# CERN Central Monitoring Overview

Borja Garrido Bear (On behalf of the monitoring team)

12.09.2024

# INTRODUCTION

- **The MONIT Infrastructure**
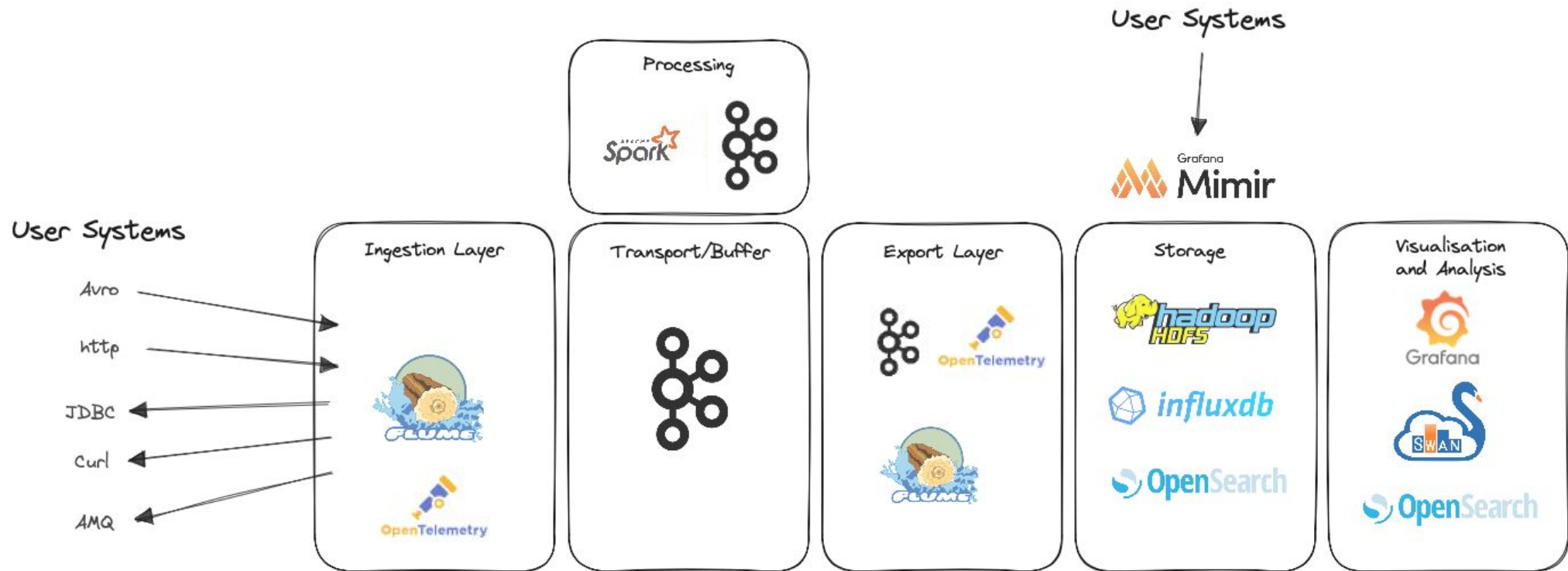- **Some users**

THE THREE USE CASES

FTS

DDM

XROOTD

SUMMARY

# The MONIT infrastructure

**The aim of the monitoring team is to provide an infrastructure allowing the usage of tools and service to ease the load of monitoring CERN data centres (host and services) and WLCG experiments**

# Some users

# INTRODUCTION

# THE THREE USE CASES

- **Same set of tools**
- **Data processing**
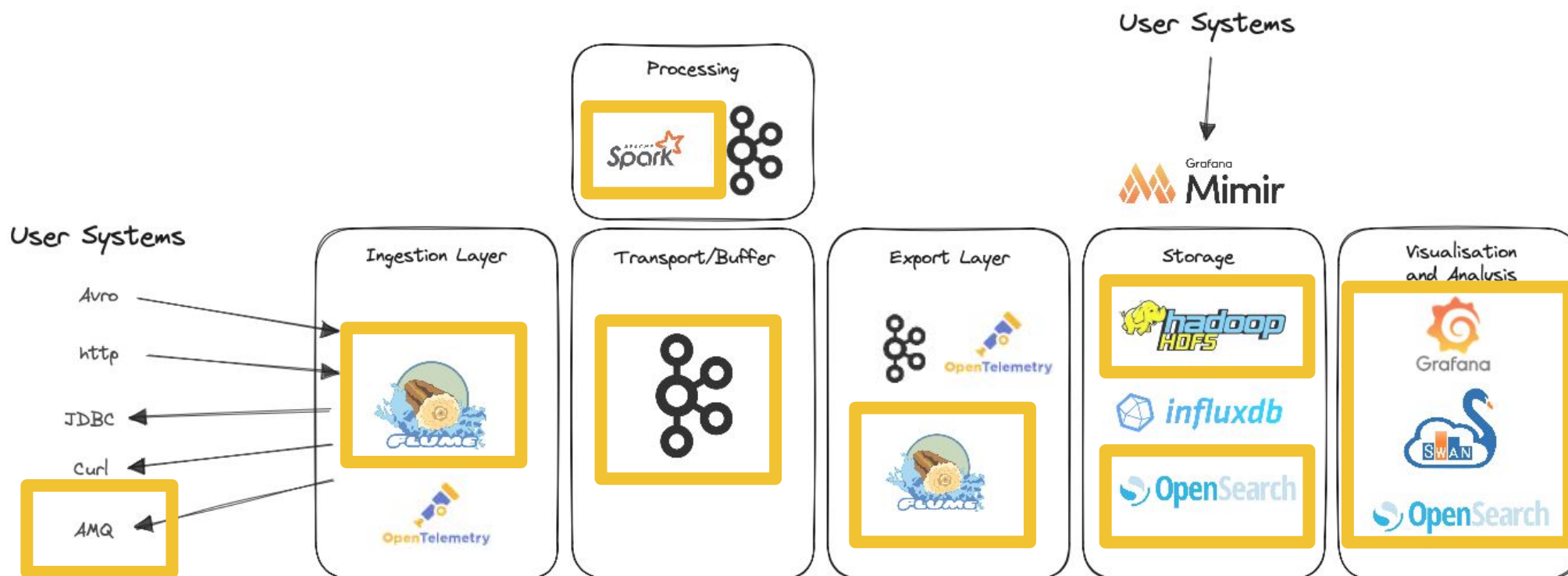- **CRIC topology document**

# FTS

# DDM

# XROOTD

# SUMMARY

# Same set of tools

- **FTS, DDM and XRootD flows share a similar processing within MONIT**
  - Ingested via Message broker (AMQ) in the shape of JSON
  - Processed with Apache Spark (enrichment, aggregation)
  - Stored in Opensearch (short term and long term) and HDFS
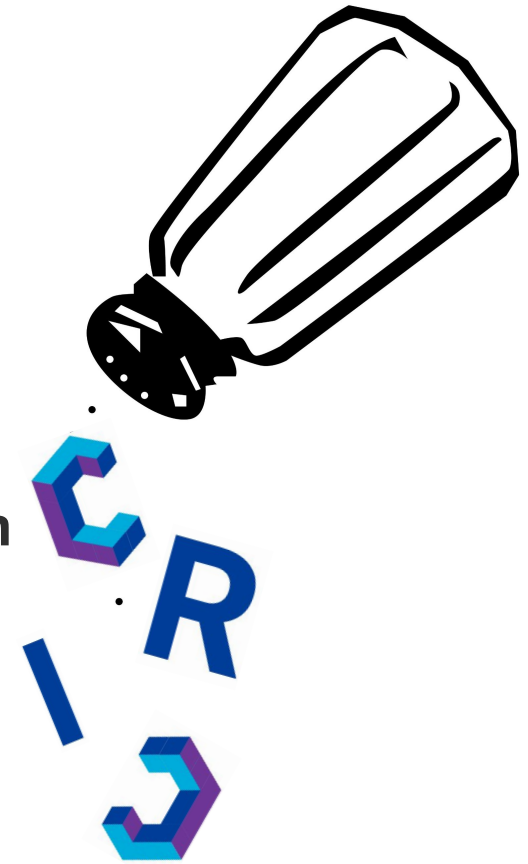  - Accessible from Grafana, Opensearch or SWAN

# Data processing

- **Data is processed with two main goals behind**

  - Real time monitoring, requires stream enrichment, stored for short-term (few months)

  - Accounting, requires batch enrichment and aggregation, stored forever

- **Enrichment process**

  - Extend the base document, in general with topology information (CRIC)

    - I.e: site, country, federation…

- **Aggregation process**

  - Create buckets based on time and a selected set of labels, usually 1h

# CRIC Topology Document

- **MONIT does a dump of CRIC endpoints every 6 hours**

- **Enrichment jobs get changes on MONIT CRIC flow every 24 hours**

- **Aggregation jobs get changes on MONIT CRIC flow every execution**

- **The CRIC document contains the following fields**

  - CRIC fields used for general enrichment

| VO | experiment_site | official_site | federation | state |
|---|---|---|---|---|
| country | country_code | description | tier | institute_name |
| netroutes | endpoint | hostname | flavour | status |

# Types of data

| | |
|---|---|
| **TRANSFER START**<br>(monit_prod_fts_raw_start)<br><br><br>Created only by the transfer agent at start of transfer | **TRANSFER COMPLETE**<br>(monit_prod_fts_raw_complete)<br><br><br>Created only by the transfer agent at end of transfer |
| **TRANSFER STATE**<br>(monit_prod_fts_raw_state)<br><br><br>Created by QoS and Transfer daemons when a file state changes | **OPTIMIZER**<br>(monit_prod_fts_raw_queue_state)<br><br><br>Created by the Optimizer when adjusting parameters |

- **Information taken from this presentation by the FTS team**

# Transfer Complete transformations (I)

- **Initial round of derivation (extract fields from other fields)**

  - Field added by MONIT, Field provided by FTS

| | | | |
|---|---|---|---|
| `src/dst_se <- src/dst_url` | `protocol <- src/dst_url` | Remove paths from `src/dst_hostname` | `job_id <- tr_id` |
| `file_id <- tr_id` | `log_link <-`<br>`https://{endpnt}:8449/fts3/ftsmon/#/job/{job_id}` | `file_size <- f_size` | `latency <-`<br>`now() - timestamp_tr_comp` |
| `operation_time <-`<br>`(tr_timestamp_complete -`<br>`tr_timestamp_start)` | `transfer_time <-`<br>`(timestamp_tr_comp -`<br>`timestamp_tr_st)` | `throughput <-`<br>`f_size / (transfer_time/1000)` | `srm_finalization_time <-`<br>`(time_srm_fin_end -`<br>`time_srm_fin_st)` |
| `srm_preparation_time <-`<br>`(time_srm_prep_end -`<br>`time_srm_prep_st)` | `srm_overhead_time <-`<br>`(srm_preparation_time -`<br>`srm_finalization_time)` | `srm_overhead_percentage` | `timestamp_checksum__diff<-`<br>`timestamp_checksum__diff -`<br>`timestamp_checksum__st` |
| `t_final_transfer_state_flag`<br>`<-`<br>`t_final_transfer_state` | `activity <-`<br>`file_metadata[activity]` | `dst/src_rse <-`<br>`file_metadata[dst/src_rse]` | `_id <-`<br>`SHA1(job_id+retry+tr_id)` |

Complete format in <u>FTS Messaging documentation</u>

# Transfer Complete transformations (II)

- **Enrichment with different sources**

    - This ends up producing so called FTS Complete Enriched dataset (monit_prod_fts_enr_complete)

| CRIC | Transfer State | Final derivations |
|------|----------------|-------------------|
| <ul><li>Match by:<ul><li>VO</li><li>Endpoint</li></ul></li><li>Added fields (src/dst):<ul><li>Hostname</li><li>Experiment_site</li><li>Site</li><li>Tier</li><li>Country</li><li>Federation</li></ul></li></ul> | <ul><li>Match by:<ul><li>Timestamp</li><li>File_ID</li></ul></li><li>Added fields:<ul><li>Staging</li><li>Staging_start</li><li>Staging_finished</li><li>User_filesize</li></ul></li></ul> | <ul><li>Remote_access:<ul><li>Src_site != Dst_site</li></ul></li><li>Transferred_volume:<ul><li>File_size if complete</li><li>Tr_bt_transferred else</li></ul></li></ul> |

# Transfer Complete transformations (III)

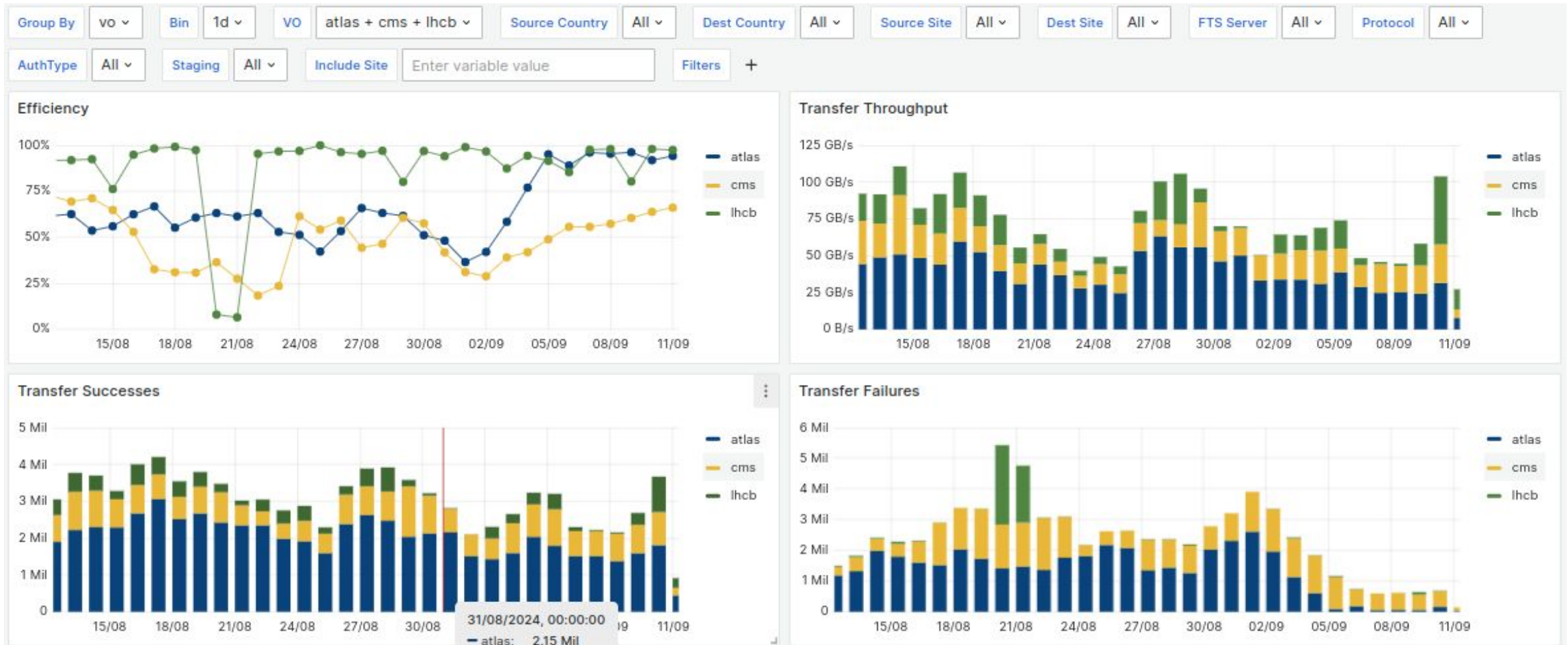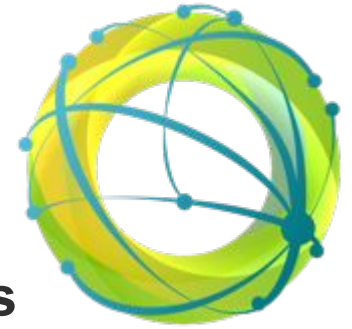- **Aggregation for accounting**

  - This ends up producing so called FTS Complete Aggregated dataset (monit_prod_fts_agg_complete)

  - Aggregates data by labels every hour and corrects produced data up to 2 days in the past using the labels as primary key, uses the tr_timestamp_complete field for the time window

| activity | channel_type | country | experiment_site | federation |
|---|---|---|---|---|
| hostname | se | site | srm_v | tier |
| endpnt | ipv6 | vo | is_recoverable | job_state |
| remote_access | srm_space_token | staging | t_failure_phase | technology |
| t_final_transfer_state | tr_error_category | tr_error_scope | rse | protocol |
| ipver | auth_method | transfer_type | count | avg_file_size |
| transferred_volume | avg\|sum operation_time | avg\|sum overhead_time | sum_transfer_state | sum_staging_duration |
| stage_transferred_volume | | | sum_transfer_time | |

# Dashboards

- **Initially accounting dashboard was the main goal for [WLCG Monitoring](#)**

- **Data was there, so people started using it to build their own dashboards**

INTRODUCTION

THE THREE USE CASES

FTS

**DDM**
- **Types of data**
- **Transformations**
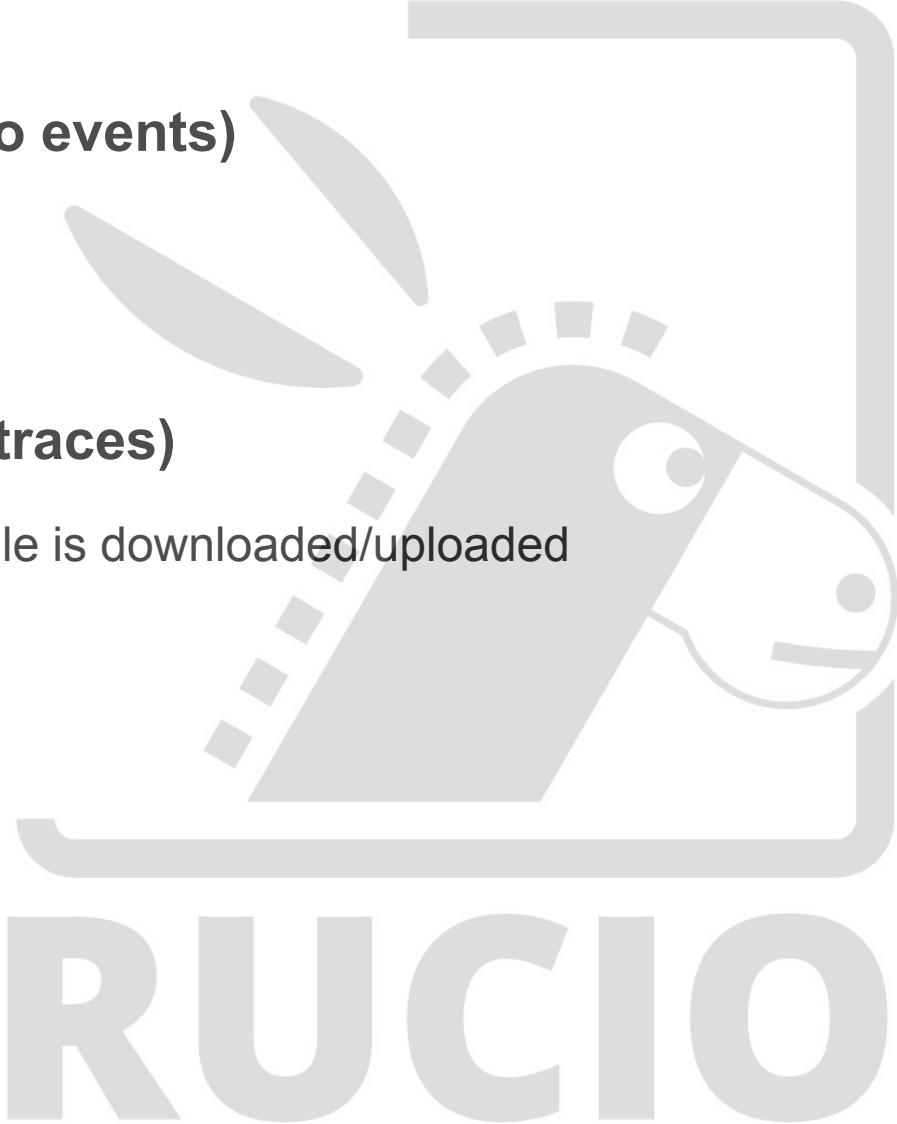- **Dashboards**

XROOTD

SUMMARY

# Types of data

- **Transfer Monitoring (Rucio events)**

  - Transfers events

  - Deletion events

- **Access Monitoring (Rucio traces)**

  - Traces sent by pilot when a file is downloaded/uploaded

# Transformations (I)

| EVENTS | TRACES |
|---|---|
| Rename fields replacing "-" by "_" | |
| Promote payload to JSON root | |
| Generate new "tf_" unix times from: created_at, submitted_at, started_at, transferred_at | |
| Generate purged_reason by regex replace of reason or stateReason field | |
| Sets VO based on producer | |
| Adds "https:" to transfer_link if needed | |
| Sets event_timestamp as tf_created_at | Sets event_timestamp as traceTimeentryUnix |

# Transformations (II)

- **Enrichment with different sources (Only for ATLAS)**

  - This ends up producing so called DDM transfer Enriched dataset (monit_prod_ddm_enr_transfer)

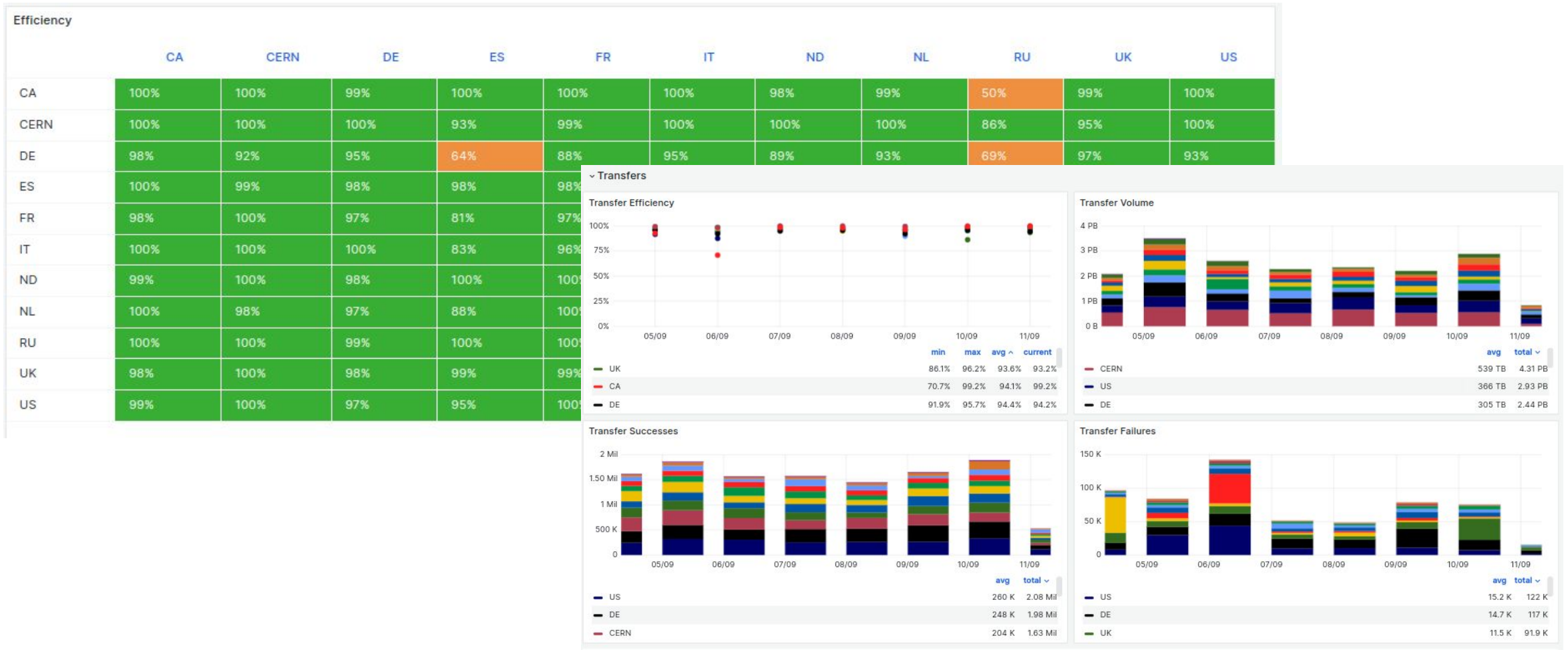| CRIC (ATLAS) | CRIC (WLCG) | Final derivations |
|---|---|---|
| <ul><li>Match by:<ul><li>Endpoint</li></ul></li><li>Added fields (src/dst):<ul><li>Experiment_site</li><li>Token</li><li>Cloud</li></ul></li></ul> | <ul><li>Match by:<ul><li>Experiment_site</li></ul></li><li>Added fields (src/dst):<ul><li>Tier</li><li>Country</li><li>Federation</li></ul></li></ul> | <ul><li>Remote_access:<ul><li>Src_site != Dst_site</li></ul></li></ul> |

# Transformations (III)

- **Aggregation for accounting**

  - This ends up producing so called DDM Aggregated dataset (monit_prod_ddm_agg_transfer)

  - Aggregates data by labels every hour and corrects produced data up to 1 day in the past using the labels as primary key, uses the event_timestamp field for the time window

| state | activity | protocol | endpoint | is_staging |
|---|---|---|---|---|
| fts_host | vo | site | tier | federation |
| country | event_type | remote_access | experiment_site | token |
| cloud | | | | |
| bytes_done | bytes_failed | bytes_planned | bytes_total | |
| files_done | files_failed | files_planned | files_total | |
| throughput_done | throughput_failed | throughput_planned | throughput_total | |

# Dashboards

- **Done in collaboration with ATLAS for their monitoring**

# New Flow

# New Flow (New components)

- **Shoveler**
  - Receives streams (UDP) and persists them to a message queue

    - It also does stream validation and IP translation (configurable map)

  - Should sit close to the XRootD server (run at sites)

    - Avoid UDP fragmentation and general unreliability

- **Collector**

  - One single collector per Message Queue (run centrally)

  - Aggregates streams into a transfer document

    - Needs to keep state due to how XRootD produces streams

  - Does some initial massaging of the transfer document

    - Extracts VOs, resolves hostnames, domains…

# Transformations (I)

- **Much simpler that in the previous cases**
  - Mainly because the collector does most of the required derivations before sending the data

```
remote_access <-
client_domain != server_domain
```

```
is_transfer <-
bytes_read + bytes_write == file_size
```

```
event_timestamp <- end_time
```

# Transformations (II)

- **Join together WLCG and US transfers (currently integrated in different flows)**

- **Enrichment with different sources**

  - This ends up producing so called XRootD transfer enriched dataset (monit_prod_xrootd_enr_transfer)

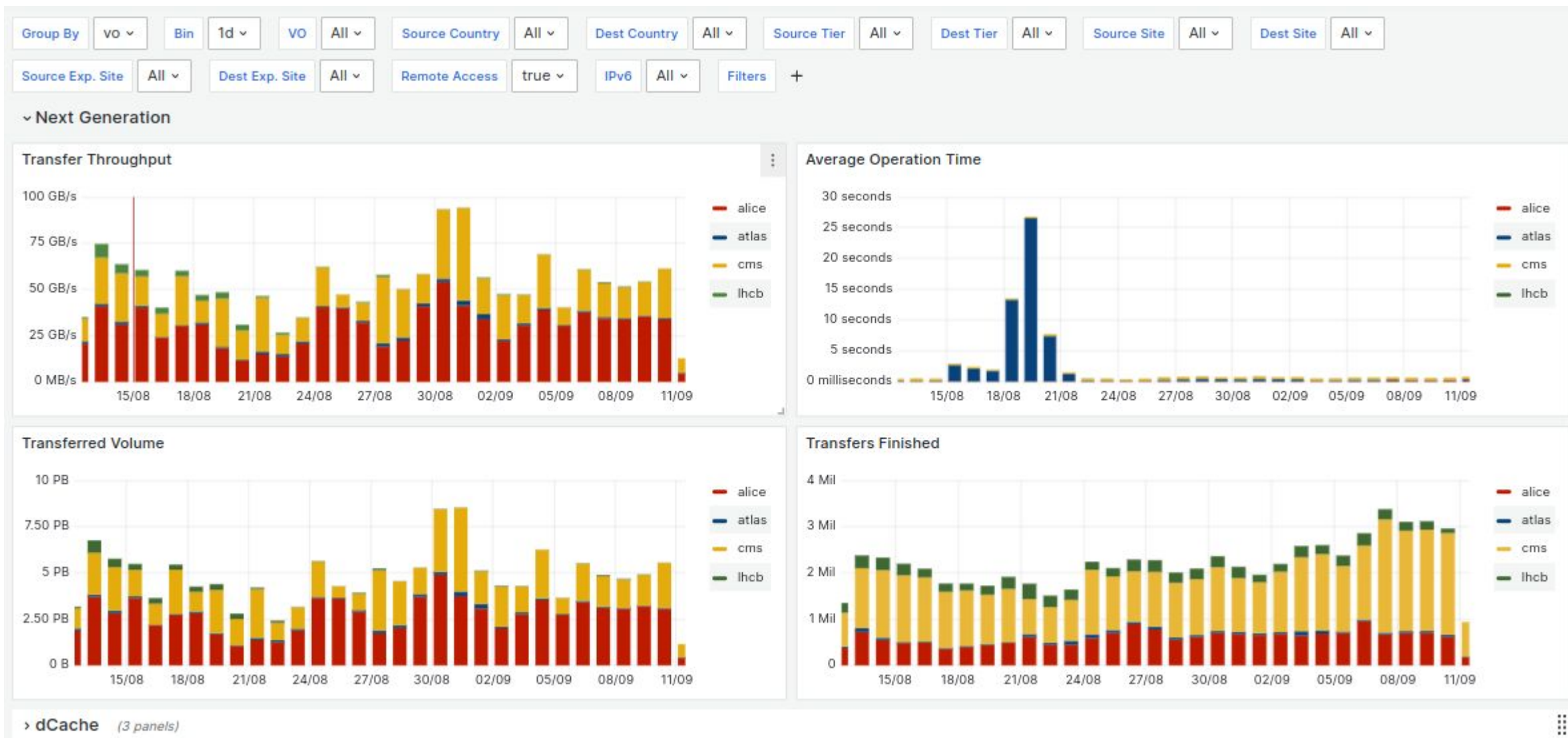| Static Topology | CRIC (WLCG) |
|---|---|
| <ul><li>Match by:<ul><li>Domain</li><li>VO</li></ul></li><li>Added fields (src/dst):<ul><li>Site</li><li>Experiment_site</li><li>Country</li></ul></li></ul> | <ul><li>Match by:<ul><li>Experiment_site</li></ul></li><li>Added fields (src/dst):<ul><li>Tier</li><li>Federation</li></ul></li></ul> |

# Transformations (III)

- **Aggregation for accounting**

  - This ends up producing so called XRootD Aggregated dataset (monit_prod_xrootdng_agg_transfer)

  - Aggregates data by labels every hour and corrects produced data up to 6 hours in the past using the labels as primary key, uses the event_timestamp field for the time window

| vo | site | ipv6 | is_transfer | operation |
|---|---|---|---|---|
| server_site | experiment_site | tier | country | federation |
| hostname* | technology | remote_access | | |
| count | avg_file_size | sum_file_size | transferred_volume | |
| avg_operation_time | sum_operation_time | | | |

# Dashboards

INTRODUCTION

THE THREE USE CASES

FTS

DDM

XROOTD

**SUMMARY**

# Summary

- **There's more to the eye from just "simple" integration into MONIT**
  - FTS and DDM flows have been curated with time (10+ years)
  - New fields are added from other datasources as requested
- **Split knowledge of the monitoring (not ideal)**
  - Tools experts know what they produce
  - MONIT experts know all the extra bits happening under stage
- **Be careful when putting different flows against each other**
  - Something like the timestamp concept difference might drive to very different plots!
- **WLCG Goal remains the same**
  - Being able to plot together WLCG transfers data (FTS + XRootD)
- **XRootD new flow still in early days**
  - Will require lots of work to arrive to a "stable" state comparable with FTS or DDM

# Thank you !

# Q & A

SNOW: Monitoring Service
Mattermost: MONIT
Docs: https://monit-docs.web.cern.ch/

home.cern