# Experience deploying xCache for CMS in Spain

Carlos Perez Dengra (CIEMAT-PIC), Josep Flix Molina (CIEMAT-PIC) and Anna Sikora (UAB) on behalf of the CMS Collaboration.

*30 of March of 2023 (XRootd+FTS Workshop) @ Ljubljana*

# Context

- **CMS jobs** have the **capability to read data remotely** using the **CMS XRootD federation** (overflow to close sites, files opened in fallback and so on).

- We have been **exploring the xCache** service at PIC Tier-1 and CIEMAT Tier-2 centers to cache data which is read from remote centers

- xCache helps **reducing data access latency**, **improving CPU efficiency** and potentially **reducing the storage** deployed in the region

- We have deployed xCache services at both sites, and dedicated **studies and performance measurements** have been performed to configure the service
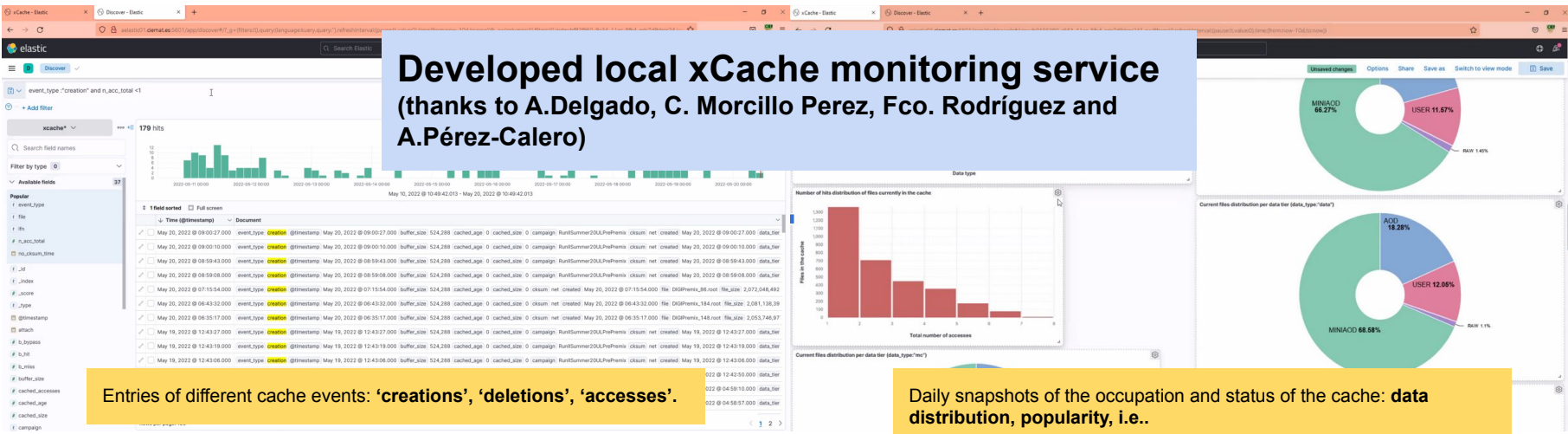
# Commissioning the xCache (initial deployment)

- **Initially, two XCache services** were deployed at both PIC and CIEMAT to understand the service optimal configuration.
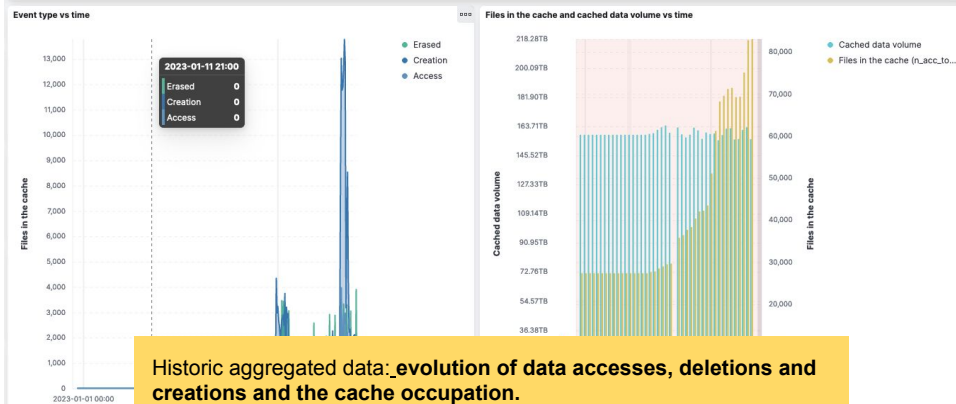
| Deployment | Storage | XrootD version (OSG) | CMS data tiers to cache | Accessibility |
|---|---|---|---|---|
| Initial (2022-Feb 2023) | T1_ES_PIC: 150 TB (mult.disk) T2_ES_CIEMAT: 22 TB (RAID6) | 5.4.2 5.4.2 | **All** **All** | All compute nodes All compute nodes |

- **Caching all CMS data tiers Trivial File Catalog** configured to cache all data. **Inefficient**, but useful to **test the service at scale → not deployed this way anymore.**

- CMS pile-up data produced **high loads in the cache and memory misuse that resulted in high jobs failure rates → reported to developers**

- **Local monitoring** system deployed based, parsing .cinfo files (next slide)

**Developed local xCache monitoring service**
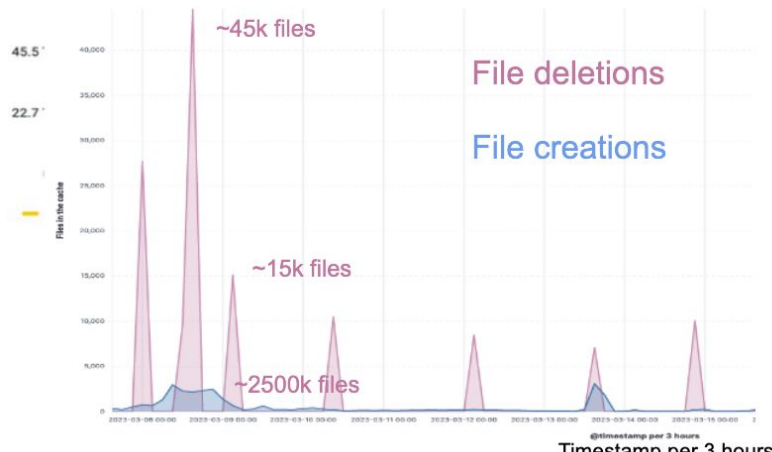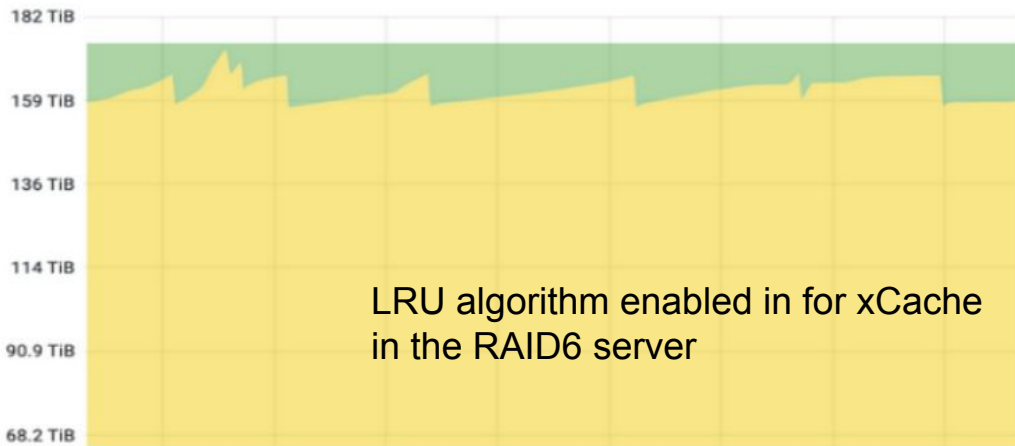**(thanks to A.Delgado, C. Morcillo Perez, Fco. Rodríguez and A.Pérez-Calero)**

Entries of different cache events: **'creations', 'deletions', 'accesses'.**

Daily snapshots of the occupation and status of the cache: **data distribution, popularity, i.e..**

Historic aggregated data: **evolution of data accesses, deletions and creations and the cache occupation.**

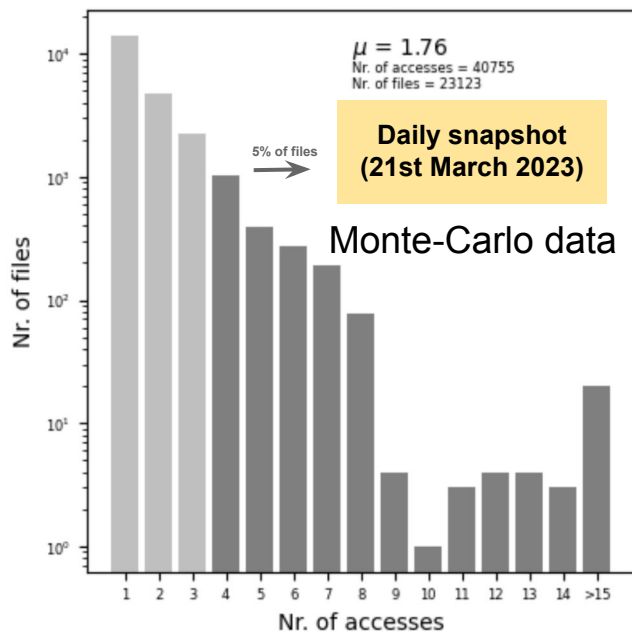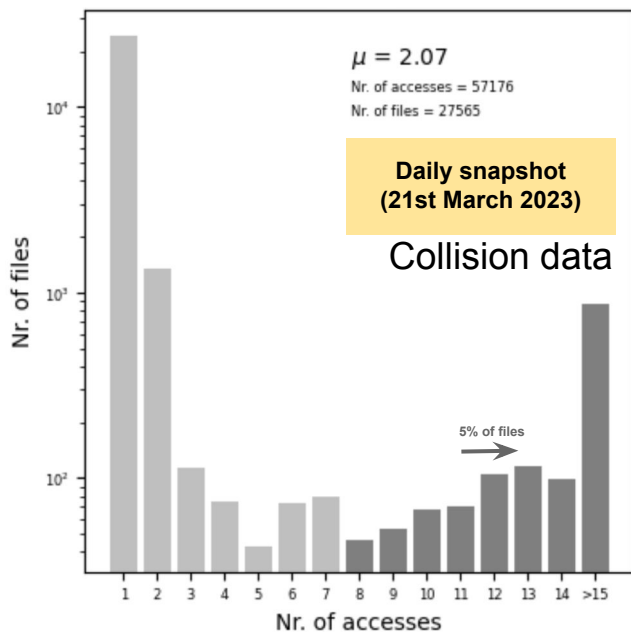Other metrics: **hit-rate computation, file lifetimes, i.e..**

# XCache deployed in PIC (current deployment)

Currently, a **single cache at PIC (180TB)** is serving data for both **PIC and CIEMAT** compute nodes → executed tasks from both sites WN's **reading remote input data are effectively running.**

| Deployment | Storage | XrootD version (OSG) | CMS data tiers to cache | Accessibility |
|---|---|---|---|---|
| Current (single cache for both sites) | T1_ES_PIC: 180 TB (RAID) | 5.5.1 | **All (except PREMIX, "*PrePremix*" and unmerged)** | All compute nodes + half of the compute nodes at CIEMAT |



LRU algorithm enabled in for xCache in the RAID6 server



~45k files

File deletions

File creations

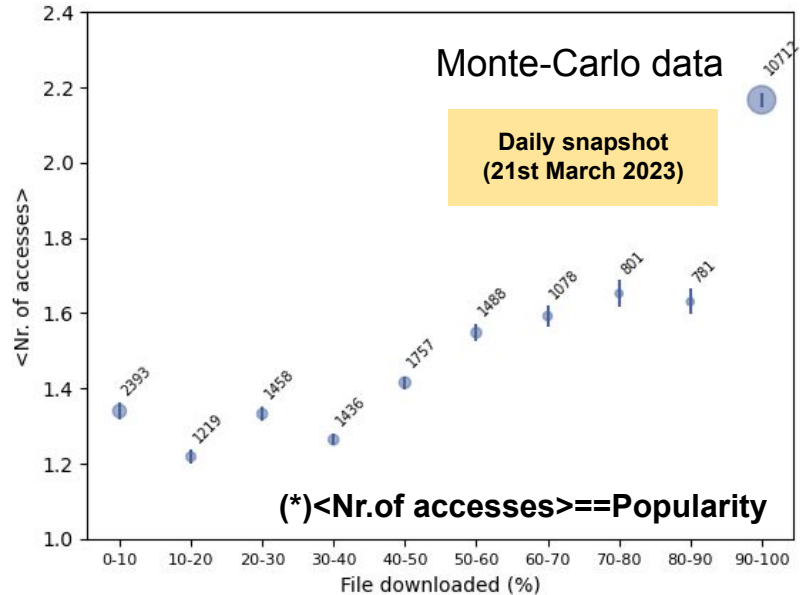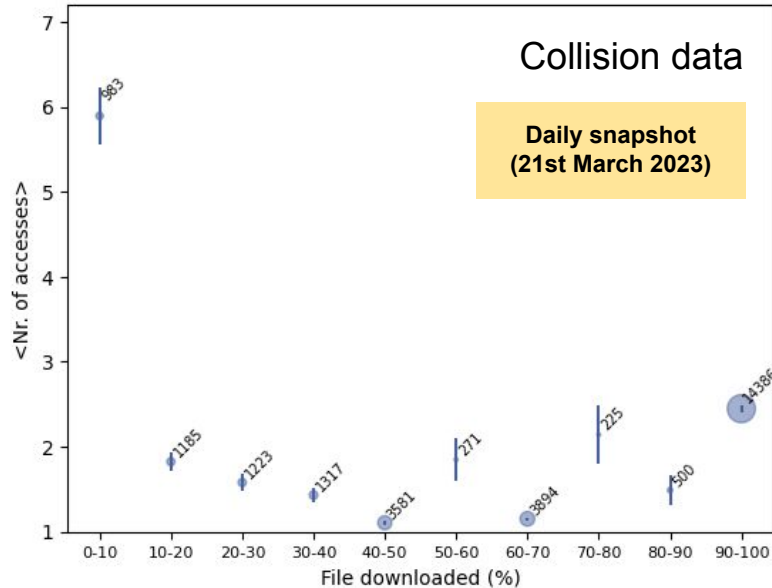~15k files

~2500k files

Timestamp per 3 hours

# Data popularity at PIC XCache



Average number of files accesses for both 'data' and 'mc' file types stored in the cache. This corresponds to a daily snapshot taken from .cinfo files

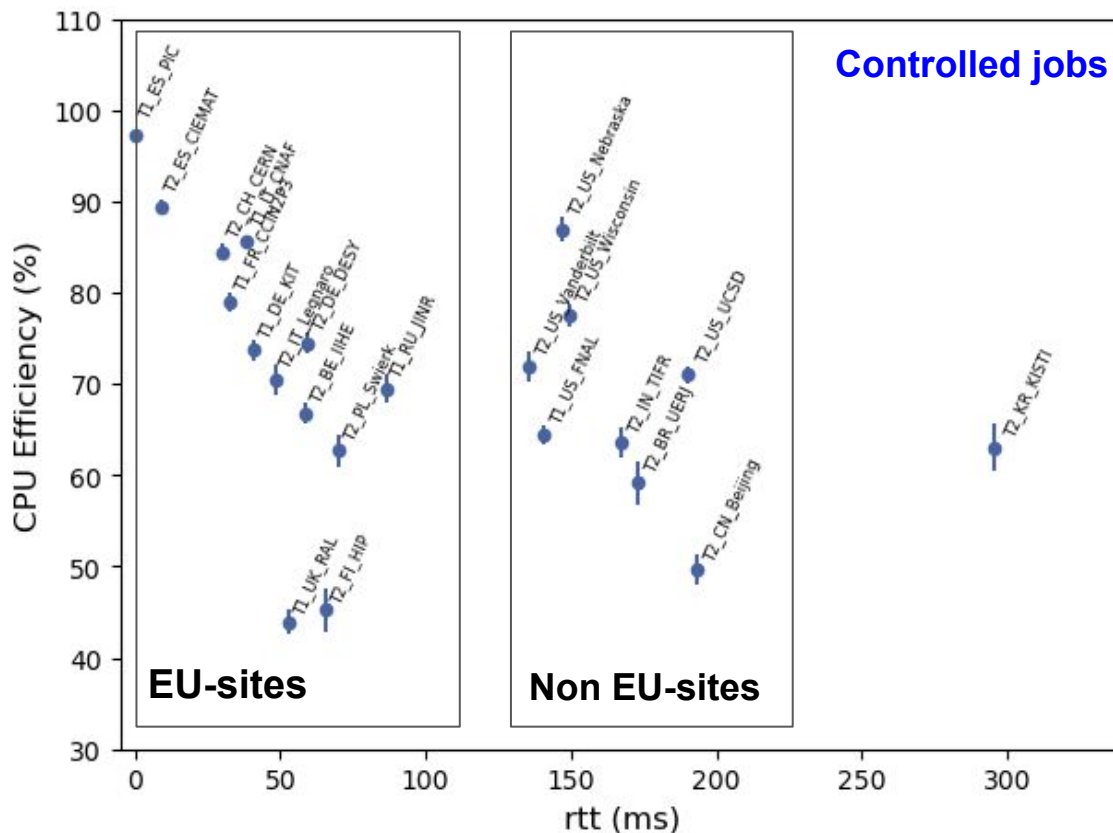# Files download percentage at PIC XCache



Collision data

**Daily snapshot (21st March 2023)**



Monte-Carlo data

**Daily snapshot (21st March 2023)**

(*)<Nr.of accesses>==Popularity

Collision and Monte-Carlo data show a **different behavior of re-accessibility.** In both cases, **the intermediate partially downloaded files are less re-accessed.**

# CPU efficiency dependence of remote reads

We conducted controlled tests to assess the qualitative **dependence of CPU efficiency on remote reads** from **various CMS sites** in and outside of Europe

• Our test analysis jobs **read MiniAOD files from several sites** world wide. Note, however, that **this is not how CMS generally undertakes analysis workflows** → CMS late-binding.
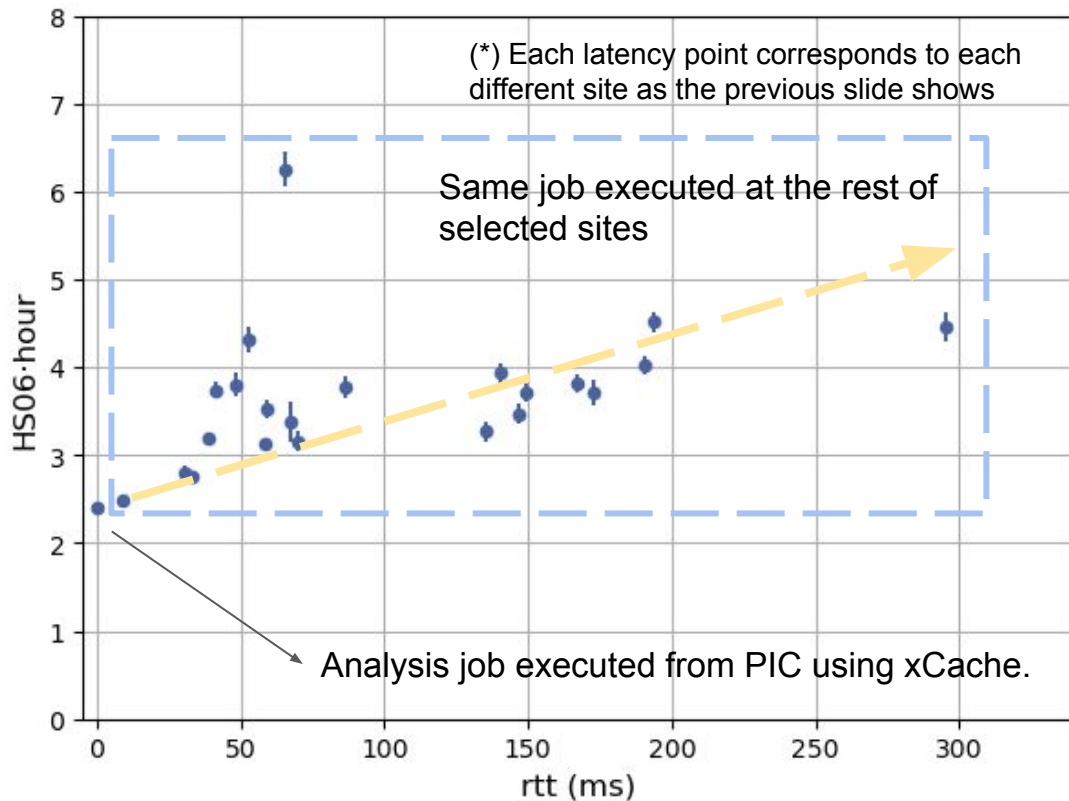
• **Controlled, uniform test** jobs allow us to **minimize errors** and make comparisons between source sites.

# CPU efficiency dependence of remote reads (2)

The cornerstone of the study is to evaluate which is the **impact over CPU efficiency of executing the same job accessing similar MINIAOD files from different sites**.

We have measured how **the CPU efficiency improves by executing** the same job reading **locally from the xCache** compared to doing it remotely -> **HS06·hours 'loss'** during the execution (see the figure)
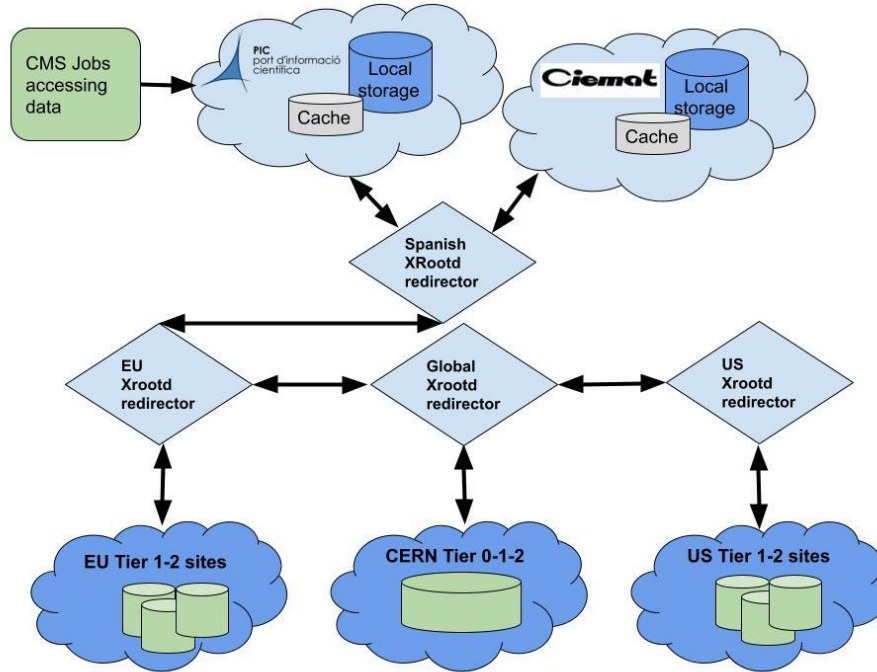


(*) Each latency point corresponds to each different site as the previous slide shows

Same job executed at the rest of selected sites

Analysis job executed from PIC using xCache.

# Conclusions

-During the initial phase of our project, we deployed the XCache from both PIC and CIEMAT separately to compare and validate the effects of different configurations on the service.

-Subsequently, we deployed the monitoring service and unified the service in PIC to efficiently and effectively serve data to the Spanish region.

- There is still some room to improve the re-accessibility of cached data →
Update the the site configuration with the actual data popularity of the service in production.

-The results of executing controlled analysis jobs remotely accessing MINIAOD at various sites reveal a significant degradation of CPU efficiency with latency due network and distance (as expected)

# Outlook

-The ongoing studies over the degradation of CPU efficiency will be complemented by dedicated studies evaluating the potential walltime saved by caching these data.

-Also, bringing the data closer to the cache increase the CPU efficiency of those jobs on our nodes. This effect has also to be well evaluated in jobs in production.

-Since xCache service is expected to alleviate the storage costs in the future, more dedicated studies will tell us how much storage we can save by keeping the most popular data for Analysis jobs within the cache.

# Backup slides

# XRootd redirectors infrastructure with the inclusion of the XCache service in PIC and CIEMAT sites

# Conclusions and actions taken in the first period (2022-Feb 2023)

Due the incidents experienced, **we decided to exclude PREMIX from the cache:** CMS jobs access the PREMIX data in non-uniform way / only used once or very infrequently → **(see the Hit-rate for PREMIX period at PIC figure).**

Similar approach was made for **PrePremix and unmerged data**, since several Processing jobs also failed due accessing these data.

**Hit rate vs time**

● hit rate 37.5

T1 PIC XCache during the caching-all period
(2022-Feb 2023)

# Distribution of data cached by XCache for PIC+CIEMAT



AOD - 5.2256% (8.9497)
AODSIM - 0.1944% (0.3329)
GEN - 0.0028% (0.0048)
GEN-SIM - 21.3933% (36.6395)
GEN-SIM-DIGI-RAW - 30.7773% (52.7110)
GEN-SIM-RAW - 1.3231% (2.2660)
GENSIM - 0.0010% (0.0017)
MINIAOD - 15.9326% (27.2870)
MINIAODSIM - 4.9939% (8.5529)
NANOAOD - 0.0009% (0.0015)
NANOAODSIM - 0.0029% (0.0050)
RAW - 20.1507% (34.5113)
RAW-RECO - 0.0015% (0.0026)

# Distribution of data cached by XCache for PIC+CIEMAT without unmerged



AOD - 5.2277% (8.9471)
AODSIM - 0.1763% (0.3018)
GEN - 0.0024% (0.0041)
GEN-SIM - 21.3859% (36.6018)
GEN-SIM-DIGI-RAW - 30.7983% (52.7110)
GEN-SIM-RAW - 1.3109% (2.2437)
GENSIM - 0.0010% (0.0017)
MINIAOD - 15.9416% (27.2839)
MINIAODSIM - 4.9914% (8.5428)
RAW - 20.1644% (34.5113)

# High failure rate Processing jobs

## Caching PrePremix and unmerged



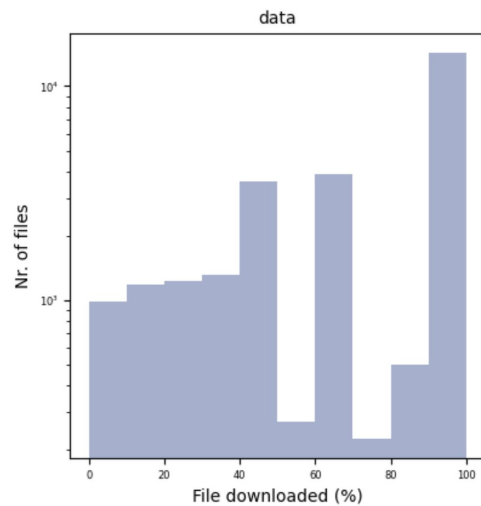## Currently excluding PrePremix and unmerged ->

Plots (PEPE)

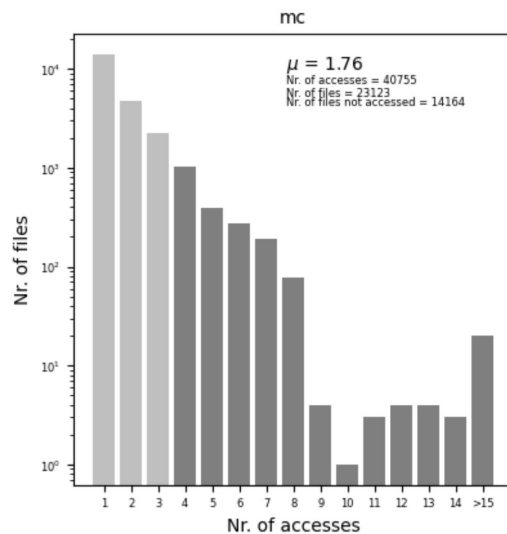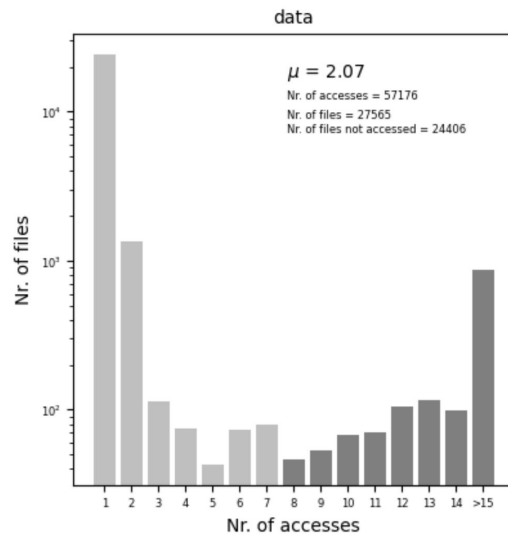Remote sites (3200 HS06·hours tests)
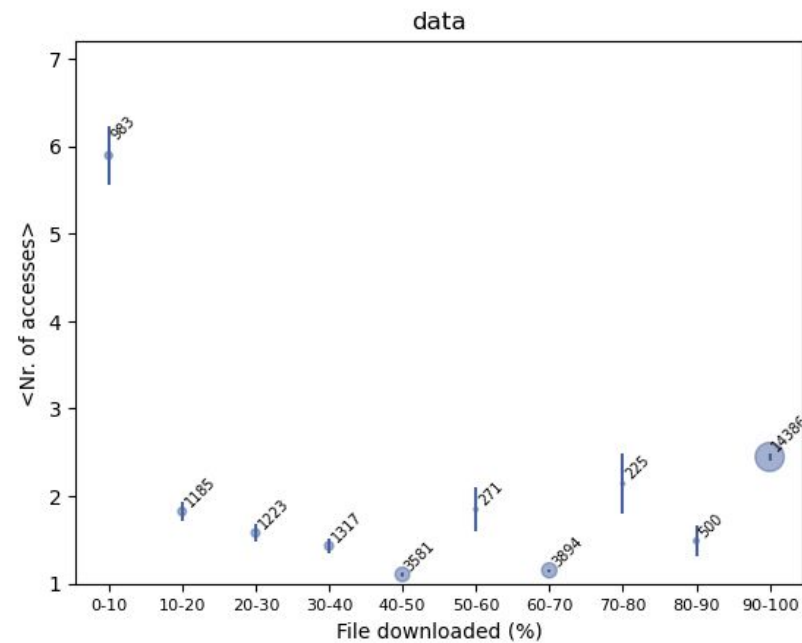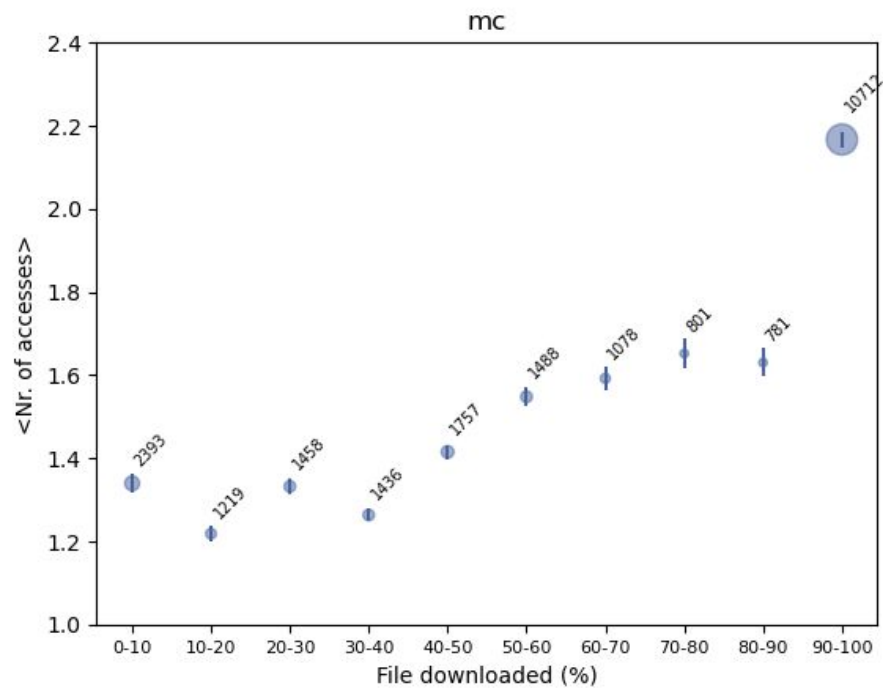
# Cache content

## He filtrado

- fichero de tests que usamos en los tests controlados de Muon Analysis
- /store/test, /store/unmerged, *PrePremix*
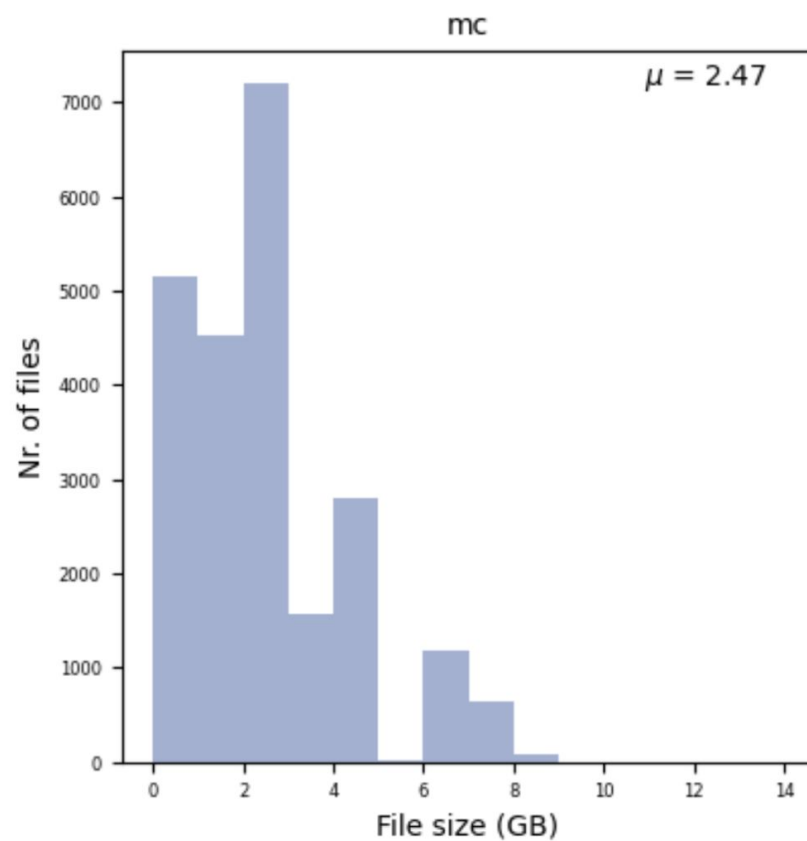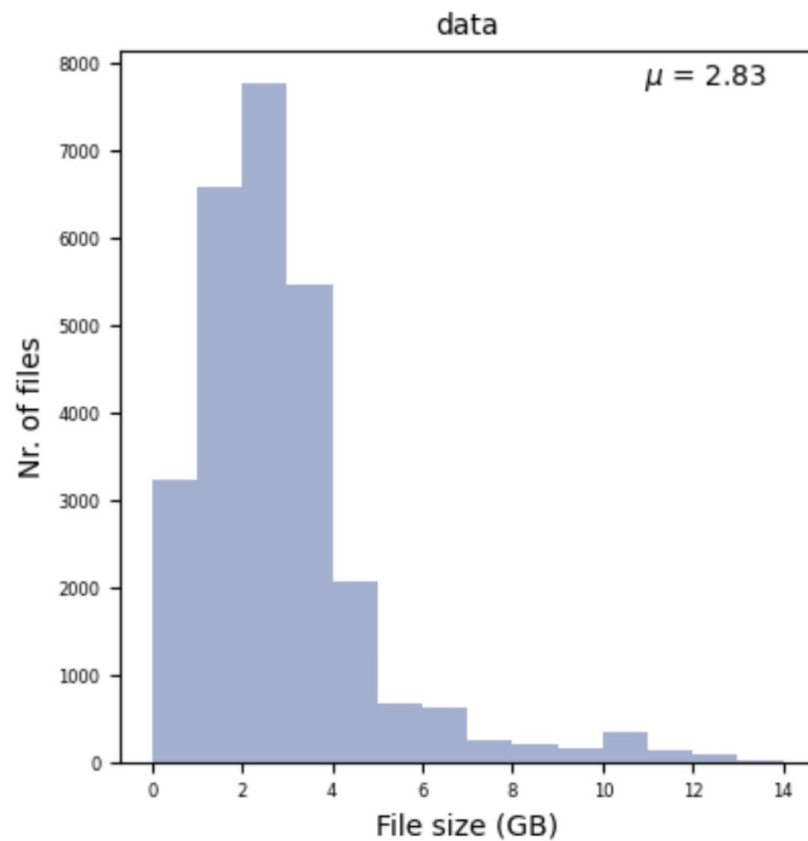- Todos aquellos ficheros en los que no se ha bajado ningún byte



Total Size = 77.95 TB (data: AOD, MINIAOD, RAW)



Total Files = 60198 (cmst3, data, group, mc, relval, user)



Total Size = 141.09 TB (cmst3, data, group, mc, relval, user)



Total Size = 57.21 TB (mc: AODSIM, GEN, GEN-SIM, GEN-SIM-RAW, MINIAODSIM)

# Cache content

# Cache content

# Cache content

# Cache content