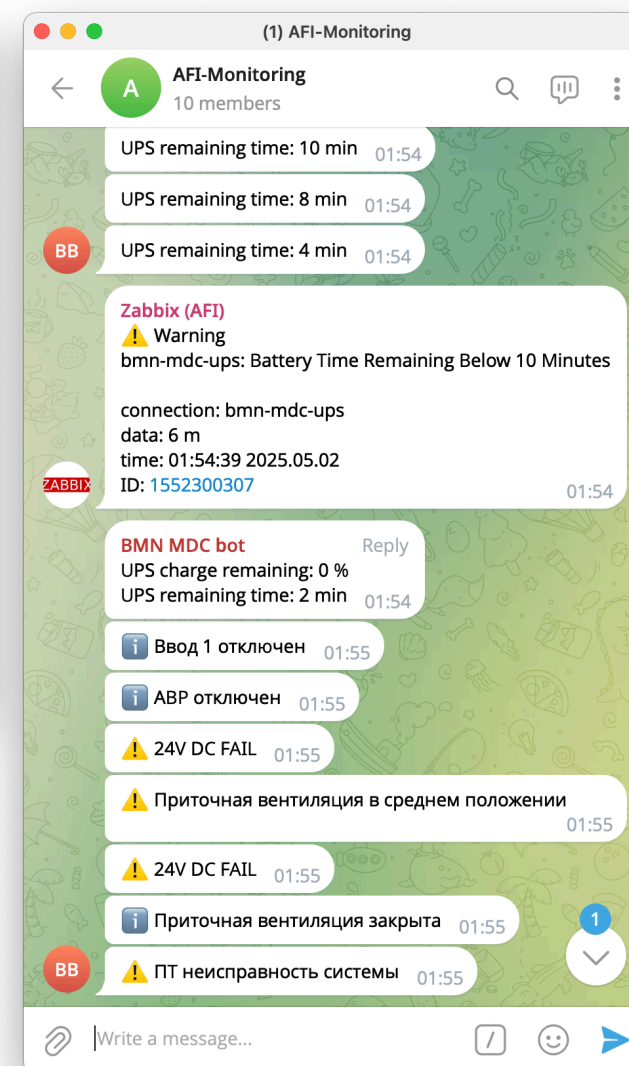


NCX Offline Cluster

# Monitoring
## Incident Resolution

**Detect**
- Node-RED
- Zabbix

→

**Alert**
- Support engineers
- Management

→

**Coordinate**
- Phone call
- Messengers

→

**Resolve**
- Switch-On
- Verify
- Postmortem

---

**SMS Alert**

BMN MDC

‼️ Кондиционер К-В3 Critical alarm

‼️ Кондиционер К-В3 Critical alarm

Пт, 2 мая, 01:14

ℹ️ Ввод 1 нет напряжения
ℹ️ Ввод 2 нет напряжения
‼️ Нет напряжения на вводах

⚠️ 24V DC FAIL

⚠️ Приточная вентиляция в среднем положении

⚠️ 24V DC FAIL

⚠️ ПТ неисправность системы

⚠️ DUPS40 Батарея неисправна

⚠️ DUPS40 DC FAIL

---

**Telegram Alert**

(1) AFI-Monitoring

AFI-Monitoring
10 members

connection: bmn-leaf-r6-1
data: errors in: 0.1326033083397308, errors out: 0
time: 01:06:09 2025.05.02
duration: 18m 30s
ID: 1552246869          01:06

BMN MDC bot          Reply
ℹ️ Ввод 1 нет напряжения
ℹ️ Ввод 2 нет напряжения
‼️ Нет напряжения на вводах          01:13

Zabbix (AFI)
⚠️ Average
bmn-spine-2-oobm: Interface Ethernet1/7 "bmn-leaf-r4-1": Link down

connection: bmn-spine-2-oobm
data: down (2)
time: 01:13:59 2025.05.02
ID: 1552270276          01:14

---

(1) AFI-Monitoring

AFI-Monitoring
10 members

UPS remaining time: 10 min          01:54
UPS remaining time: 8 min          01:54
UPS remaining time: 4 min          01:54

Zabbix (AFI)
⚠️ Warning
bmn-mdc-ups: Battery Time Remaining Below 10 Minutes

connection: bmn-mdc-ups
data: 6 m
time: 01:54:39 2025.05.02
ID: 1552300307          01:54

BMN MDC bot          Reply
UPS charge remaining: 0 %
UPS remaining time: 2 min          01:54

ℹ️ Ввод 1 отключен          01:55
ℹ️ АВР отключен          01:55
⚠️ 24V DC FAIL          01:55
⚠️ Приточная вентиляция в среднем положении          01:55
⚠️ 24V DC FAIL          01:55
ℹ️ Приточная вентиляция закрыта          01:55
⚠️ ПТ неисправность системы          01:55

---

BMN MDC bot

BMN MDC bot
bot

May 13

/status          17:33

ℹ️ UPS Ok
UPS on battery: 0 s
UPS remaining time: 353 min
UPS charge remaining: 100 %          17:33

ℹ️ Ввод 2 отключен
ℹ️ Ввод 1 включен          17:33

ℹ️ АВР автомат          17:33

✅ Напряжение на вводах в норме          17:33

ℹ️ Дверь 1 закрыта          17:33

ℹ️ Приточная вентиляция открыта
ℹ️ Дымоудаление выключено
ℹ️ Вытяжная вентиляция закрыта          17:33

✅ ПТ система в норме          17:33

ℹ️ Холодный коридор: 20.5°C
ℹ️ Горячий коридор: 23.58°C
ℹ️ Снаружи: 18.11°C          17:33

ℹ️ Кондиционер 1 работает
ℹ️ Кондиционер 2 работает
ℹ️ Кондиционер 3 работает
ℹ️ Кондиционер 4 работает
ℹ️ Кондиционеры в норме
ℹ️ Кондиционеры в норме          17:33

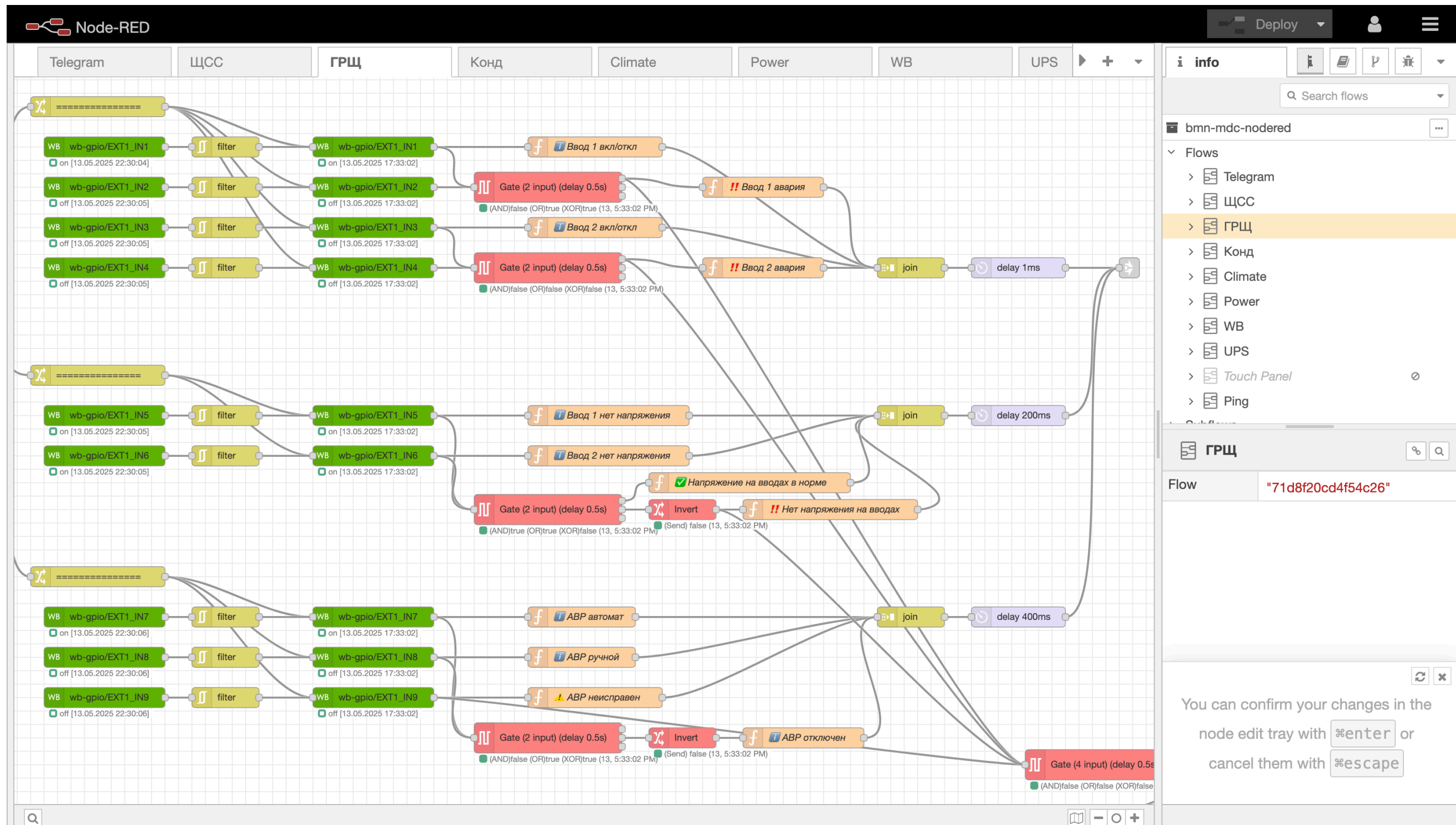✅ Питание ЩСС 24V в норме          17:33

Menu

---

- Real-time alerts via SMS and Telegram bots
- Automated escalation from monitoring systems
- Critical events resolved via coordinated response
- Status recovery confirmed by bot interaction

# Monitoring
## Node-RED Automation

- Node-RED automates low-level alert processing
- Monitors power, HVAC, ventilation, security
- Custom logic triggers SMS and Telegram alerts instantly
- Each flow handles real-time decision-making
- Part of end-to-end automated monitoring system

# Monitoring
## Grafana: DDC dashboard

# Acknowledgements

**Thank you!**

# Extra slides

# Infrastructure Management
## Infrastructure-as-Code

| Tool | Method | Approach, our usage | Tasks |
|------|--------|---------------------|-------|
| Puppet | Pull | functional (declarative) | configure services, settings |
| Ansible | Push | procedural (imperative) | updates, one-time tasks |
| — | — | manual admnistration | other complex tasks |

- Machine-readable, version-controlled configuration files (YAML, Ruby)
- Puppet modules:
  - provision, configure, manage OS and application components
  - supported by community or our custom solution
- Hierarchical design: roles, profiles, classes are assigned to groups of computers. Dev and Prod environments.
- Documentation of IT Infrastructure configuration

📄 **daq.yaml** 📋 2.34 KB

```
1   ---
2   classes:
3     - apel
4     - apel::testing
5     - autofs
6     - profile::service::cephfs_automount
7     - sysctl::base
8     - ssh::client
9     - ssh::server
10
11  apel::testing::enabled: '1'
12  daq_vncserver::home_manage: true
13  daq_fedora::homedir::desktop_bg: '#1b3324'
14  desktop::desktop: 'LXDE'
15
16  autofs::mounts:
17    net:
18      mount: '/net'
19      mapfile: '-hosts'
20    ceph:
21      mount: '/-'
22      mapfile: '/etc/auto.ceph'
23      options: '--timeout=120'
```

📄 **init.pp** 📋 344 Bytes

```
1   #
2   class cvmfs(
3     String $package_release,
4     String $package_release_url,
5     String $package_ensure,
6     Boolean $package_manage,
7     Array[String] $package_name,
8   ) {
9     contain cvmfs::repo
10    contain cvmfs::install
11    contain cvmfs::config
12
13    Class['::cvmfs::repo']
14    -> Class['::cvmfs::install']
15    -> Class['::cvmfs::config']
16    ~> Service['autofs']
17  }
```

```
[root@bmn-evb ~]# puppet agent -vt
Info: Using environment 'production'
Info: Retrieving pluginfacts
Info: Retrieving plugin
Info: Retrieving locales
Info: Loading facts
Info: Caching catalog for bmn-evb.he.jinr.ru
Info: Applying configuration version '1684344730'
Notice: /Stage[main]/Autofs::Service/Service[autofs]/ensure: ensure changed 'stopped' to 'running' (corrective)
Info: /Stage[main]/Autofs::Service/Service[autofs]: Unscheduling refresh on Service[autofs]
Notice: Applied catalog in 9.74 seconds
[root@bmn-evb ~]#
```
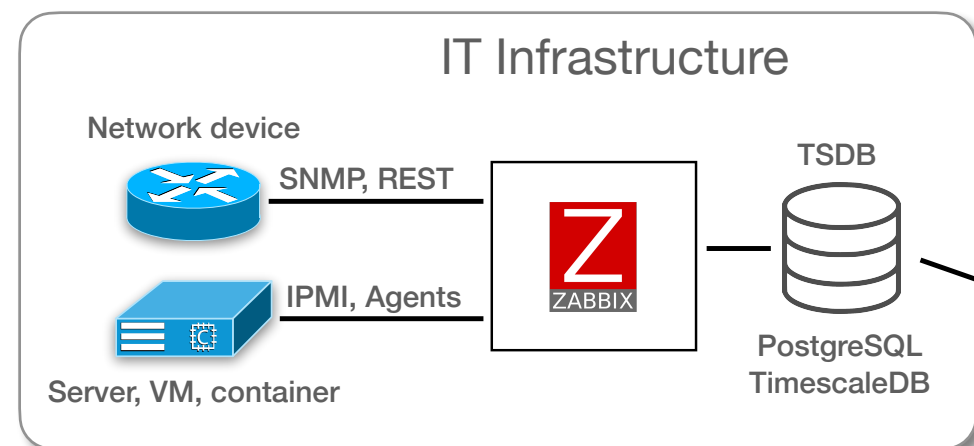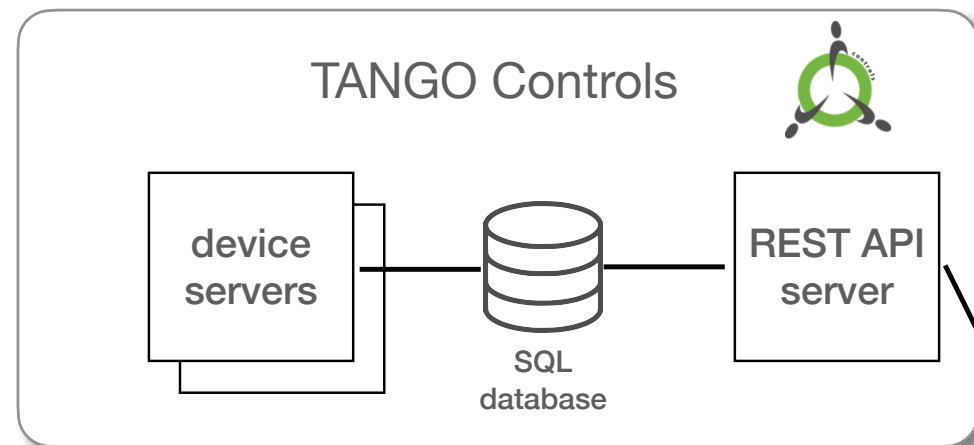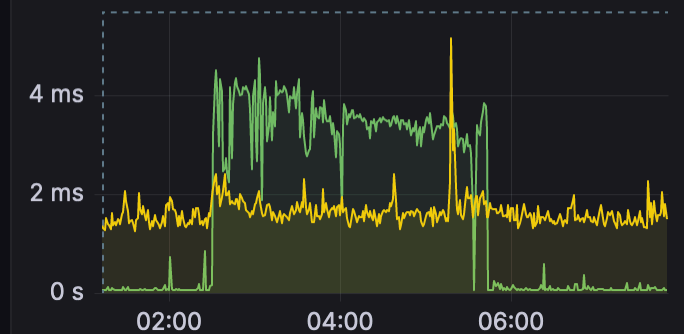
# Monitoring Architecture

## TANGO Controls

device servers — SQL database — REST API server

## IT Infrastructure

Network device
SNMP, REST — ZABBIX — TSDB — PostgreSQL TimescaleDB
IPMI, Agents
Server, VM, container

## DAQ Metrics

DAQ software — TSDB — InfluxDB
Redis key-value broker

Query data
Visualise
Analyse
Alert

Unified dashboards
Data source plugins

Time-Series Data panels

Throughput
4 GB/s
3 GB/s
2 GB/s
1 GB/s
0 B/s
02:00    04:00    06:00

AVG Op Latency
4 ms
2 ms
0 s
02:00    04:00    06:00

Slow Extraction
1.4
1.2
1
0.8
0.6
0.4
0.2
0
0    1    2    3    4
relative time

BMN Beam Counter
[0] BC1L — 1 094 373 hits
700k
600k
500k
400k
300k
200k
100k
0
4.5    5    5.5    6    6.5    7    7.5    8
time

*Plotly* panel for arbitrary data
Graphs and Histograms

# Log Message Processing and Analysis

## Elasticsearch, Logstash, Kibana