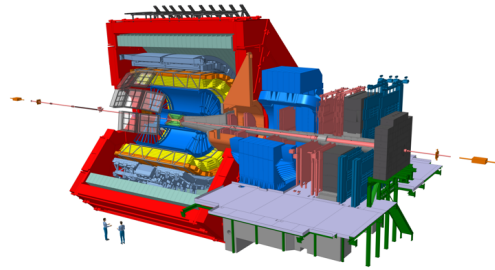# Towards the integrated ALICE Online-Offline monitoring subsystem

Adam Wegrzynek

for the ALICE Collaboration

CHEP 2018
Sofia, Bulgaria

# ALICE O²

19 detectors

3.3 TB/s      9000 fibers

**First Level Processors**      270 nodes

Synchronous

500 GB/s

**Event Processing Nodes**      1500 nodes

100 GB/s

Asynchronous
post-processing      **Storage**      60 PB storage
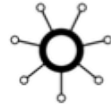
# Comparison

**1.**

*Modular stack*

**2.**

**3.**

- ▶ Performance requirements

- ▶ Functional architecture

- ▶ Experience at CERN

# 1. *Modular stack*

1. **collect**d
   ▸ System performance metrics
   ▸ Hardware monitoring

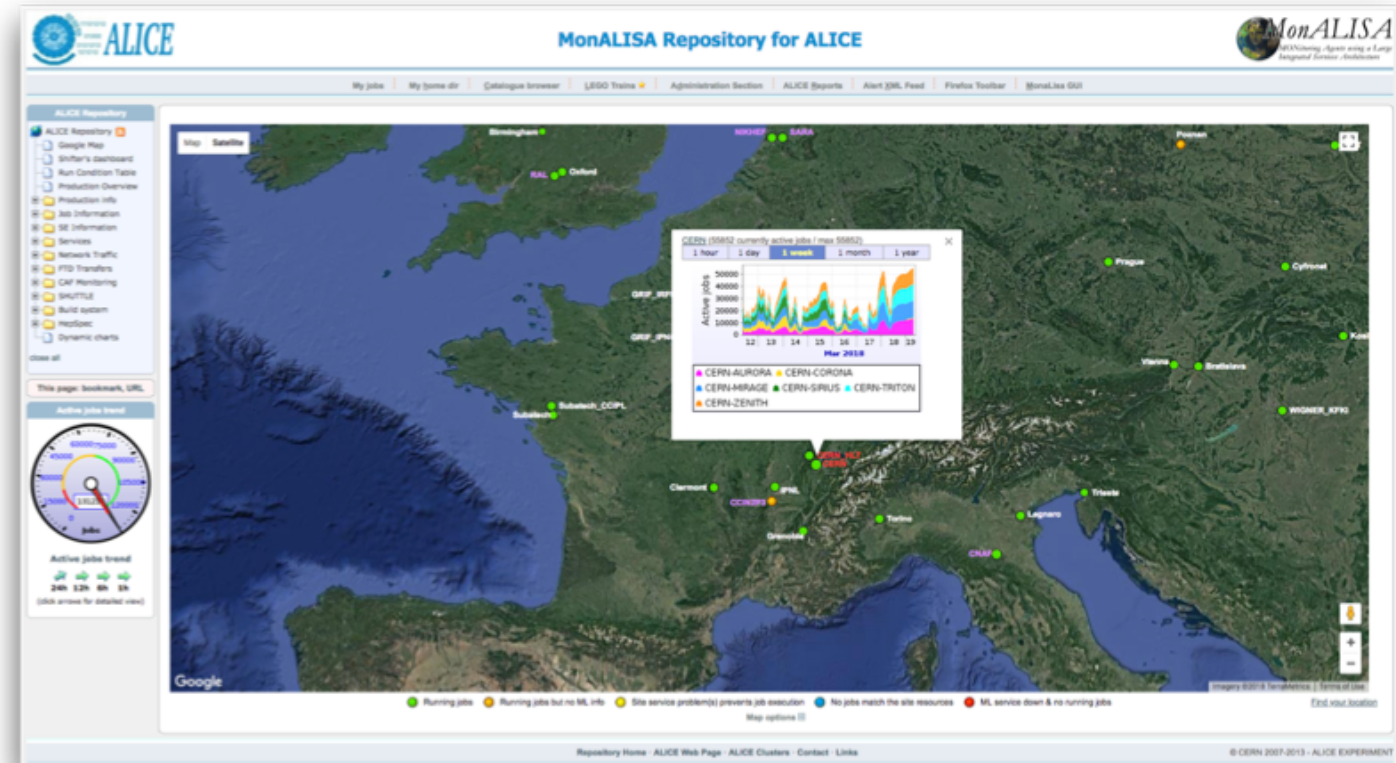2. 
   ▸ Metric routing

3. 
   ▸ In memory data processing

4. 
   ▸ Time series database

5. 
   ▸ Visualization tool

6. 
   ▸ Alarming

   ▸ Currently used at INFN Bari, CERN IT

# 2. MonALISA

▶ Distributed data collector infrastructure

▶ Discovery mechanism

▶ Aggregation, filtering, alerts

▶ Real-time data distribution

▶ In memory buffers

▶ SQL database

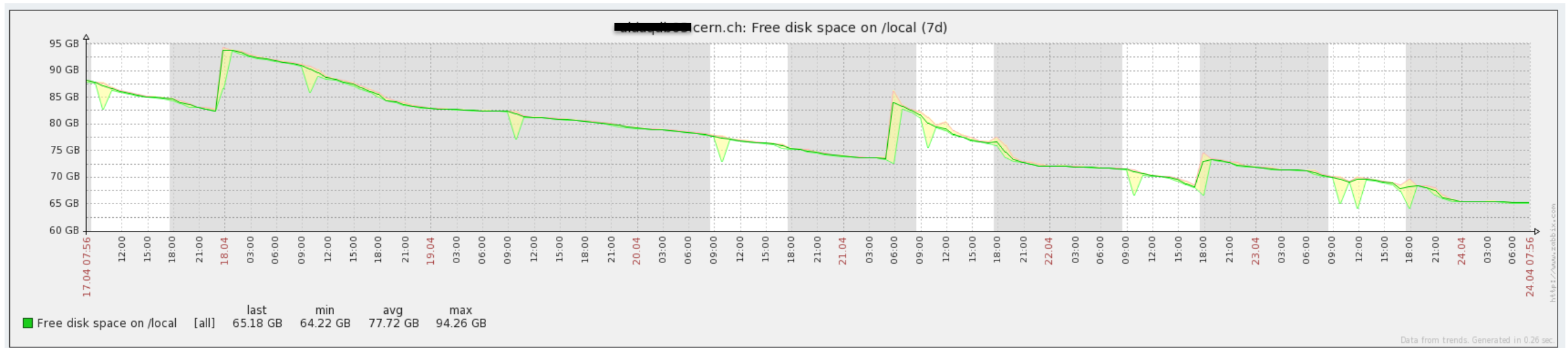▶ Currently used by ALICE Offline



Courtesy of **Costin Grigoraş**

# 3. Zabbix

- ▶ Agent-server

- ▶ Push via Zabbix protocol

- ▶ Community support

- ▶ Currently used in ALICE **HLT** and **DAQ** for computing infrastructure monitoring

# Comparison table (1)

| | Modular Stack | | MonALISA | | Zabbix |
|---|---|---|---|---|---|
| Reference OS (CC7) | | Yes | | Yes | | Yes |
| Documentation | | Good | | Insufficient | | Good |
| Support and maintenance | | Yes | | Yes | | Yes |
| Running in isolation | | Yes | | Yes | | Yes |
| 600 kHz rate | | Yes | | Yes | | No |
| Scalable >>600 kHz | | Yes | | Yes | | No |
| Handle 100k sources | | Yes | | Yes | | No |
| Storage size | | ~30 bytes | | ~75 bytes | | 90-500 bytes |

# Comparison table (2)

| | Modular Stack | | MonALISA | | Zabbix | |
|---|---|---|---|---|---|---|
| Functional arch. | | | | | | |
| System sensors | | Yes | | Yes | | Yes |
| Metric processing | | Batch and stream | | Stream | | Batch |
| Historical dashboard | | Yes | | Yes | | Yes |
| Real-time dashboard | | No (RFC) | | Yes (obsolete) | | No |
| Alarming | | Yes | | Yes | | Yes |
| Storage downsampling | | Yes | | Yes | | Yes |

# Selection

**1.**

*Modular stack*

**2.**

Remains for Grid job monitoring

**3.**

# *Modular stack* metric flow

**Computing node**

**Monitoring backend**

System sensors

Application monitoring

Collection, processing

Storage

Visualization
Real-time, historical

Alarming

CollectD

Processing device

Processing device

Processing device

...

Processing device

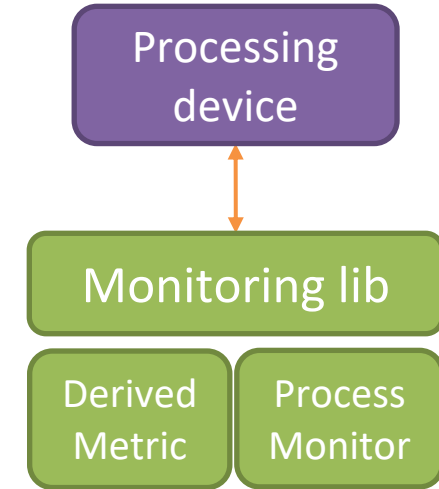Monitoring lib

Derived Metric

Process Monitor

Flume

InfluxDB

Spark

Riemann

Grafana

# Monitoring library

▶ Push metrics to a backend

▶ Monitors the process

▶ Derived metrics

▶ Tags

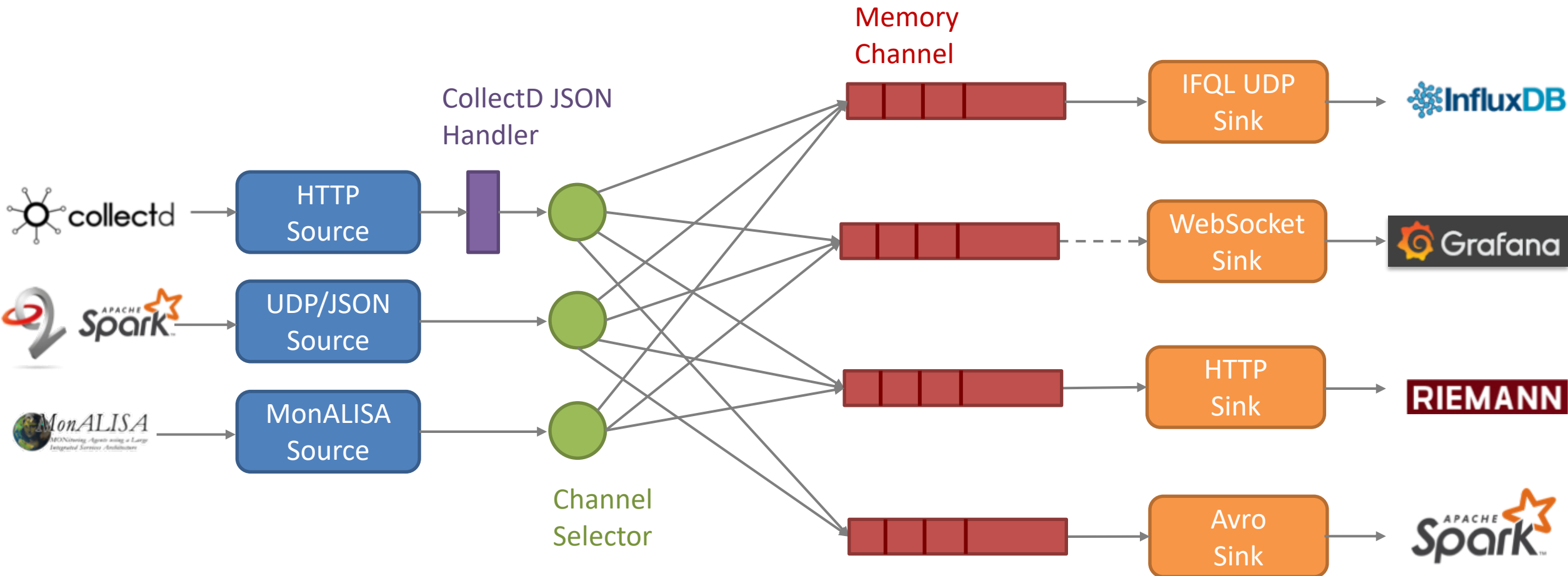▶ AliceO2Group/Monitoring



**myMetric**,0 **10** 1530099250985     hostname=test.cern
                                        role=readout
                                        detector=TPC
                                        ....

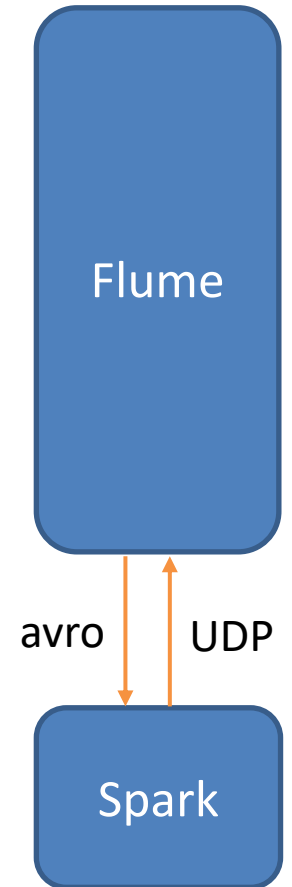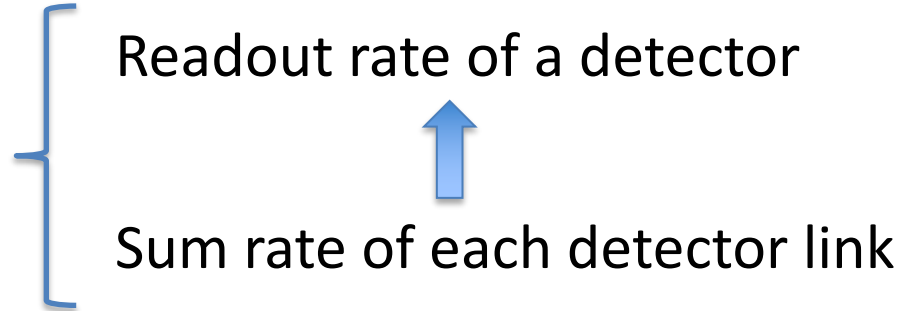name     type     value     timestamp     tags

# Flume routing



Courtesy of **Gioacchino Vino**

# Spark jobs

▶ Higher level metrics

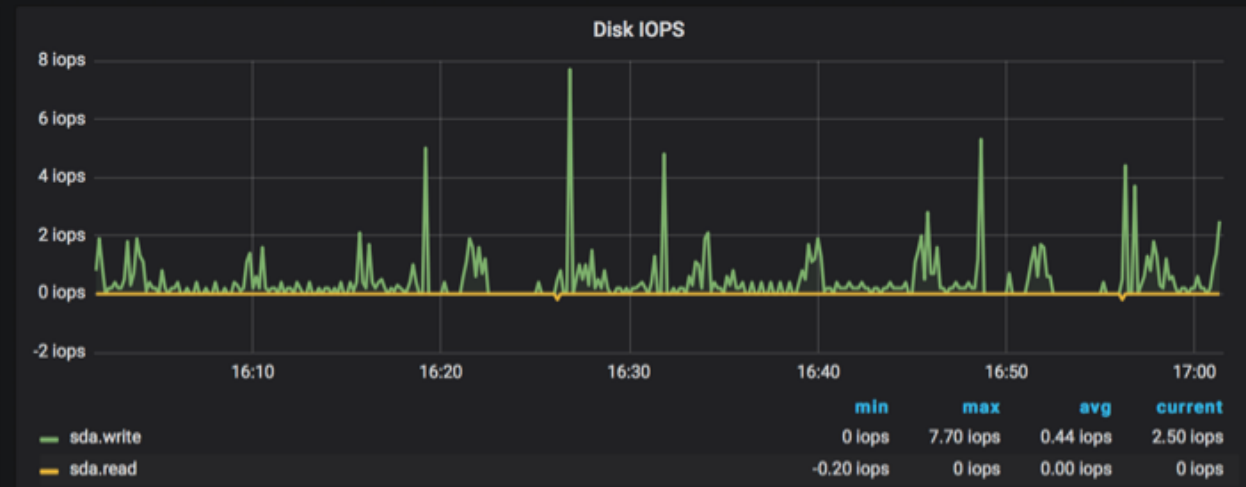▶ Written in Scala

▶ Operates on Flume events

▶ Configurable

Readout rate of a detector

Sum rate of each detector link

Flume

avro    UDP

Spark

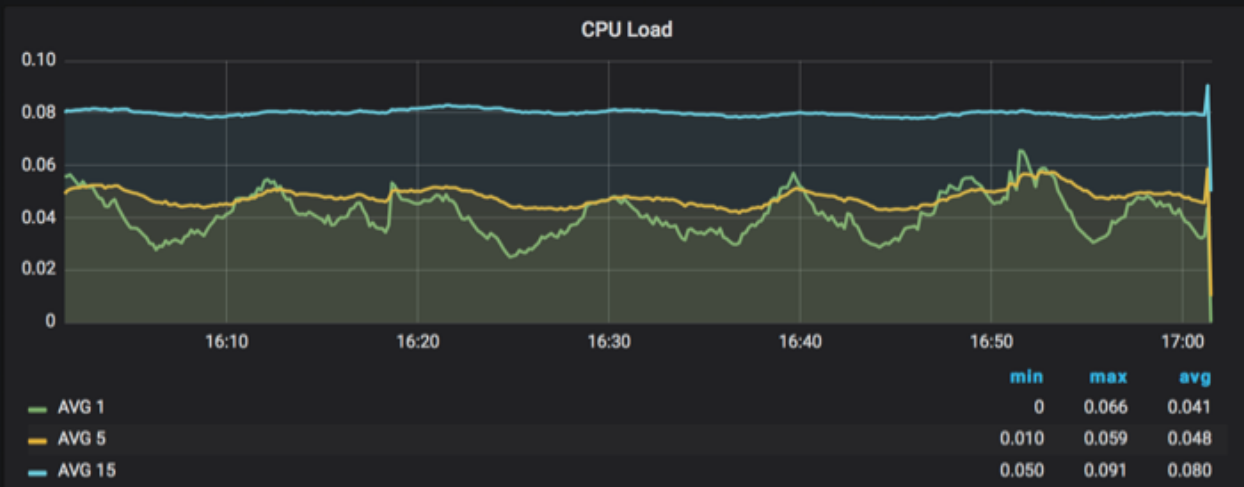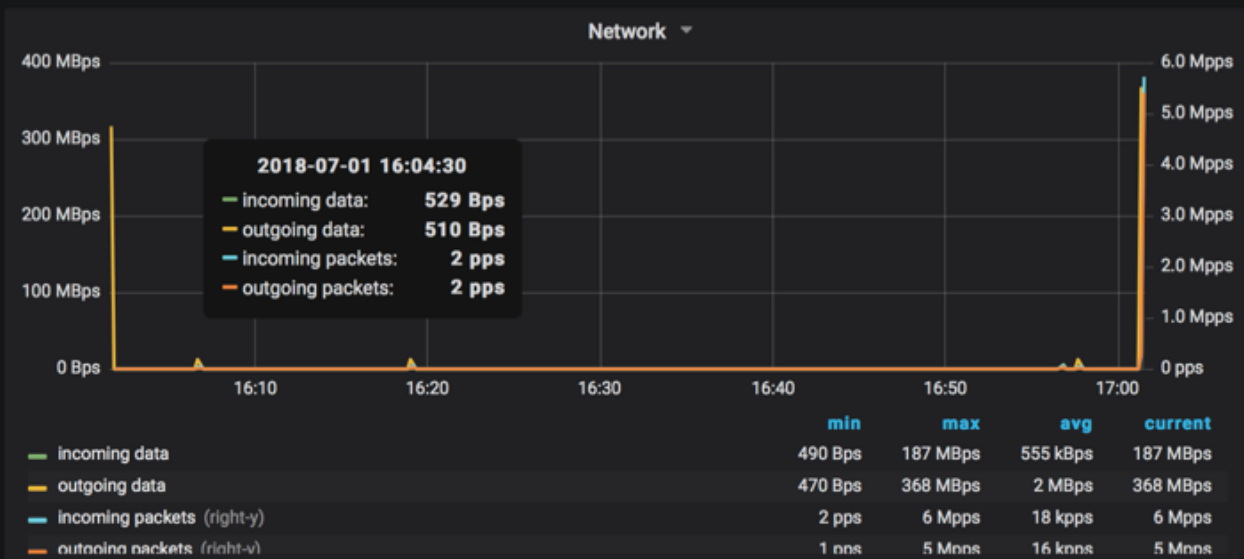# InfluxDB timeseries storage

▶ Up to 700 kHz writes

8 data streams

2x SSD drives RAID0, 25 GbE

▶ Continuous Queries

Downsample high resolution data

(merge 12 points into 1 by applying average)

▶ Retention Policies

Drop high resolution data after 30 days

Keep low resolution data for 1 year

# Grafana

# Integration with O$^2$ Software

▶ Quality Control

▶ Data Processing Layer

Evolution of the ALICE Software Framework for LHC Run 3

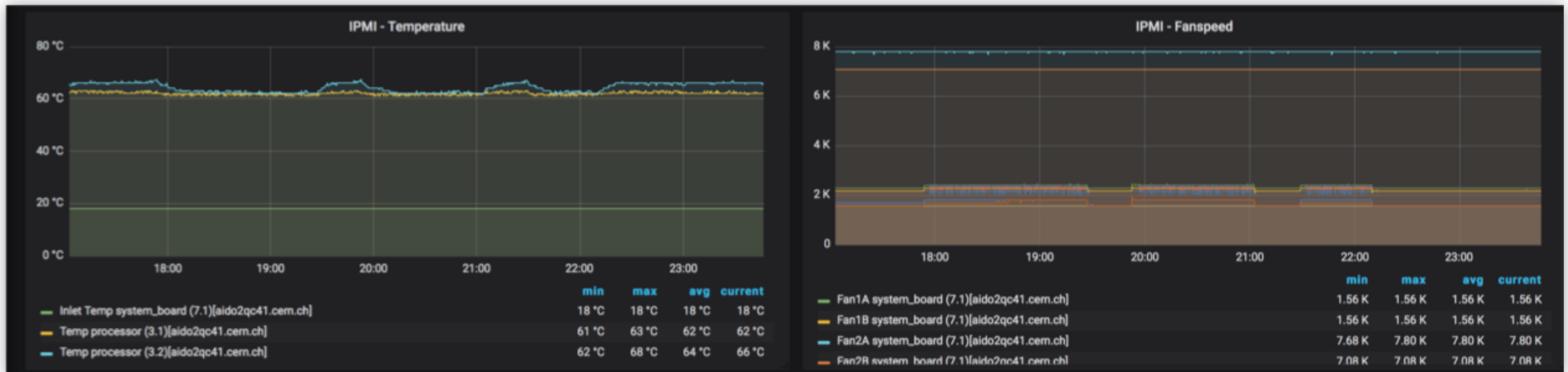Giulio Eulisse, Tuesday 10 July 14:15, Hall 3

▶ Readout

Readout software for the ALICE integrated Online-Offline (O2) system

Filippo Costa, Thursday 12 July 11:00, Hall 3.1

# Conclusion

- ▶ 3 options compared

- ▶ Modular Stack selected for $O^2$ farm monitoring

- ▶ Defined interfaces between tools

- ▶ Deployed in the detector commissioning facilities

# Future steps

- ▶ Alarming

  - ▶ Define thresholds and patterns

- ▶ Grafana real-time data source

  - ▶ Display critical metrics in real time

- ▶ Sensors to custom hardware

  - ▶ Monitor status of custom FPGA board