

ATLAS EventIndex in Hadoop

Hadoop User Forum

4 December 2023

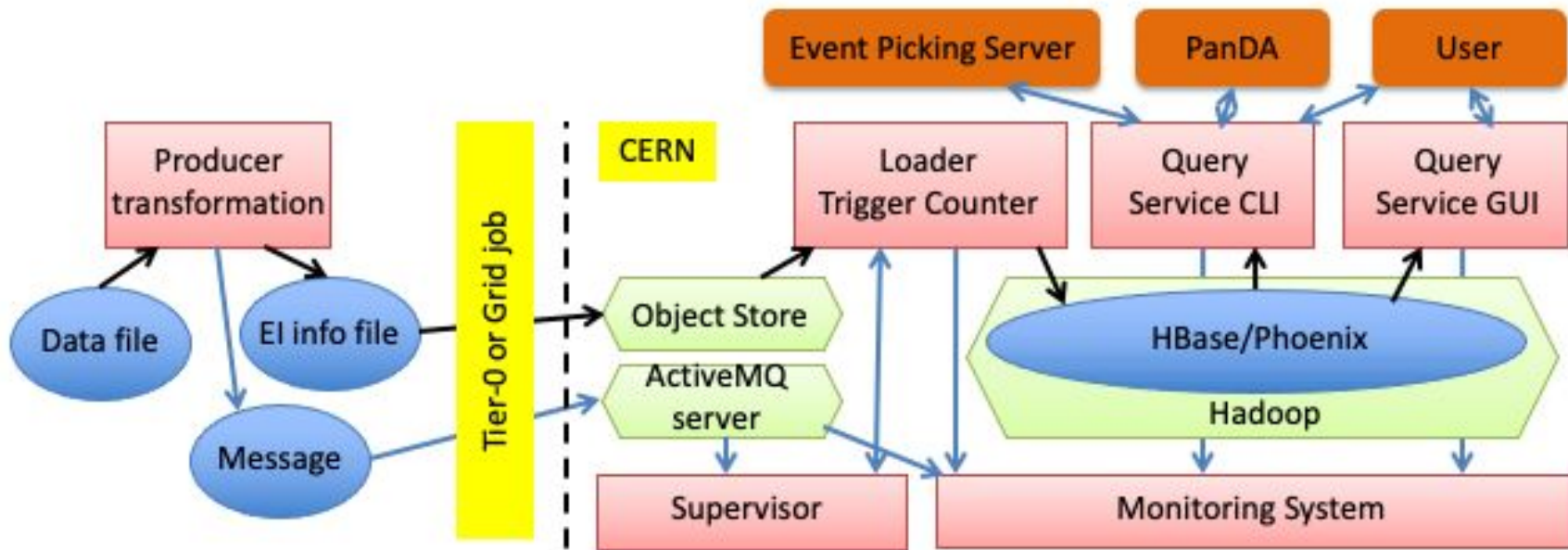
Grigori Rybkine

IJCLab Orsay (France)

on behalf of the ATLAS EventIndex team



- The EventIndex is the global catalogue of all ATLAS events
- For each event, each data format and each processing version, it contains:
 - Event identifiers (run and event number, plus other stuff)
 - Location (GUID of the file containing it) and provenance
 - Trigger and other useful metadata
- Main use case is event picking for detailed analysis and/or displays
 - Also production checks and dataset overlap counts
- The first EventIndex version was designed in 2013, implemented in 2014 and started operations in 2015
 - Software tools evolved considerably since then!
- The partitioned architecture allowed the replacement of individual components during the years



- Grid jobs extract EventIndex data from each data file and send it to CERN
- Data is aggregated per dataset and stored in HBase tables:
 - “Datasets” table with dataset-level information
 - “Events” table with event-level information
 - Auxiliary tables with other information useful to retrieve and decode event records



- The first implementation in 2014 stored events in HDFS MapFiles, one per dataset
 - HBase was used for the dataset catalogue, and later also to store event records with a subset of the information (no trigger)
 - For read performance reasons, a copy of a subset of events (real data only) with a subset of information (no Trigger) was made to Oracle too
- The core data storage system was reimplemented during 2021 and deployed in 2022 for the start of LHC Run 3
 - HBase for the “datasets” and “events” tables
 - Phoenix interface for SQL queries
 - New client query service CLI implemented for optimal performance
- Old data (Run 1 and Run 2) were transferred to the new structures and new Run 3 data written directly to HBase
 - Simulated data as well!



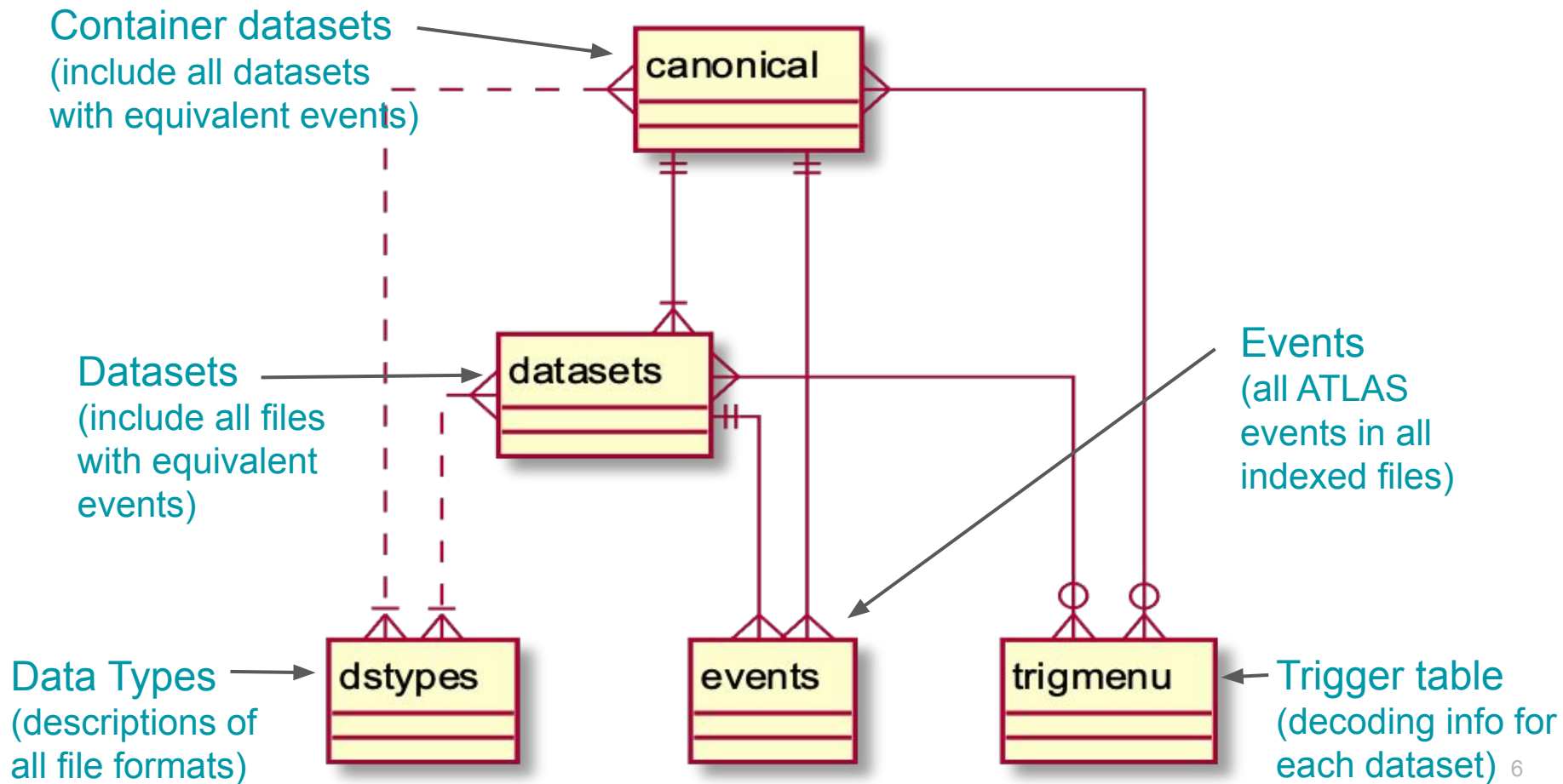
events
<ul style="list-style-type: none">● dspid : integer● dstypeid : smallint● eventno : bigint● seq : smallint
a.tid : integer a.sr : binary(32) a.mcc : integer a.mcw : float
b.pv : binary(34) ARRAY[]
c.lb : integer c.bcid : integer c.lpsk : integer c.etime : timestamp c.id : integer c.tbp : smallint[] c.tap : smallint[] c.tav : smallint[]
d.lb1 : integer d.bcid1 : integer d.hpsk : integer d.lph : smallint[] d.ph : smallint[]

datasets
<ul style="list-style-type: none">● runno : integer● project : varchar(200)● datatype : varchar(200)● streamname : varchar(200)● prodstep : varchar(200)● version : varchar(200)● tid : integer● dspid : integer● dstypeid : smallint
smk : integer events_rucio : bigint files : integer events : bigint events_uniq : bigint events_dup : bigint files_dup : integer rank : smallint status : varchar(200) rucio_at : timestamp updated_at : timestamp dups_at : timestamp trigger_at : timestamp is_open : boolean is_deleted : boolean has_raw : boolean has_trigger : boolean prov_seen : smallint ARRAY[] sr_cnt : varchar(200) sr_clid : varchar(200) sr_tech : varchar(200) name : varchar(250)

canonical
<ul style="list-style-type: none">● runno : integer● project : varchar(200)● datatype : varchar(200)● streamname : varchar(200)● prodstep : varchar(200)● version : varchar(200)● dspid : integer● dstypeid : smallint
smk : integer events_rucio : bigint files : integer events : bigint events_uniq : bigint events_dup : bigint files_dup : bigint rank : smallint status : varchar(200) rucio_at : timestamp updated_at : timestamp dups_at : timestamp trigger_at : timestamp is_open : boolean is_deleted : boolean has_raw : boolean has_trigger : boolean prov_seen : smallint ARRAY[] is_rucio_eq : boolean sr_cnt : varchar(200) sr_clid : varchar(200) sr_tech : varchar(200) name : varchar(250)

dstypes
<ul style="list-style-type: none">● type : tinyint● name : varchar(20)
id : smallint

trigmenu
<ul style="list-style-type: none">● smk : integer● type : tinyint● name : varchar(200)
id : smallint





- Server-client architecture
- Server in Java, making use of Spring Framework
- Database - Phoenix/HBase tables - access via Phoenix JDBC driver
- Provides required performance in terms of both time and number of events queried
- Achieved by heavy use of INNER JOIN in SQL queries
 - `SELECT ... FROM EVENTS INNER JOIN (SELECT ... FROM CANONICAL WHERE ...)`
`ON ...`
- lhs INNER JOIN rhs
 - rhs will be built as hash table in server cache



- With default Phoenix configuration, had server-side HashJoinCacheNotFoundException
- Had to tune Phoenix configuration, in particular, on the server side
 - phoenix.coprocessor.maxServerCacheTimeToLiveMs
360000
Default
30000



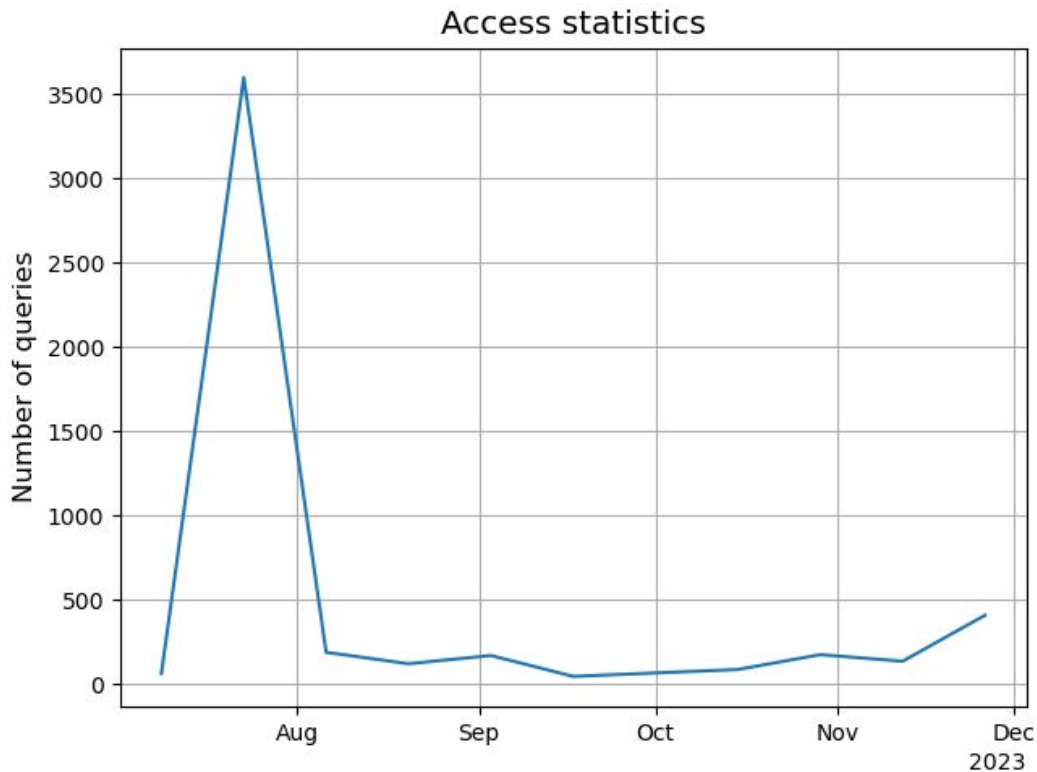
- Putting in production required additional work
 - normally unnecessary since done by Maven/Gradle tools
- Phoenix release deployed on the Hadoop cluster and clients
 - **custom** build of Phoenix-5.0 against **HBase-2.2**
 - artifacts not distributed standard Java way to <https://mvnrepository.com> or similar repository



- (Container) Datasets: 280k
 - Run 1 data: 2300 only real data produced at Tier-0 in AOD format
 - Run 2 data: 88400 includes reprocessing and many derivation formats (DAODs)
 - Run 3 data: 6600 reduced DAOD formats and no repro yet for 2023 data
 - Run 2 MC: 170k many small datasets with variants of event generation
 - Run 3 MC: 12k analysis of 2023 data just starting as stats for real data is low
- Events table entries: 850 billion
 - Largest MC dataset has 2 billion events
 - Largest real data dataset has 1.75 billion events (minimum bias triggers)
- Data volume in HBase: 78 TB (234 TB after 3x replication)



Data query rates the last few months

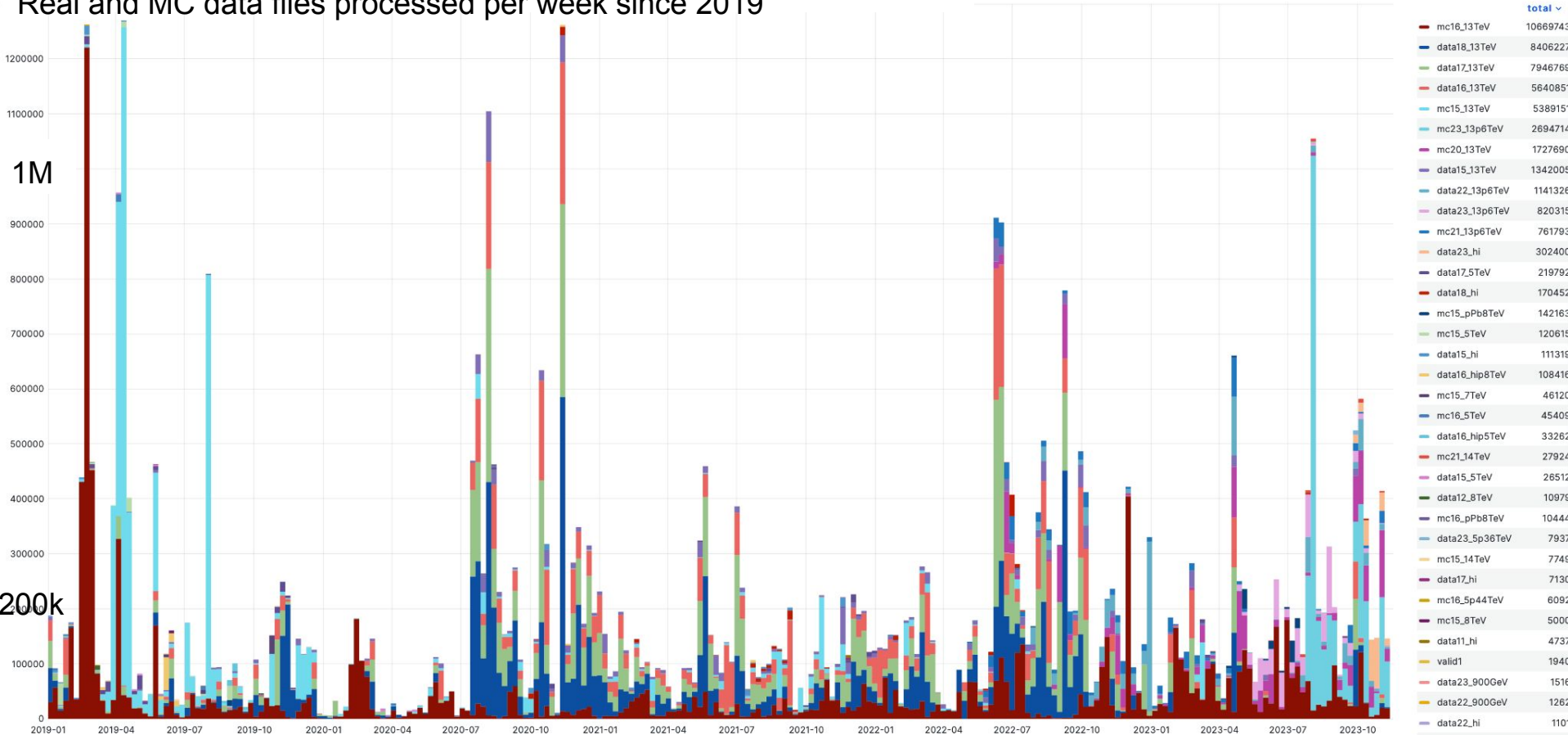




EventIndex in HBase: ingestion rates



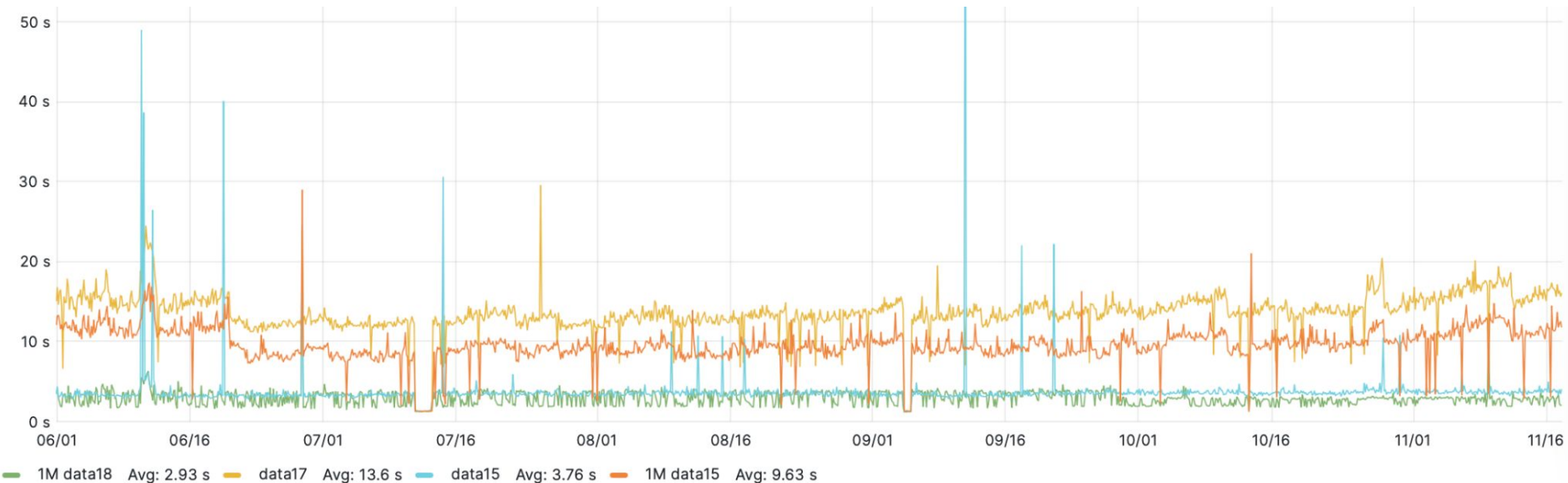
Real and MC data files processed per week since 2019





Event lookup response times for queries of 1000 events from datasets of different sizes between June and November 2023.

Queries generally succeed but there are occasional hiccups. Retries always work.





- The system seems to run stably and we don't plan any major change till the end of Run 3 (end of 2025)
 - The priority for the coming months is to adapt all our scripts to run on AL9 instead of CC7
 - We are short of technical personpower and the situation won't improve any time soon
- Some ideas start being discussed targeting Run 4 (x3 data and MC rates)
 - Tests done with exporting and querying data in Parquet (and the Delta format) look promising, but still partial (notably no ingestion tests)
 - Considerable work would be needed for a complete system



Thank you!