

# Monitoring Infrastructure for the CERN Data Centre

Asier Aguado

17/05/2018



# Monitoring Mission

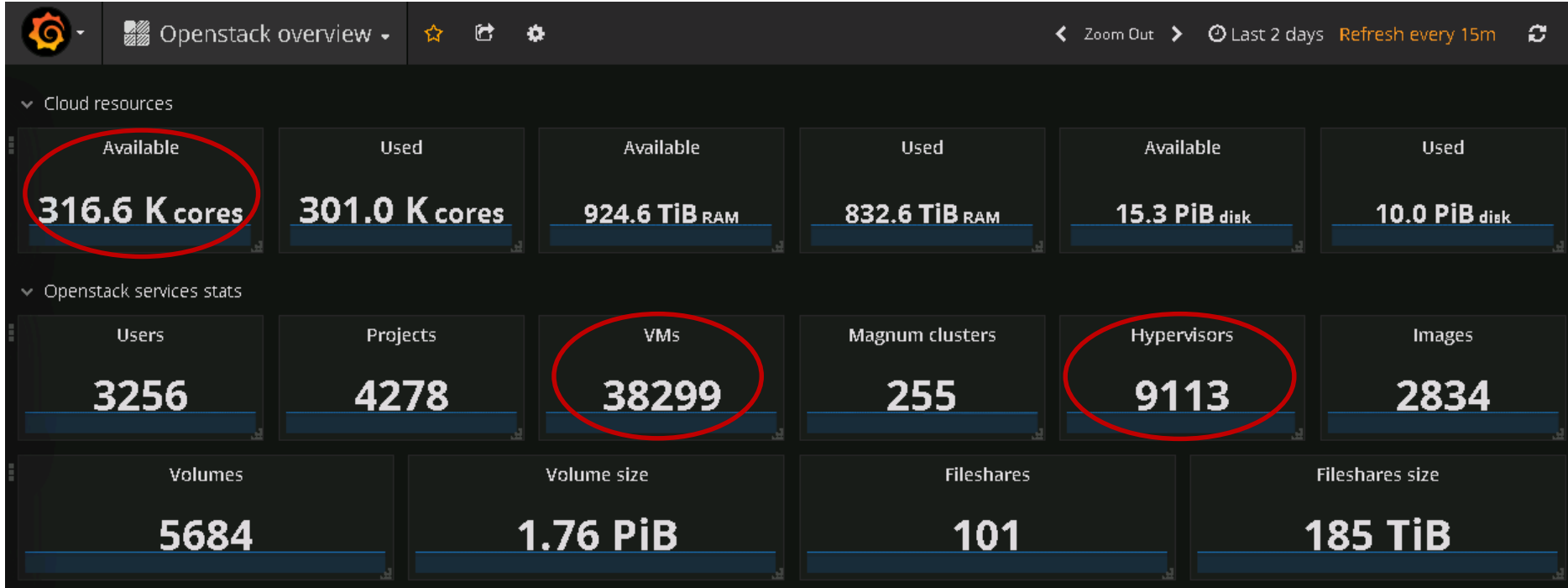
- Provide Monitoring as a Service for CERN Data Centre and the WLCG collaboration
  - e.g. Dashboards, Alarms, Search, Archive
- Collect, transport, store and process metrics and logs for applications and infrastructure

# History

Two use cases at CERN had their own monitoring solution

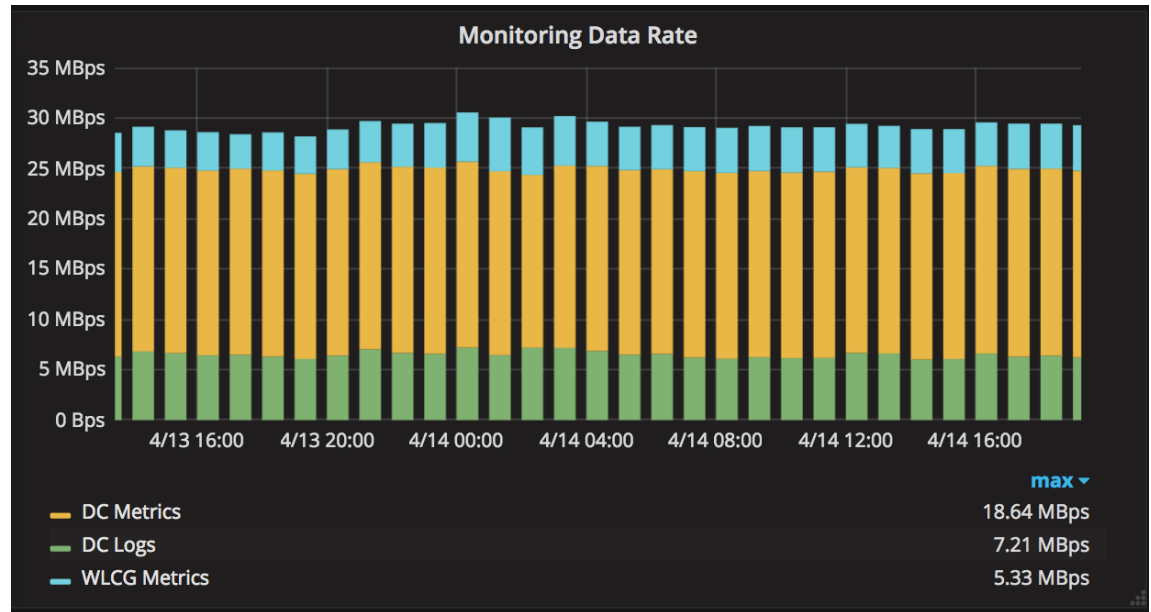
- LEMON and SLS for Data Centre and IT services
- Monitoring tools for WLCG jobs and transfers
- Mission: integrating all of them into a single monitoring solution (MONIT)

# Data Centre



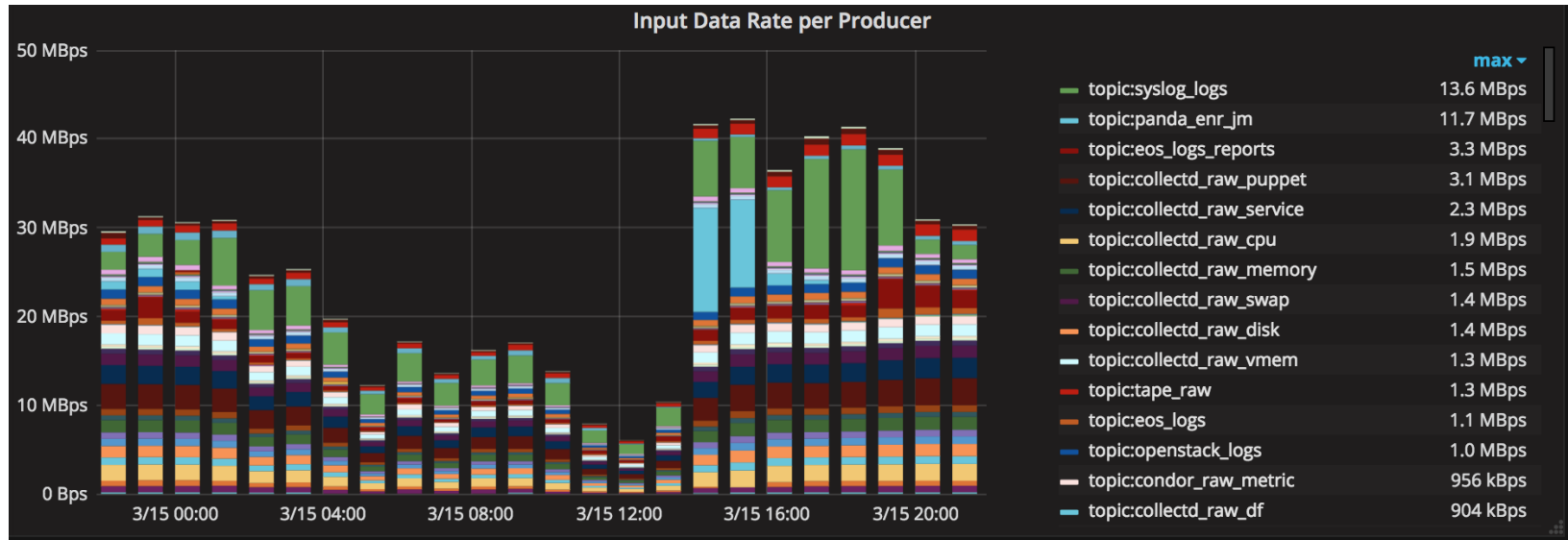
# Monitoring Data

- DC Metrics include:
  - OS metrics
  - Hardware metrics
  - Application metrics
- from > 35k machines
- ~ 3 TB/day (compressed)
- > 80kHz rate



# Monitoring Data - Workload

- Spikes in rate and volume
- Temporary outages (e.g. network outages) cause spikes



# Architecture

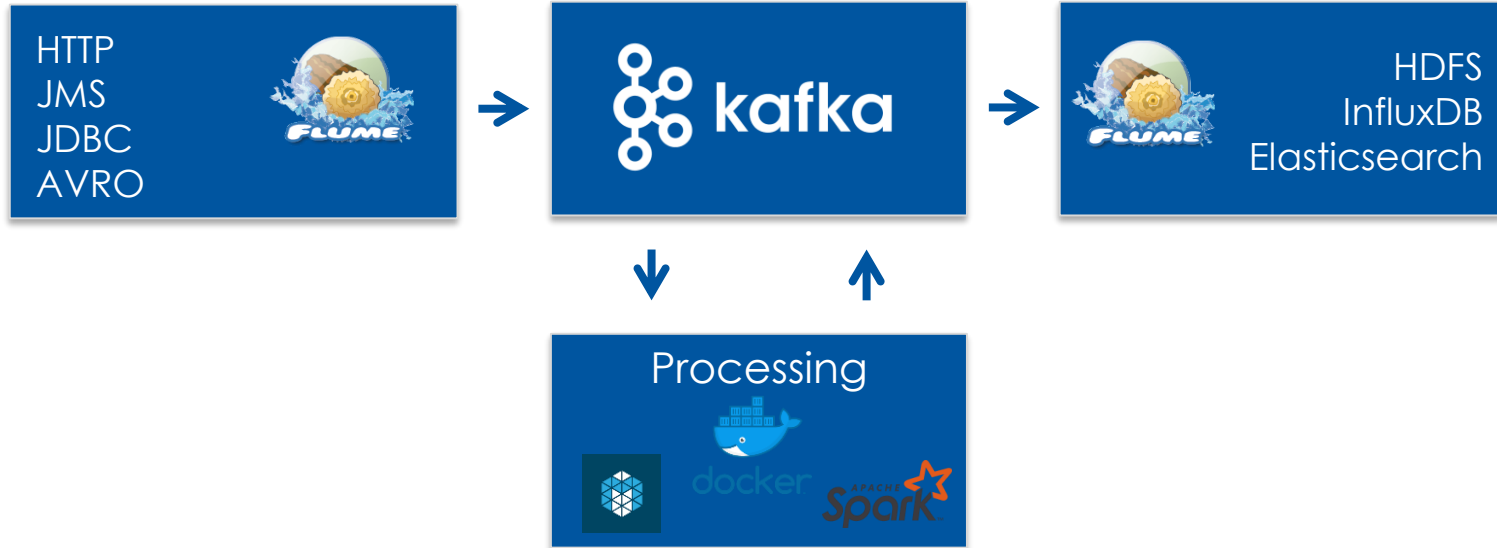
# Technologies

- Moving to community built technologies
- Reduce in-house components when possible
- Collaboration with the community:
  - Development, bug reports...

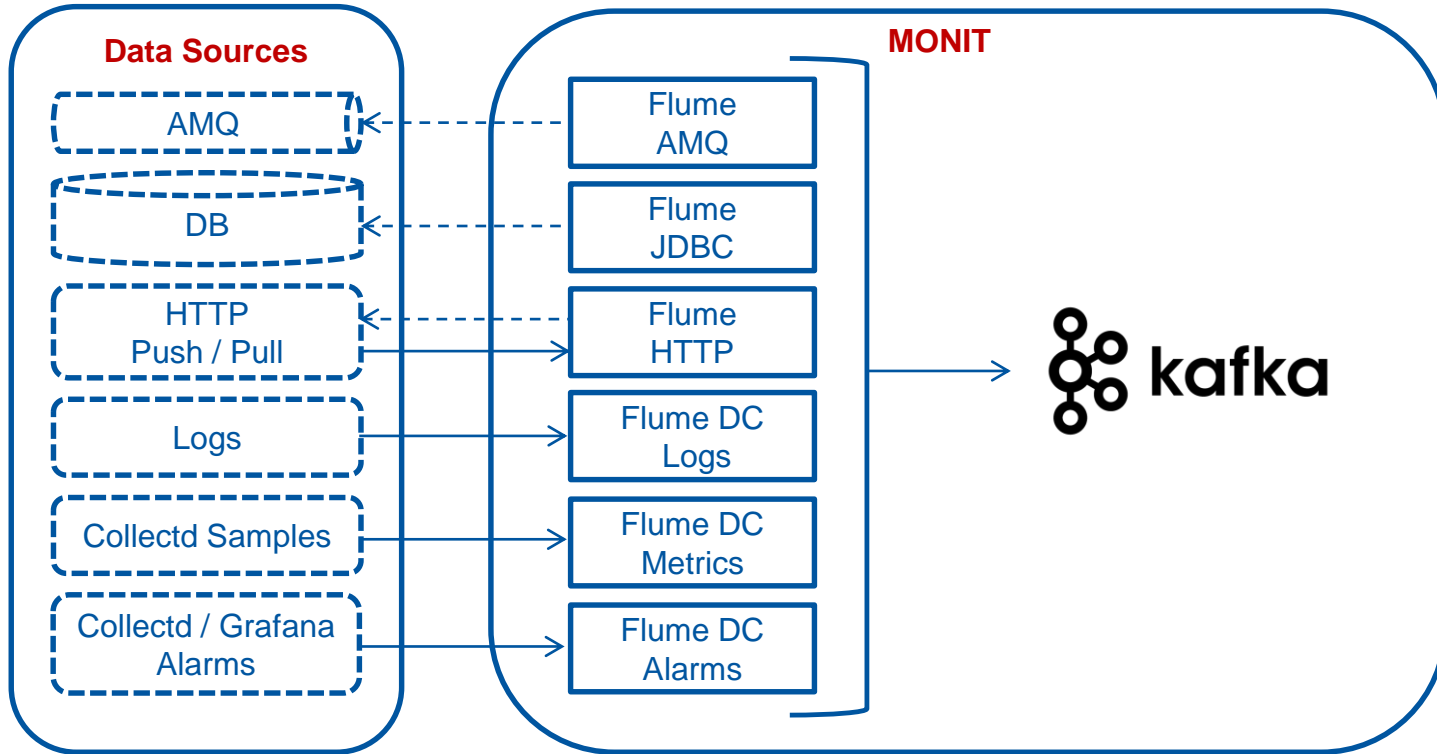




# Monitoring Architecture



# Data Sources



# Data Sources: collectd

- What is collectd

*collectd is a **daemon** which **collects system and application performance metrics** periodically and provides mechanisms to store the values in a variety of ways*



- Why we use collectd

- Modular and easy to deploy
- Community built plugins
- Easy to develop new plugins
- Continuously improving
  - (but documentation could be better)

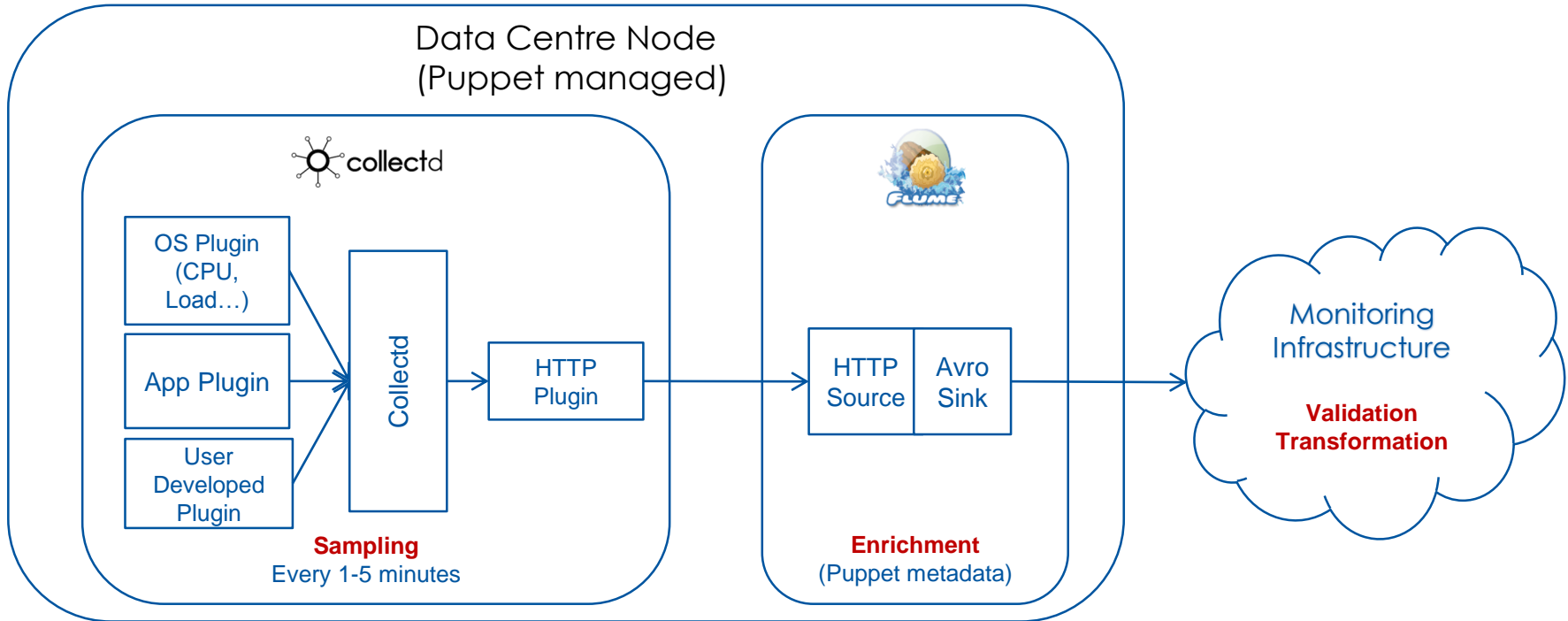
# Our case

Official support available for most of the metrics

Lemon Metric Classes	Lemon Sensors	Collectd Support	Collectd Plugin
file.filecount	file	Official	<a href="https://collectd.org/wiki/index.php/Plugin:FileCount">https://collectd.org/wiki/index.php/Plugin:FileCount</a>
file.size	file	Official	<a href="https://collectd.org/wiki/index.php/Plugin:FileCount">https://collectd.org/wiki/index.php/Plugin:FileCount</a>
file.spaceUsed	file	Official	<a href="https://collectd.org/wiki/index.php/Plugin:FileCount">https://collectd.org/wiki/index.php/Plugin:FileCount</a>
file.sslmtime	file	?	
log.Parse	parseLog	Official	<a href="https://collectd.org/wiki/index.php/Plugin:Tail">https://collectd.org/wiki/index.php/Plugin:Tail</a>
log.Parse	parseLog	Official	<a href="https://collectd.org/wiki/index.php/Plugin:Tail">https://collectd.org/wiki/index.php/Plugin:Tail</a>
cmd.ParseCmd	parse-cmd	Official	<a href="https://collectd.org/wiki/index.php/Plugin:Exec">https://collectd.org/wiki/index.php/Plugin:Exec</a>
system.bootTime	linux	Official	<a href="https://collectd.org/wiki/index.php/Plugin:Uptime">https://collectd.org/wiki/index.php/Plugin:Uptime</a>
system.contextSwitches	linux	Official	<a href="https://collectd.org/wiki/index.php/Plugin:ContextSwitch">https://collectd.org/wiki/index.php/Plugin:ContextSwitch</a>
system.CPUCount	linux	Official	<a href="https://collectd.org/wiki/index.php/Plugin:CPU">https://collectd.org/wiki/index.php/Plugin:CPU</a>
system.CPUInfo	linux	?	
system.CPUUtil	linux	Official	<a href="https://collectd.org/wiki/index.php/Plugin:CPU">https://collectd.org/wiki/index.php/Plugin:CPU</a>
system.CPUUtilization	linux	Official	<a href="https://collectd.org/wiki/index.php/Plugin:CPU">https://collectd.org/wiki/index.php/Plugin:CPU</a>
system.createdProcesses	linux	Official	<a href="https://collectd.org/wiki/index.php/Plugin:Processes">https://collectd.org/wiki/index.php/Plugin:Processes</a>
system.diskStats	linux	Official	<a href="https://collectd.org/wiki/index.php/Plugin:Disk">https://collectd.org/wiki/index.php/Plugin:Disk</a>
system.existingProcesses	linux	Official	<a href="https://collectd.org/wiki/index.php/Plugin:Processes">https://collectd.org/wiki/index.php/Plugin:Processes</a>
system.exitCode	linux	Official	<a href="https://collectd.org/wiki/index.php/Plugin:Exec">https://collectd.org/wiki/index.php/Plugin:Exec</a>
system.fullLoadAvg	linux	Official	<a href="https://collectd.org/wiki/index.php/Plugin:Load">https://collectd.org/wiki/index.php/Plugin:Load</a>
system.interrupts	linux	Official	<a href="https://collectd.org/wiki/index.php/Plugin:IRQ">https://collectd.org/wiki/index.php/Plugin:IRQ</a>
system.loadAvg	linux	Official	<a href="https://collectd.org/wiki/index.php/Plugin:Load">https://collectd.org/wiki/index.php/Plugin:Load</a>
system.meminfo	linux	Official	<a href="https://collectd.org/wiki/index.php/Plugin:Memory">https://collectd.org/wiki/index.php/Plugin:Memory</a>
system.memoryShared	linux	Official	<a href="https://collectd.org/wiki/index.php/Plugin:Memory">https://collectd.org/wiki/index.php/Plugin:Memory</a>
system.memoryStats	linux	Official	<a href="https://collectd.org/wiki/index.php/Plugin:Memory">https://collectd.org/wiki/index.php/Plugin:Memory</a>
system.networkInterfaceDropped	linux	Official	<a href="https://collectd.org/wiki/index.php/Plugin:Interface">https://collectd.org/wiki/index.php/Plugin:Interface</a>
system.networkInterfaceInfo	linux	?	
system.networkInterfaceIO	linux	Official	<a href="https://collectd.org/wiki/index.php/Plugin:Interface">https://collectd.org/wiki/index.php/Plugin:Interface</a>
system.numberOF.Sockets	linux	Official	<a href="https://collectd.org/wiki/index.php/Plugin:TCPConns">https://collectd.org/wiki/index.php/Plugin:TCPConns</a>
system.numberOF.Users	linux	Official	<a href="https://collectd.org/wiki/index.php/Plugin:Users">https://collectd.org/wiki/index.php/Plugin:Users</a>

# Data Sources: collectd

Data Centre Node  
(Puppet managed)



# Transport: Kafka

Kafka as rock-solid core of our pipeline

- Distributed message broker
- Enables stream processing
- 72h retention policy
- Gzip compressed
- Replica factor: 3



# Processing



## Stream processing

### Data enrichment

- Join information from several sources (e.g. network topology)

### Data aggregation

- Over time (e.g. summary statistics for a time bin)
- Over other dimensions (e.g. compute a cumulative metric for a set of machines hosting the same service)

### Data correlation

- Advanced Alarming: detect anomalies and failures correlating data from multiple sources (e.g. data centre topology-aware alarms)

## Batch processing

- Reprocessing, data compression, historical data, periodic reports

# Storage

- HDFS for long-term archive
  - Compressed JSON or Parquet
  - Data kept ~ forever
- ES for data exploration & discovery
  - 2 large instances: 30 nodes
  - Data kept for 1 month
- InfluxDB for time-series dashboards
  - > 20 instances
  - From 8GB up to 128GB memory per instance
  - Data kept for ~5y, down-sampling
- Provided by other IT services



elasticsearch





# Visualization

- Grafana for dashboards
  - Users can create their own
- Kibana for data exploration
  - Data discovery and logs
- Swan for analytics (notebooks)



Grafana



Kibana



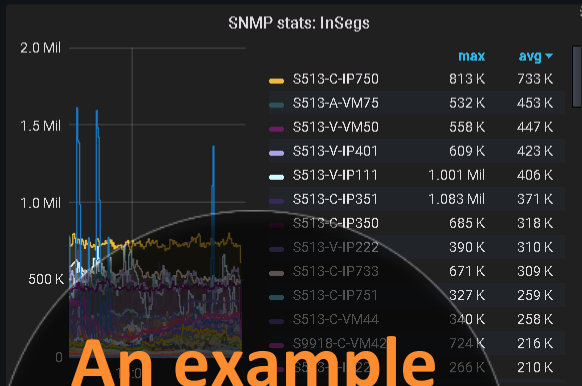
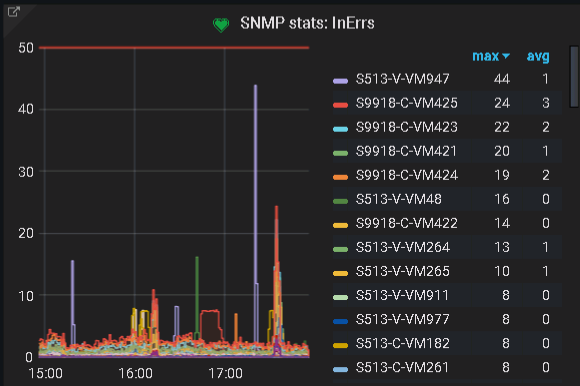
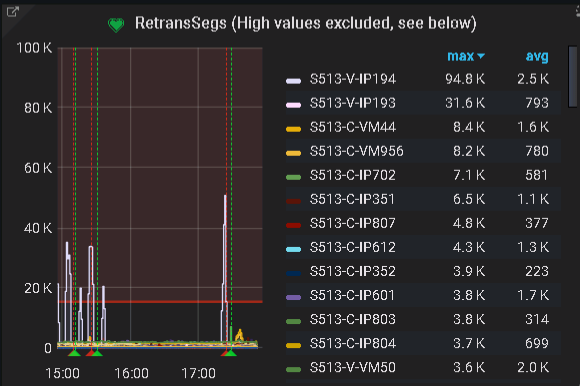
# Visualization: Grafana

- Open-source platform for dashboards
- Support multiple back ends: ES / InfluxDB
- Advanced features
  - Templates, ad-hoc filters, auto-completion
  - Advanced query syntax, organizations, ACLs
  - Alarms

Environment All
Top Hostgroup All
Hostgroup All
Host All
LanDB Service Name All
Filters +

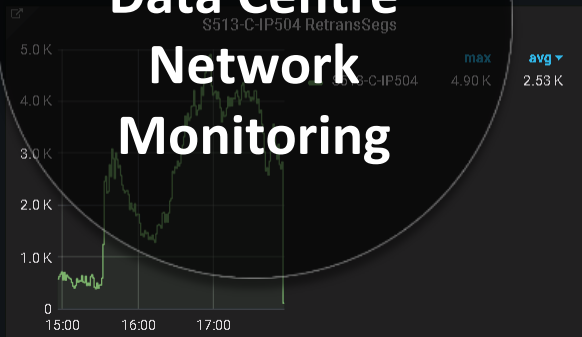
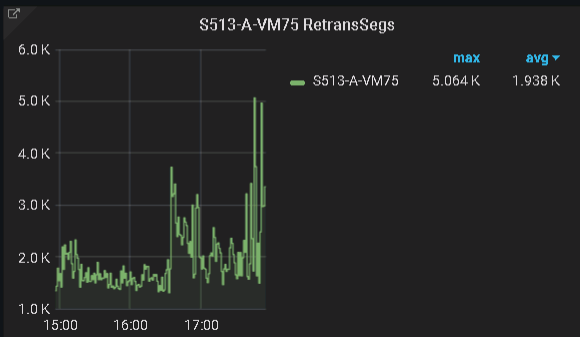
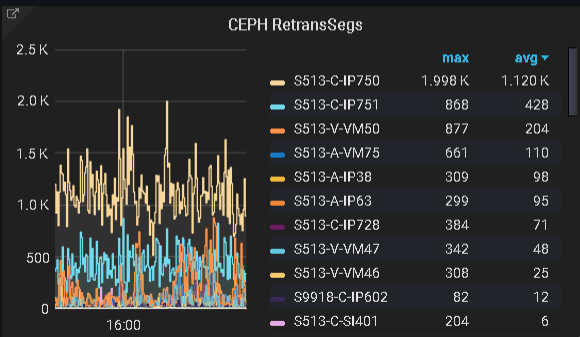
Dashboard Row (1 hidden panels) ⚙️ 🗑️

Dashboard Row



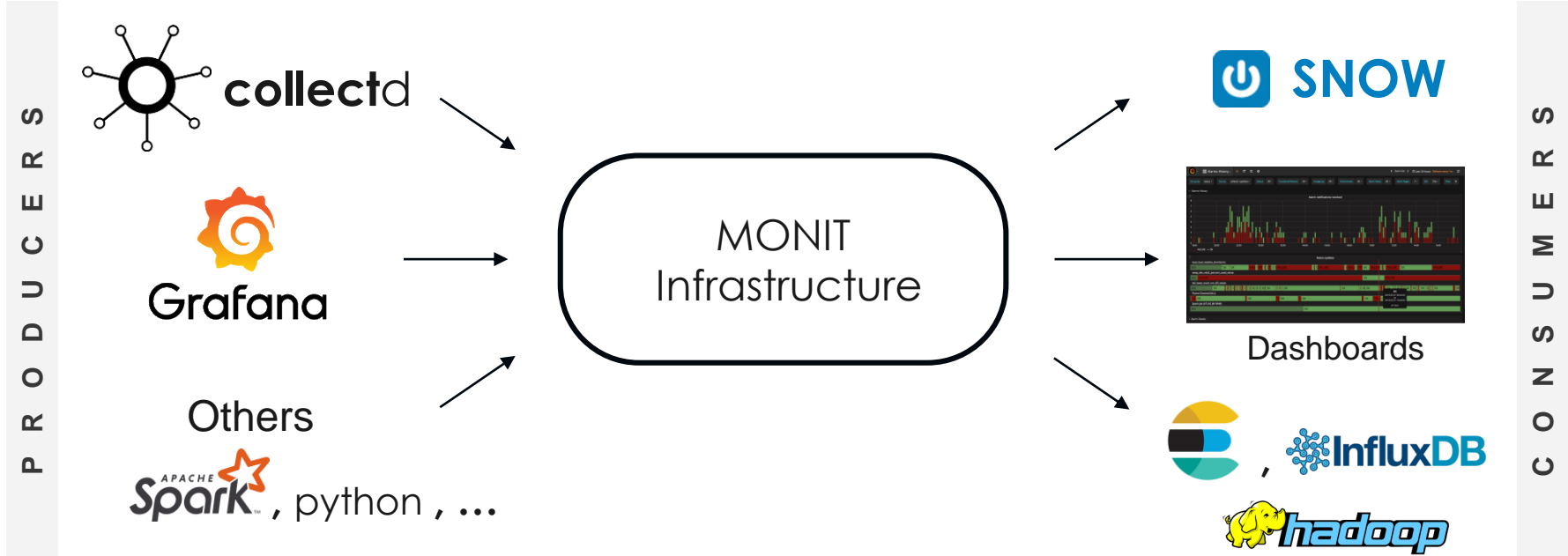
An example  
Data Centre  
Network  
Monitoring

Dashboard Row



Alarms (WIP)

# Alarms: Overview



# Alarms

- Standard representation as any other *MONIT document*
- Integrated notification system
  - All alarms can produce: support tickets, emails, etc.
- Collectd alarms
  - Use collectd local metrics with *threshold* plugin
  - *Low level*, created at the (virtual) machine
- Grafana alarms
  - Defined using the Grafana GUI
  - Integrated *automatically* into the MONIT workflow





Functional Element LX BATCH Alarm Name Regex All Alarm Name All

Live view

*alarm.batch\_vm\_error\_state*  
Status - OK

*df\_pool\_percent\_bytes\_used\_value*  
Status - OK

*exception.batch\_vm\_error\_state*  
Status - FAILURE

*load\_load\_relative\_shortterm*  
Status - OK

*swap\_dev\_vda2\_percent\_used\_value*  
Status - FAILURE

*tail\_base\_count\_filesystem\_error\_value*  
Status - OK

*tail\_base\_count\_vm\_kill\_value*  
Status - OK

> Statistics

+ ADD ROW

**Alarms**  
Live view of  
components  
health

# Operations



# Using MONIT at CERN DC

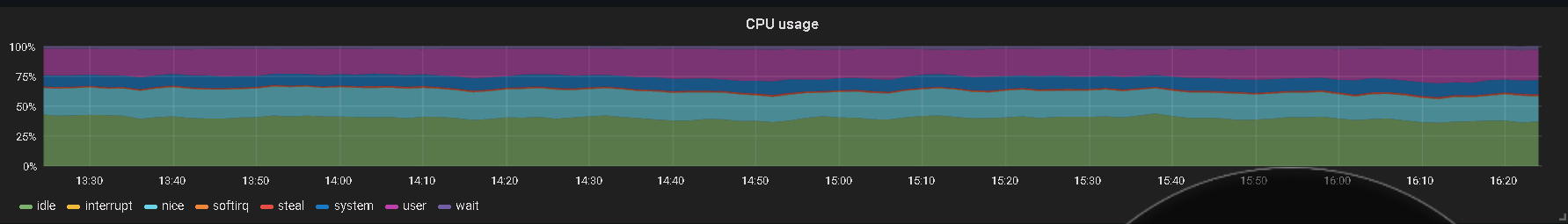
1. Configure the service using the Agile Infrastructure (**Puppet**)
  - Basic monitoring configuration provided by MONIT
2. Configure service specific metrics and logs
  - Configure **collectd plugins** in Puppet
  - Configure local **log collector**
3. Monitor the service
  - General monitoring dashboards available in Grafana
  - Browse logs and discover service-specific data in Kibana
  - Create service-specific dashboards in Grafana
  - More advanced use cases



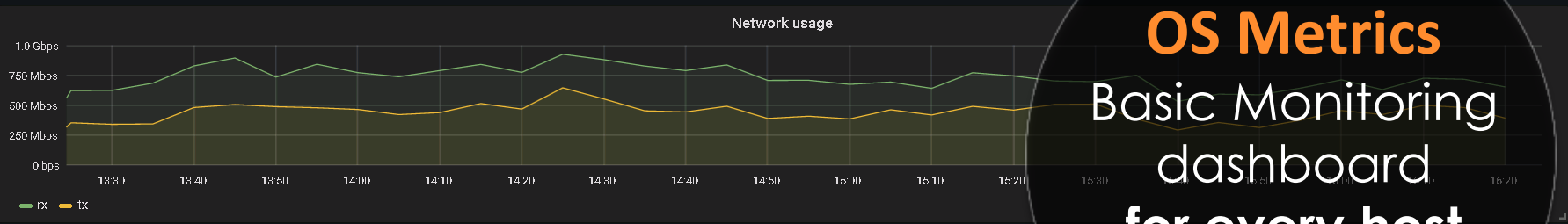


Environment All ▾ Hostgroup lxplus/nodes/login ▾ Host All ▾ Availability Zone All ▾ Data Range one\_week ▾ Bin auto ▾ Filters +

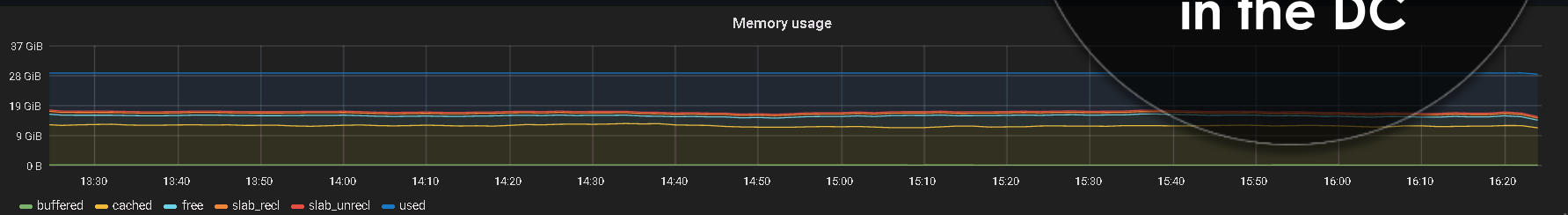
▼ CPU



▼ Network



▼ Memory



**OS Metrics**  
 Basic Monitoring  
 dashboard  
 for every host  
 in the DC



partition batch pool a + b environment All frontend frontend-puppet-batch4 balancers punch/puppet/hap/a + punch/puppet/hap/b

↓ KPIs

0.0004% t-outs

0.3% ISEs

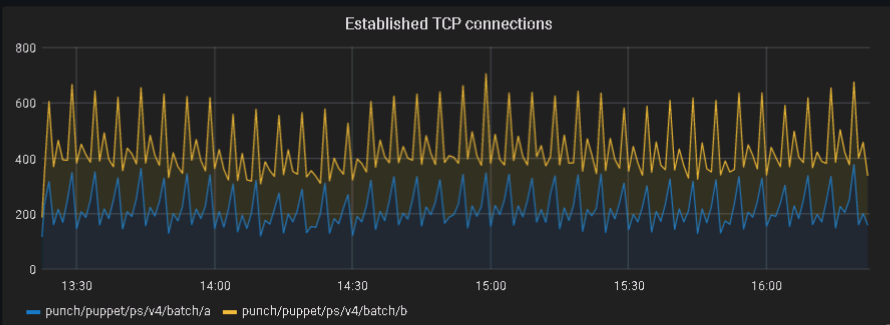
0.053% 404s

99.6619% OKs

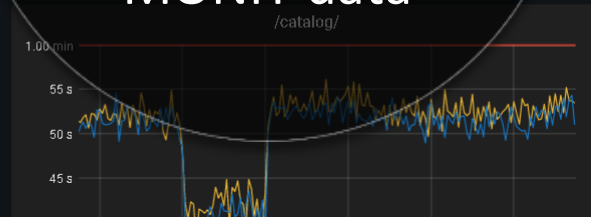
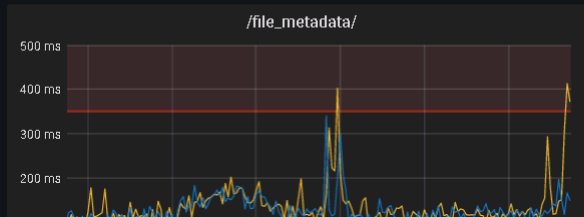
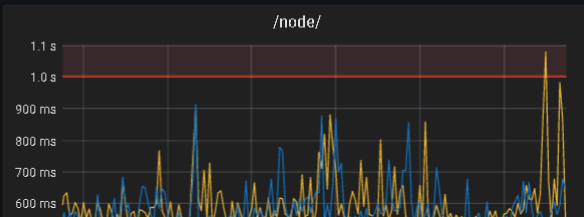
> %CPU (3 hidden panels) ⚙️ 🗑️

> Memory (1 hidden panels) ⚙️ 🗑️

↓ Networking & Load



↓ Latencies of the entrypoints



**IT Services**  
 Create custom  
 dashboards  
 with metrics and  
 MONIT data

# MONIT Operations

- Based on Openstack VMs (CentOS7)
- Puppet for service configuration
- Separate production/QA clusters
- Collectd for infrastructure metrics
- MONIT alarms



openstack.



# Lessons learned

- Migration to open source technologies pays back
- Building a reliable pipeline
  - Decouple producers from back ends
  - Producers isolation (1 producer, 1 topic)
  - Rely on technology that scales (Kafka over Flume)

# Lessons learned

- Component based approach over monolithic solutions
- Get to know your storage systems
  - InfluxDB
    - Works well with time-series data
    - SQL-like language with built-in aggregations
    - Requires care with high cardinality in number of series
      - Needs some schema design
  - Elasticsearch
    - Good for discovery and search of logs
    - Complex engine that requires special care with configuration
  - HDFS
    - Only for batch or high latency scenarios

# Summary

- Production infrastructure in place
  - Including metrics, logs and alarms
- Based on mainstream and open source technologies
- We provide monitoring as a service

# Thank you

- Our docs:  
<http://monit-docs.web.cern.ch/monit-docs/>
- Questions / comments are welcome!



