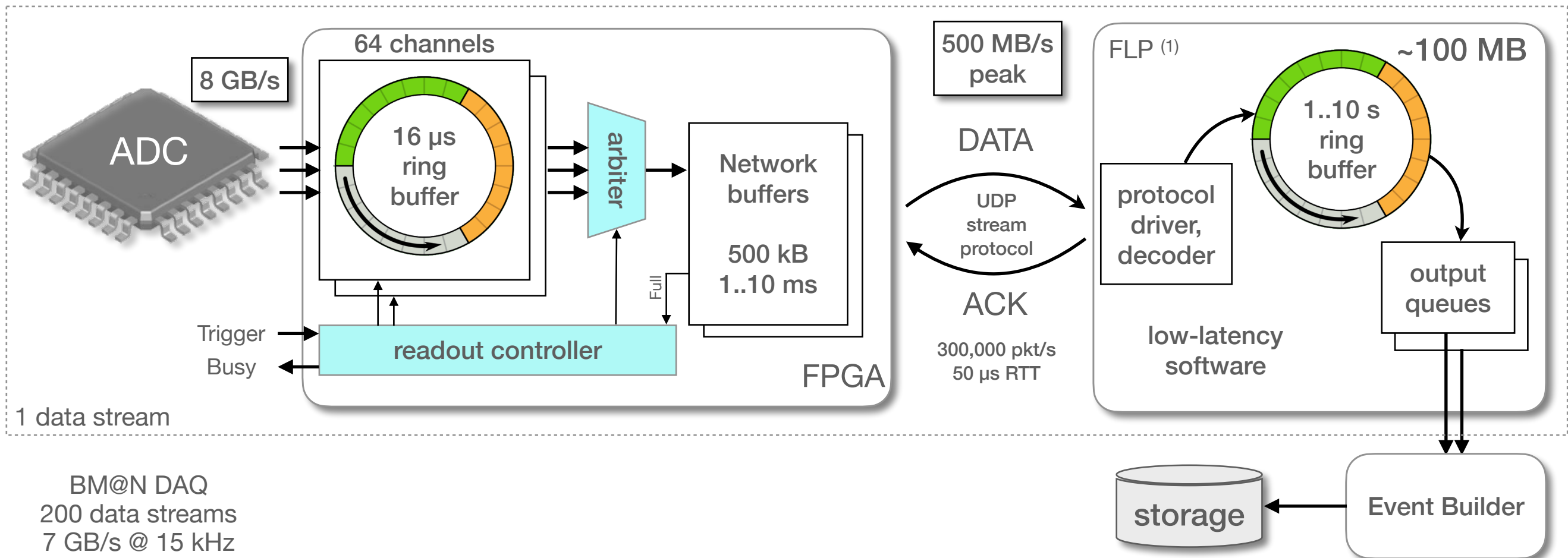# BM@N Data Acquisition IT Infrastructure

**BM@N Experiment at the NICA Facility**
**10th Collaboration Meeting**
**St Petersburg, May 14 – 19, 2023**

**ILIA SLEPNEV, JINR**

# Data Acquisition

## Data transfer from detector to storage system



**8 GB/s**

ADC

64 channels

16 µs ring buffer

arbiter

Network buffers

500 kB
1..10 ms

Full

Trigger

Busy

readout controller

FPGA

500 MB/s peak

DATA

UDP stream protocol

ACK

300,000 pkt/s
50 µs RTT

FLP [1]

~100 MB

protocol driver, decoder

1..10 s ring buffer

low-latency software

output queues

1 data stream

BM@N DAQ
200 data streams
7 GB/s @ 15 kHz

storage

Event Builder

Decouple fast microsecond-scale synchronous acquisition from slow second-scale software processing

Input queue full condition suspends data taking process causing missed physical events
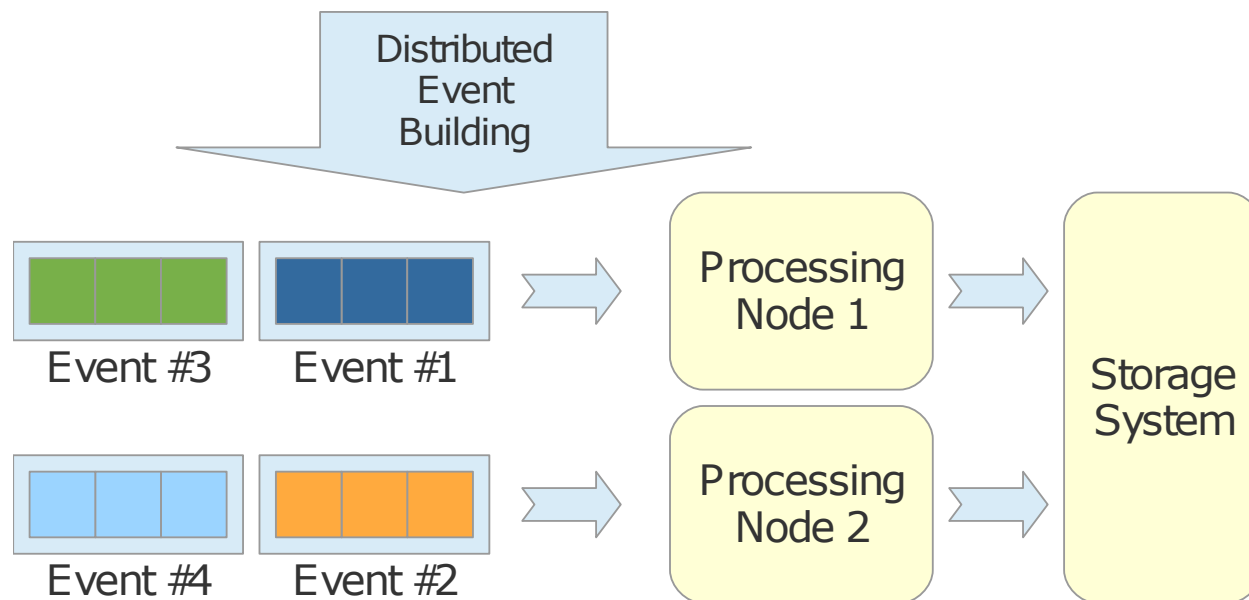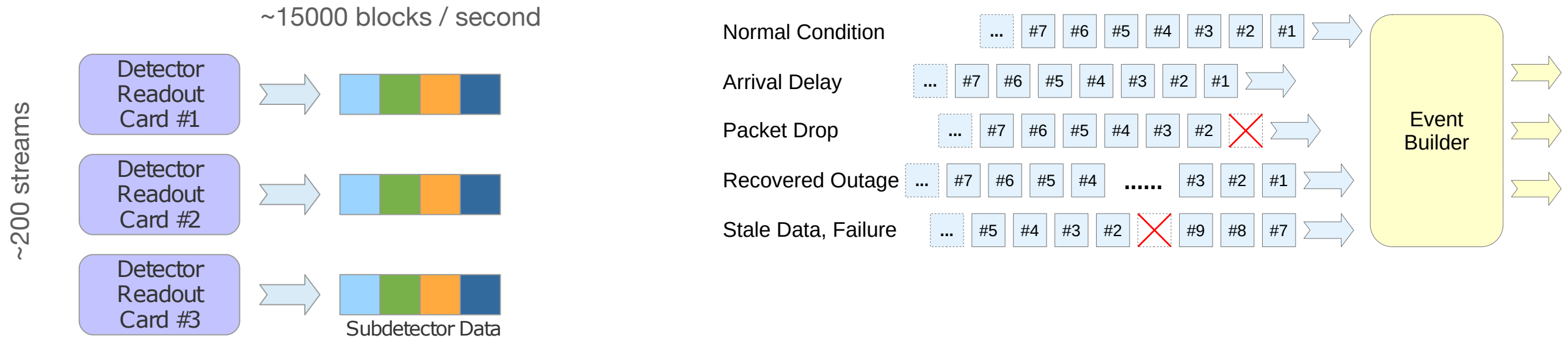Possible solutions
- discrete on-board DRAM memory chips
- custom PCIe data processing cards in FLP
- Ethernet based readout with commodity hardware

(1) FLP — First Level Processor (bare metal server)
Primary task: to receive data from readout electronics, buffer, validate, format and enqueue data blocks ready to be transferred to event building network.

FLP in BM@N is part of synchronous processing, it directly affects readout efficiency

# Data Acquisition

## Distributed Event Building

~15000 blocks / second

~200 streams

Detector Readout Card #1

Detector Readout Card #2

Detector Readout Card #3

Subdetector Data

Distributed Event Building

Event #3   Event #1

Event #4   Event #2

Processing Node 1

Processing Node 2

Storage System

Normal Condition    ...  #7 #6 #5 #4 #3 #2 #1

Arrival Delay    ...  #7 #6 #5 #4 #3 #2 #1

Packet Drop    ...  #7 #6 #5 #4 #3 #2 ✗

Recovered Outage    ...  #7 #6 #5 #4 ...... #3 #2 #1

Stale Data, Failure    ...  #5 #4 #3 #2 ✗ #9 #8 #7

Event Builder

Event Building – process of sorting data fragments from subdetectors and assembling complete event data ready for physical analysis

Reliability: handle single errors, data dropouts, corrupted data, timeouts, detector electronics restarts or servers restarts without process interruption

Event Building in BM@N is part of asynchronous processing, it does not affect readout efficiency under normal conditions
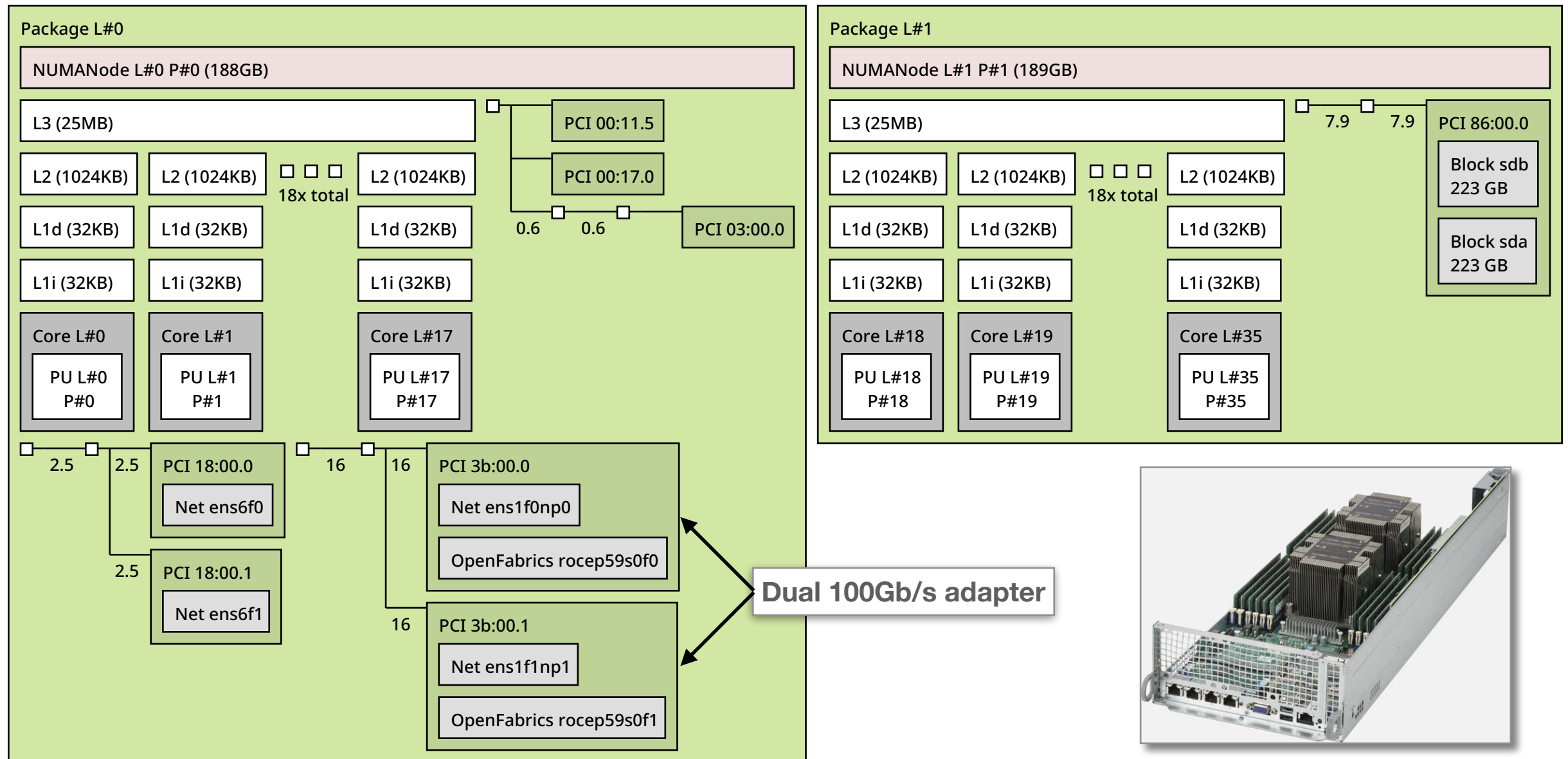
# Data Acquisition
## FLP Hardware Topology

Typical NUMA topology of dual-CPU server used in BMN DAQ Data Center



**stream processing CPU cores**
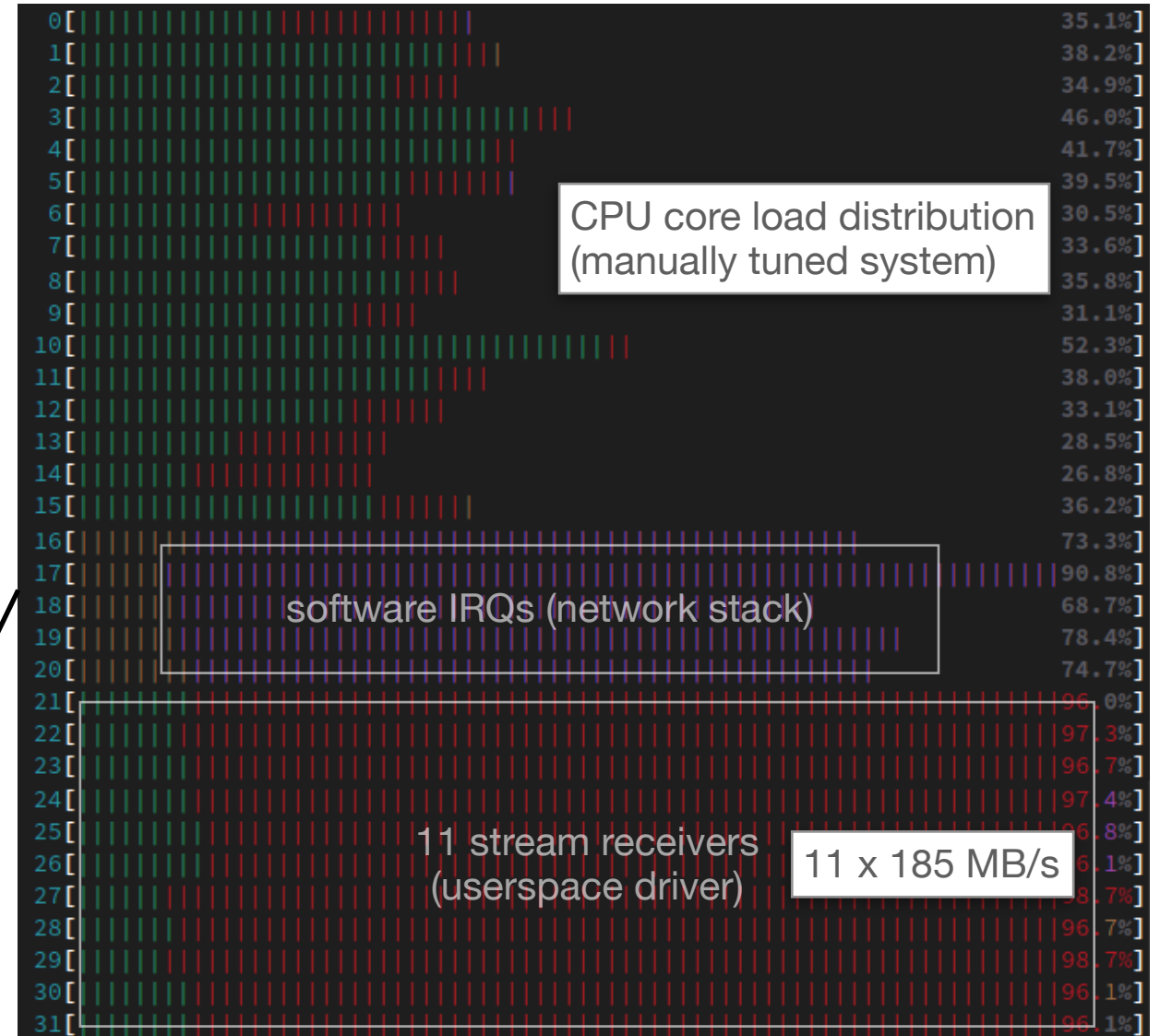
**service CPU cores**

Machine (377GB total)

**Package L#0**

NUMANode L#0 P#0 (188GB)

L3 (25MB)

PCI 00:11.5

L2 (1024KB) | L2 (1024KB) | ☐ ☐ ☐ 18x total | L2 (1024KB)

PCI 00:17.0

L1d (32KB) | L1d (32KB) | L1d (32KB)

0.6 | 0.6 | PCI 03:00.0

L1i (32KB) | L1i (32KB) | L1i (32KB)

Core L#0 | Core L#1 | Core L#17

PU L#0 P#0 | PU L#1 P#1 | PU L#17 P#17

2.5 | 2.5 | PCI 18:00.0 | 16 | 16 | PCI 3b:00.0

Net ens6f0 | Net ens1f0np0

OpenFabrics rocep59s0f0

2.5 | PCI 18:00.1

Net ens6f1

16 | PCI 3b:00.1

Net ens1f1np1

OpenFabrics rocep59s0f1

**Dual 100Gb/s adapter**

**Package L#1**

NUMANode L#1 P#1 (189GB)

L3 (25MB)

7.9 | 7.9 | PCI 86:00.0

L2 (1024KB) | L2 (1024KB) | ☐ ☐ ☐ 18x total | L2 (1024KB)

Block sdb 223 GB

L1d (32KB) | L1d (32KB) | L1d (32KB)

Block sda 223 GB

L1i (32KB) | L1i (32KB) | L1i (32KB)

Core L#18 | Core L#19 | Core L#35

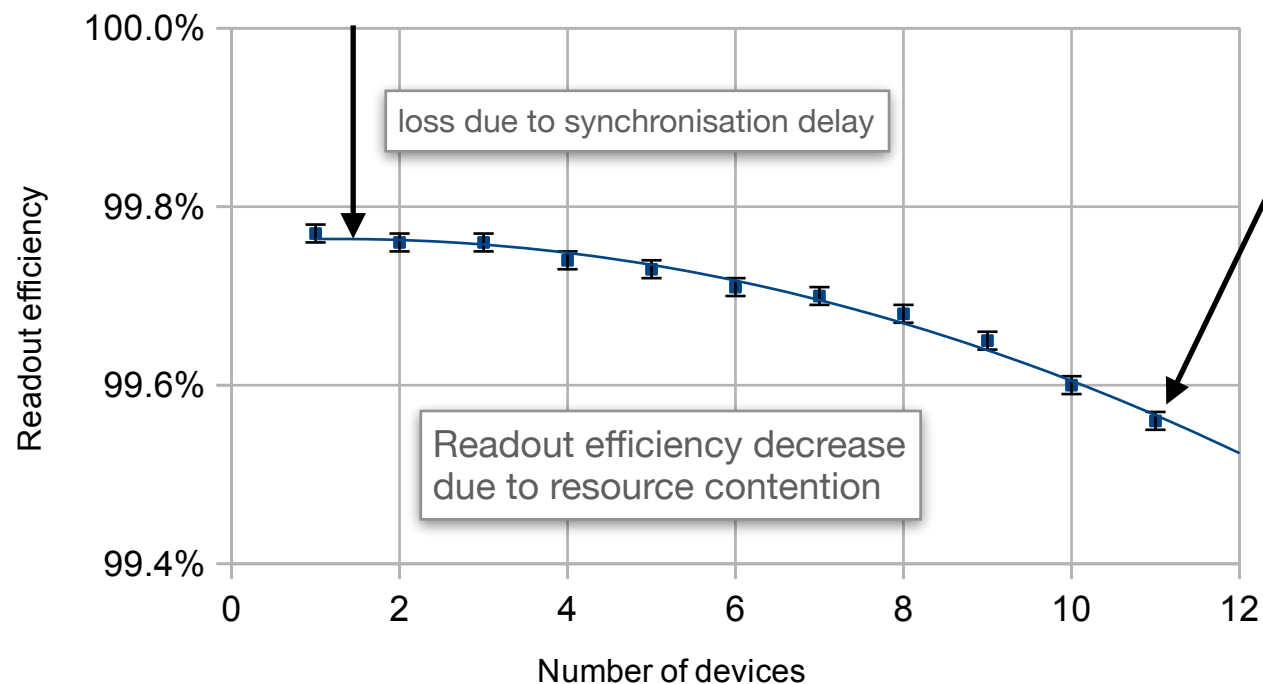PU L#18 P#18 | PU L#19 P#19 | PU L#35 P#35

Host: c5n20.he.jinr.ru
Date: Wed 03 May 2023 11:15:52 PM MSK

# Data Acquisition

## OS Tuning for Real-Time

50 kHz random trigger
4.8 µs programmed dead time
…
10 kHz effective rate

Default system task scheduler is not optimal. Manual tuning is required.



CPU core load distribution (manually tuned system)

software IRQs (network stack)

11 stream receivers (userspace driver)

11 x 185 MB/s



loss due to synchronisation delay

Readout efficiency decrease due to resource contention

Readout efficiency (y-axis): 100.0%, 99.8%, 99.6%, 99.4%
Number of devices (x-axis): 0, 2, 4, 6, 8, 10, 12

CPU cores assignment (11 streams test case)

0-15 — available to system scheduler
16-20 — software IRQ (network adapter jobs)
21-31 — stream receivers (userspace driver)

Platform: Supermicro X10DRT-PT, Dual Intel Xeon E5-2697A v4
Network: Mellanox ConnectX-5 100G
Readout modules: TQDC16VS-10G, f/w v2.9

Tuning
# tuned-adm profile network-latency

4 CPU cores per data stream

Linux kernel 6.0.15
mitigations=off
intel_idle.max_cstate=1
idle=poll
skew_tick=1
isolcpus=16,17,…,31

https://access.redhat.com/documentation/en-us/red_hat_enterprise_linux_for_real_time/9

# IT Infrastructure
## Compute Resource Virtualisation

| Class | Task | Requirements | Where to run |
|---|---|---|---|
| Near real-time | FLP, data stream receiver | UDP protocol processing minimal latency, non-interruptible | Bare Hardware or LXC Manual tuning required |
| High Throughput | Event Building | High TCP/IP throughput, large buffers in RAM | Bare Hardware or LXC |
| Lightweight | Detector Control, Web services | Ease of Management | VM or container *KVM, LXC, Docker-in-KVM* |

What should be tuned for real-time
- CPU power and frequency management
- Network adapter interrupt coalescing
- Disable firewall
- Process and soft-IRQ affinity, CPU core isolation

# IT Infrastructure

## BM@N DAQ Network

Network bandwidth
External: 200 Gb/s (400 planned)
Fabric: 200 Gb/s (extension possible)
Servers: 100 Gb/s (partially redundant)
Electronics: 1 and 10 Gb/s, no redundancy



Technological Network
core routers
2 x Nexus 9336

Cisco ACI Fabric

Spine switches
2 x Nexus 9364C

100Gb/s

100G Leaf switches
2 x Nexus 9336

10G Leaf switches
7 x Nexus 93180

100Gb/s

10Gb/s

10Gb/s

SSD Ceph
10 nodes

HDD Ceph
12 nodes

FLP, EvB
10 nodes

Proxmox
10 nodes

service
6 nodes

10 x Aruba 3810M

Slow Control,
Control Room

Readout electronics
10G links

Software Defined Storage: CephFS
HDD: 2.2 PB (EC-replicated)
SSD: 100 TB (triple replicated)

Compute resources
CPU: 720 cores total
RAM: 7500 GiB

Network Redundancy and High Availability
• Network is base component of data taking process and is absolutely critical
• Thousands of devices communicate, loss of control or monitoring is critical
• Long data taking runs, no maintenance windows for critical upgrades and improvements
• Automated response to hardware and software failures, minimal operator intervention

# Infrastructure Management

## Data sources

### NetBox, openDCIM

- Models hardware and networks
  - IP address space allocation
  - MAC, IP address, VLAN, User
  - Hardware items management
  - Cable and connections tracking
- Single source of truth. Not a monitoring
- Data source for automation tools
  - DHCP, DNS, RADIUS configs
- Documentation and reporting
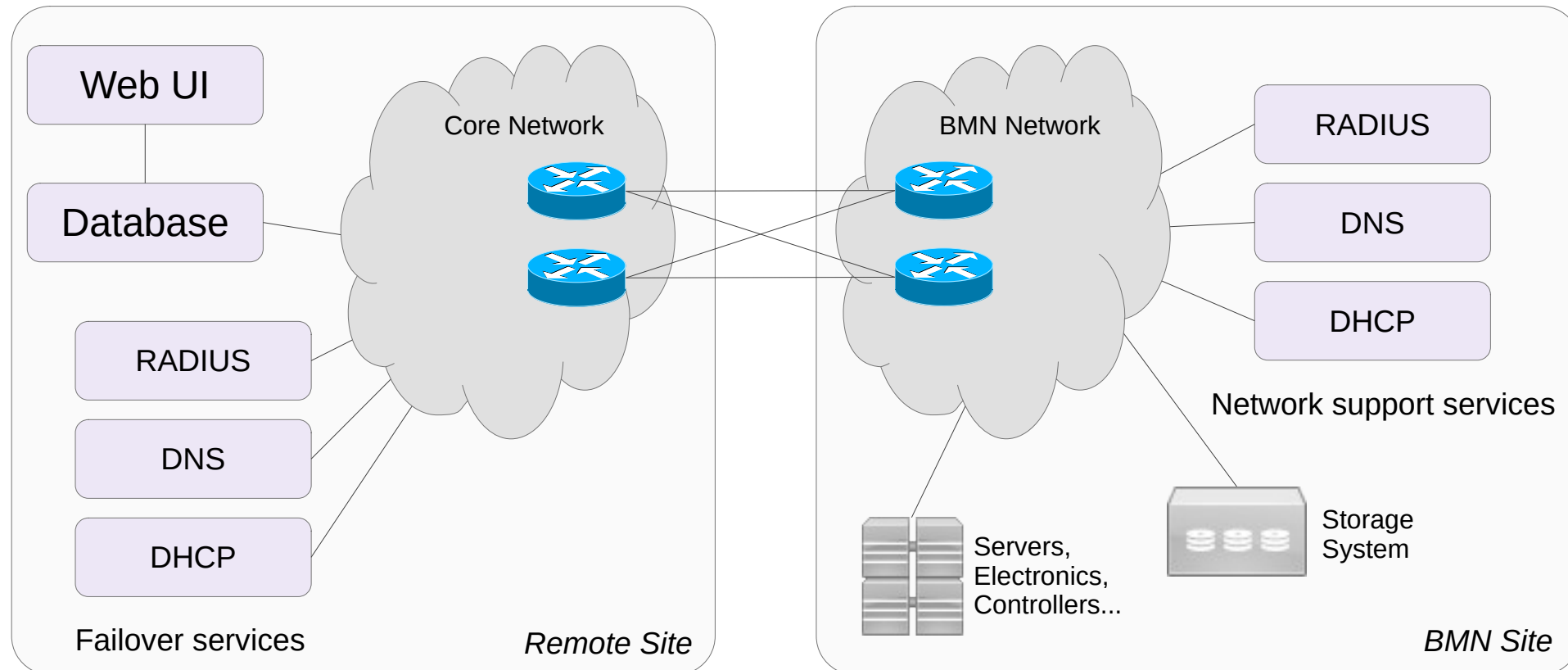- Existing and planned infrastructure

# Infrastructure Management
## High Availability of Network Services

- Database is source of configuration data, but not critical for operation
- Service configuration is updated periodically from database
  - Configuration is validated and service is restarted
  - Running service is announced to routers with IP Anycast address
- Failed service is excluded and traffic is routed to next available service automatically



Network support services
- DHCP: assignment of IP address by looking up database with hardware address
- RADIUS: dynamic assignment of virtual network, IEEE 802.1x
- DNS: hostname to IP address translation

# 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.
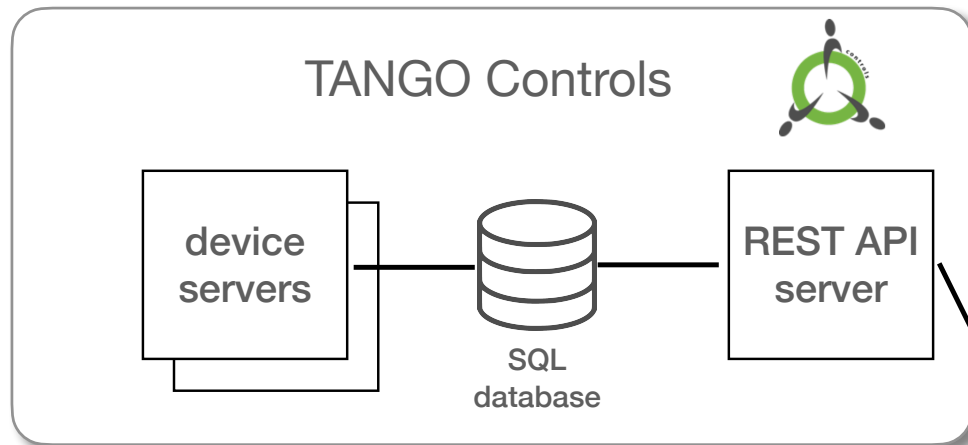- Documentation of IT Infrastructure configuration

📄 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  }
```
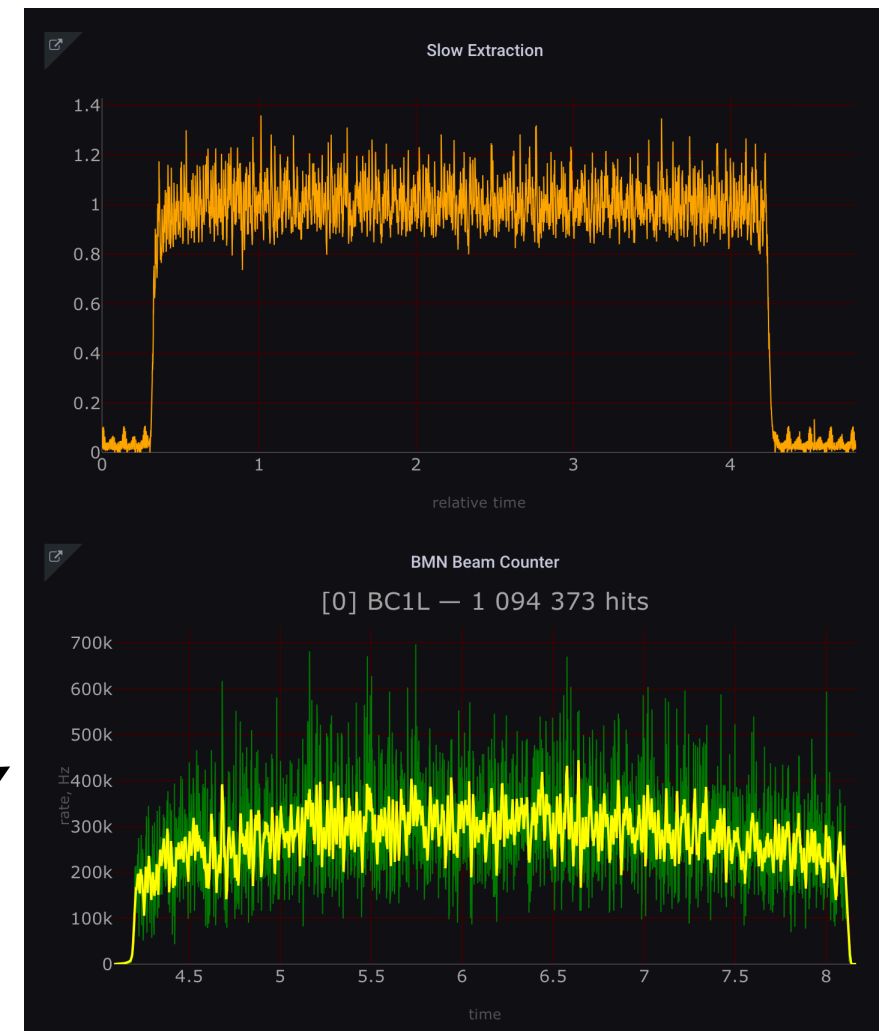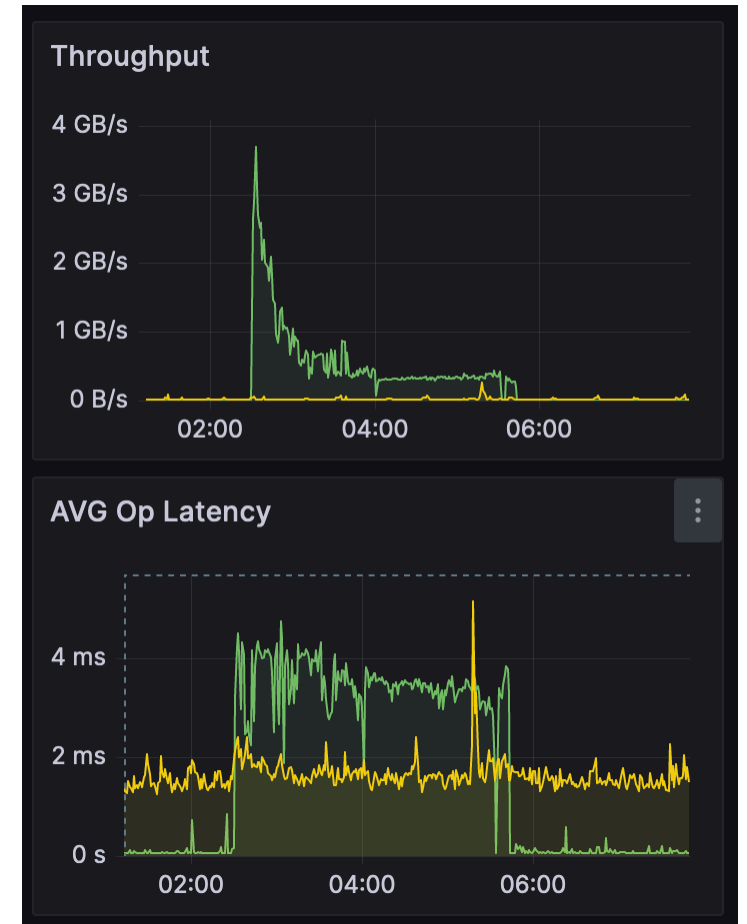
```
[[root@bmn-evb ~]# puppet agent -vt
Info: Using environment 'production'
Info: Retrieving pluginfacts
Info: Retrieving plugin
Info: Retrieving locales
Info: Loading facts
Info: Caching catalog for bmn-evb.he.jinr.ru
Info: Applying configuration version '1684344730'
Notice: /Stage[main]/Autofs::Service/Service[autofs]/ensure: 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 ~]#
```

# Observability: metrics
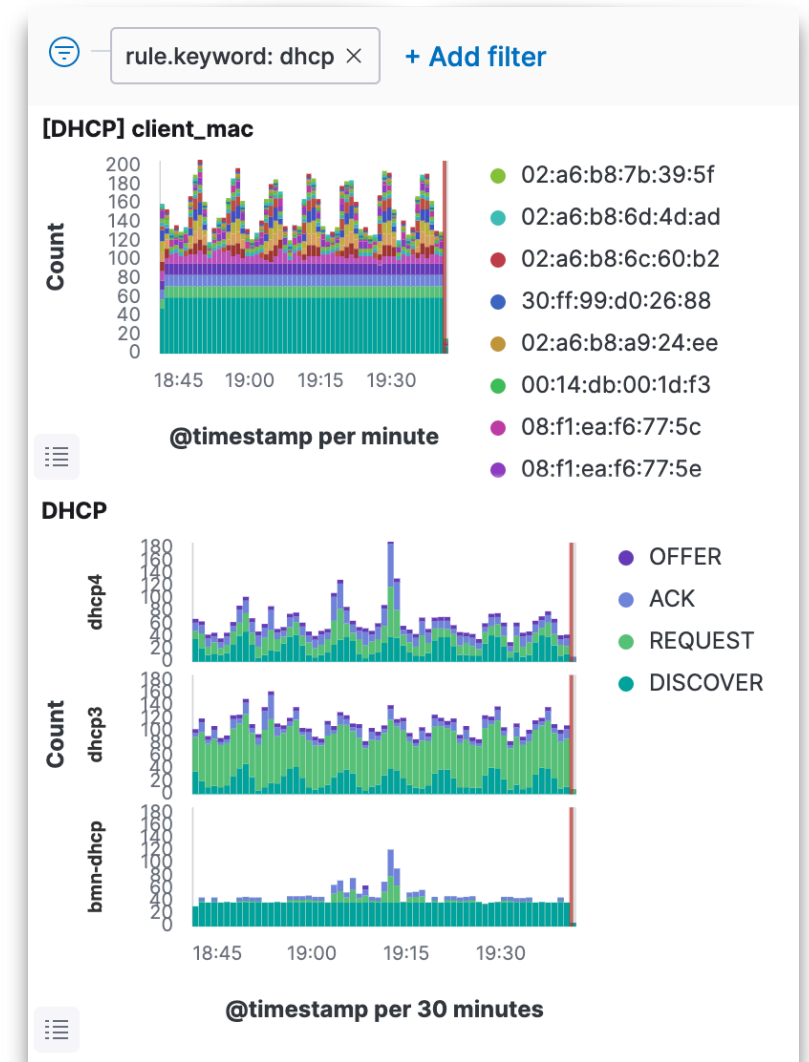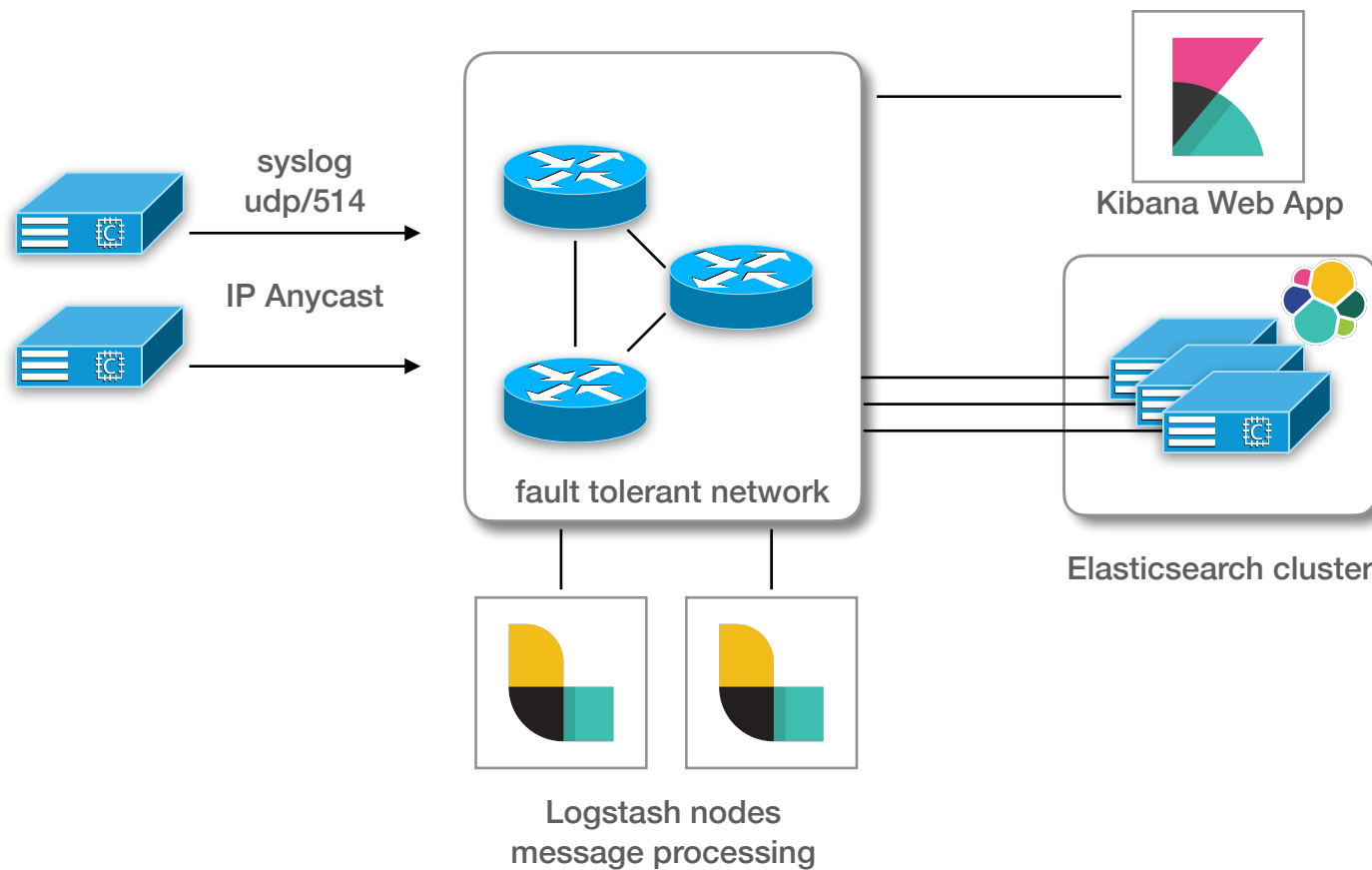## Grafana Unified Dashboards



TANGO Controls

device servers — SQL database — REST API server

IT Infrastructure

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

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    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

# Observability: logs
# Elasticsearch, Logstash, Kibana

# #TODO

DAQ Infrastructure was adequate for past BM@N runs, no critical issues. However…

- Improve DAQ readout efficiency. Considering intermediate hardware buffer based on CRU16.
- New detectors, thus more data streams, additional FLP nodes are required.
- Trigger rate is expected to rise on next run, increasing total data throughput. More Event Builder nodes are required.
- Extend storage space for autonomous operation.
- BMN DAQ Data Center is at limit of cooling power. Considering to move storage system to building 14
- Routine maintenance for all information systems. OS and application upgrades.
- Last, not least. Finish integration with JINR SSO authentication.