

# Storage modeling

Markus W. Schulz

Andrea P. Sciabà

HOW Workshop  
JLab, 20-03-2019

# Outline

- Current activities and plans
  - Relationship with DOMA Access
- Alternative storage models
  - Site caches
  - Experimental measurements
  - Modeling
  - Dealing with data loss
- Impact of latency and bandwidth limitations
  - Measurements and emulation
- Data popularity and storage utilization
  - Access frequencies, lifetime of data replicas

# Current activities and plans

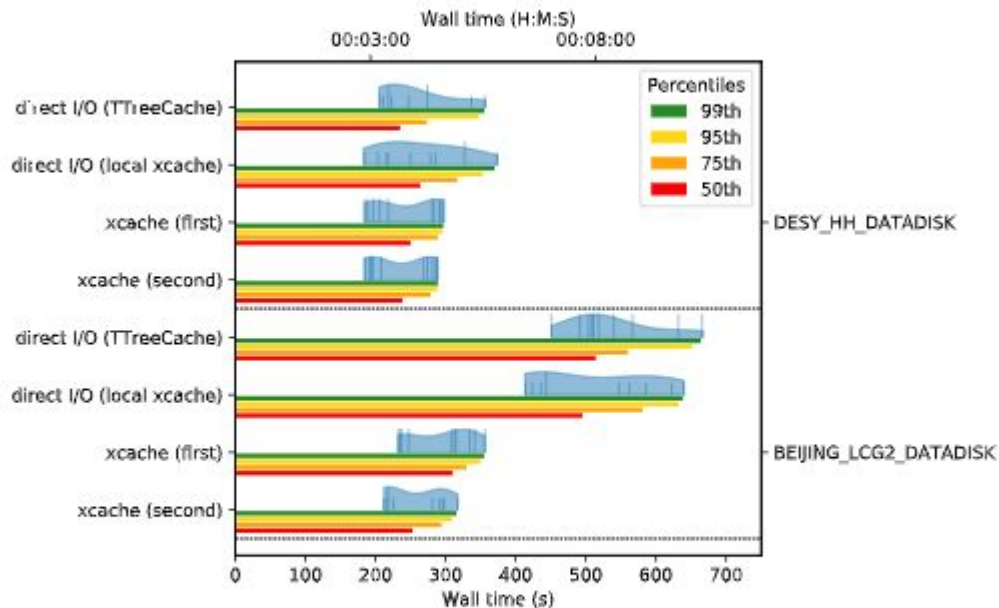
- Last summer we started to investigate the feasibility of  **caches**  and the impact of  **latency**  on workloads, and created tools
- Work continued and broadened by our WG
  - Now relevant work at a much larger scale by many players is done in the framework of DOMA
  - Example: <https://indico.cern.ch/event/769502/>
  - Please take a look at what has been presented at DOMA-access for an overview
- Should we continue independent work in this area or link this directly to DOMA activities?
  - In any case results and conclusions have to be presented to both activities
  - Discussions on investigations will find a larger audience of storage experts at DOMA meetings
  - The impact on cost is better addressed at the cost model meetings
  - **Can it be agreed?**  Voice your disagreement in case :)

# Example: LMU

- Only one of many detailed measurements
- Also studied load on cache server etc.
- Low end hardware used (2012 server)

## Processing from different sites

Derivation Jobs ( $\approx 3\text{MB/s}$ ) - process 500 Events



- Differences for direct I/O and cached visible for far away sites
- Local XCache (on each node) can serve as alternative to TTreeCache

# Alternative storage models (Data Lake Strawman)

- Studies of the **impact of data losses** in systems with low or no local data redundancy
- Based on the Data Lake strawman model
  - Based on CMS analysis model
  - Spreadsheet to evaluate different scenarios, usag patterns, replication rates .....
    - <https://docs.google.com/spreadsheets/d/12bmAPWUzsZrDtptJTfGyR-Rw8wPv5D8U3tbNSwDNbJo/edit?usp=sharing>
  - Part of the Data Access on a Data Lake straw man model document
- These studies are best done within DOMA
- **Impact of cost** has to be taken into account for our working group

# Example from DOMA access

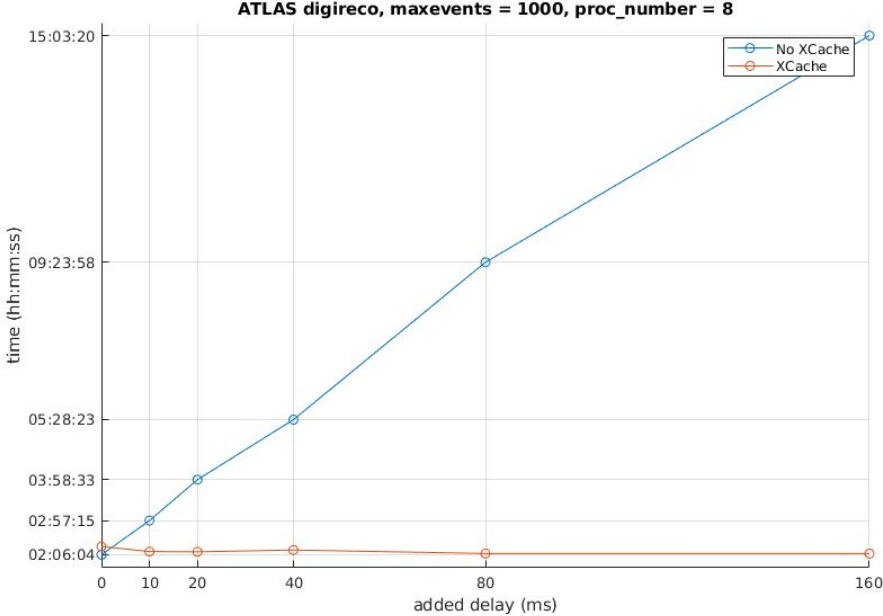
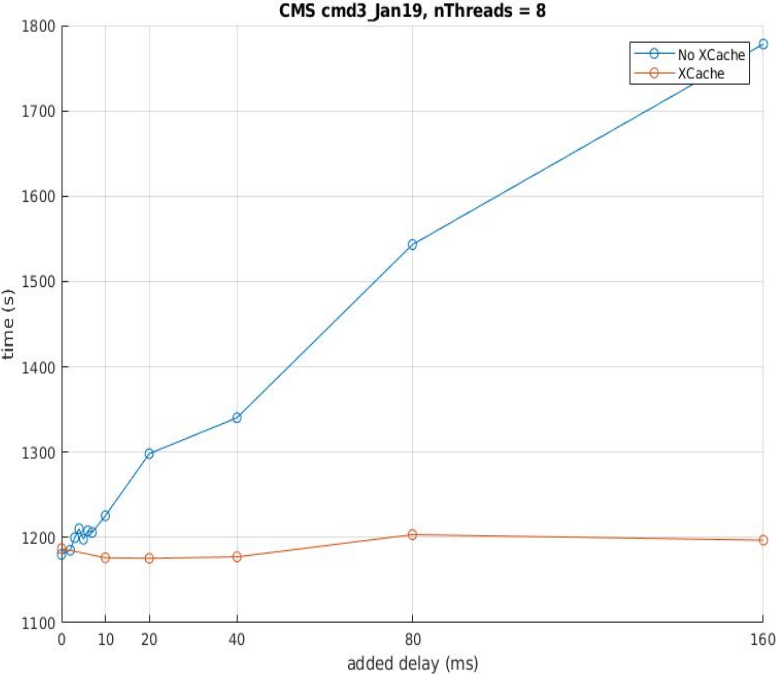
- Example: Disk fails, data is replicated from other sites (at low rate)
  - Based on known failure rates and the straw man model (see Xavier's presentation)
    - Can't be understood without the document and the spreadsheet
  - Impact is minimal, compared to normal rate of failed jobs.

Total Number of files on site	66666666
number of files accessed during 1 hour [1/h]	277778
fraction of total number of files accessed in an hour [1/h]	0.0042
fractional size of the failed disk	0.0001
number of files accessed on the failed disk per hour [1/h]	28
files missed during locating replicas	10
files missed during replication (files are gradually moved)	77
<b>total number of files missed during recovery (6h)</b>	<b>87</b>
total number of files accessed during recovery	1646090
Fraction of files missed during reco period	0.000053
Above in ‰	0.53
<b>Average file miss rate in ‰</b>	<b>0.036</b>

# Example 1: effect of latency and the impact of XCache

- Ingredients
  - Reference workloads
  - Coentini's tool (see Serhan's presentation)
    - Tool for latency, bandwidth and memory restrictions
  - XCache instances on standard CERN nodes
- Workload is run on a node, reading data from another node
  - With **added latency, without a cache**
  - With **added latency and a cache**
  - Results are preliminary, but everything indicates that **Xcache is very good at latency hiding**
    - Even when data is read for the first time/once!

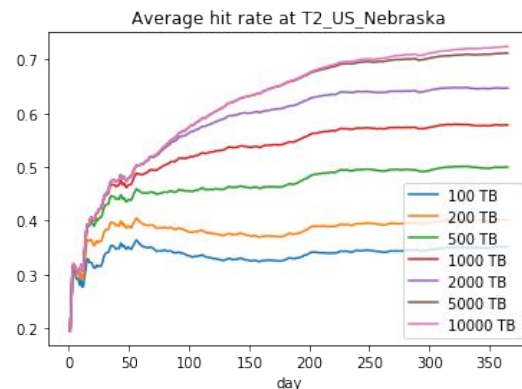
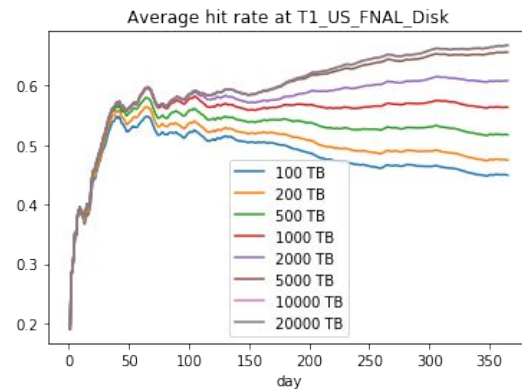
# Example: CMS RECO and ATLAS digireco





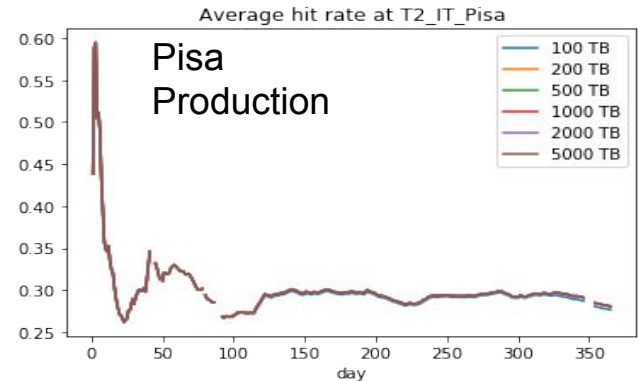
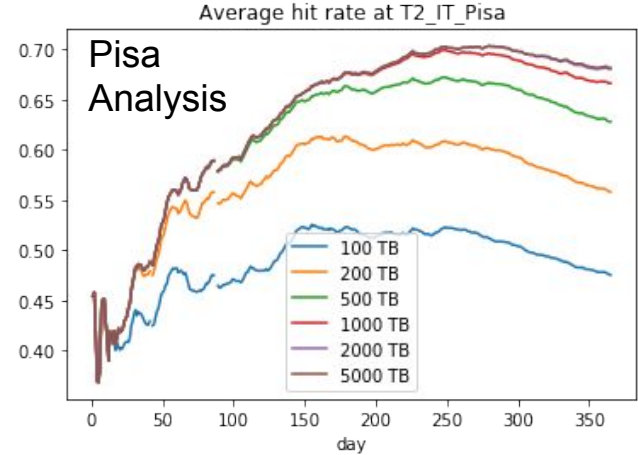
# Example 2: simulating a site cache at CMS sites

- Inspired by previous studies on ATLAS popularity data
- Ingredients
  - CMSSW popularity data
    - Site, file name, file size, access time
- Data provided by ATLAS and CMS is much richer than required by these studies
  - Preprocessing is required in all cases
  - Will propose a common intermediate format for people doing studies



# Example 2: simulating a site cache at CMS sites

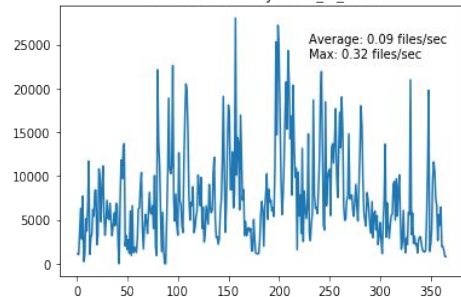
- Difference between Analysis and Production
  - But can push analysis data out of the cache
  - Hit rate independent from size
- Therefore:
  - One larger cache for analysis data
  - One smaller cache for the production files
    - Still provides latency hiding
    - Big enough to cover the load of a few days (for failed jobs)



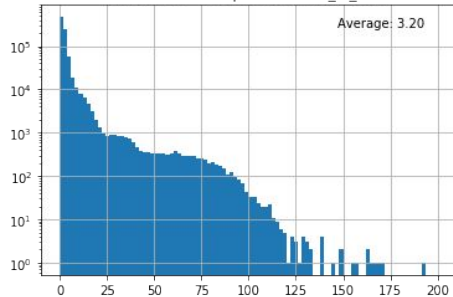
# Example 2: frequency and number of file accesses

- For input files registered in the DBS, measured
  - distribution of number of accesses
  - Files read / sec (day by day)

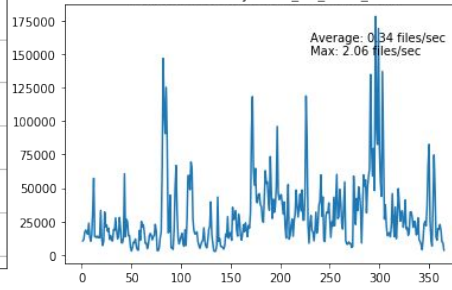
Files read vs day at T2\_IT\_Pisa



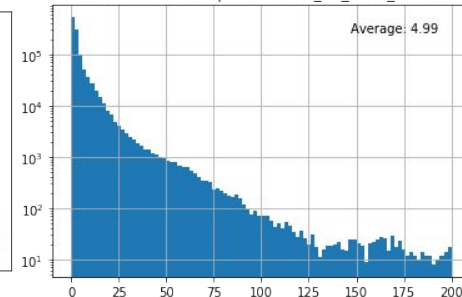
No. of accesses per file at T2\_IT\_Pisa



Files read vs day at T1\_US\_FNAL\_Disk



No. of accesses per file at T1\_US\_FNAL\_Disk

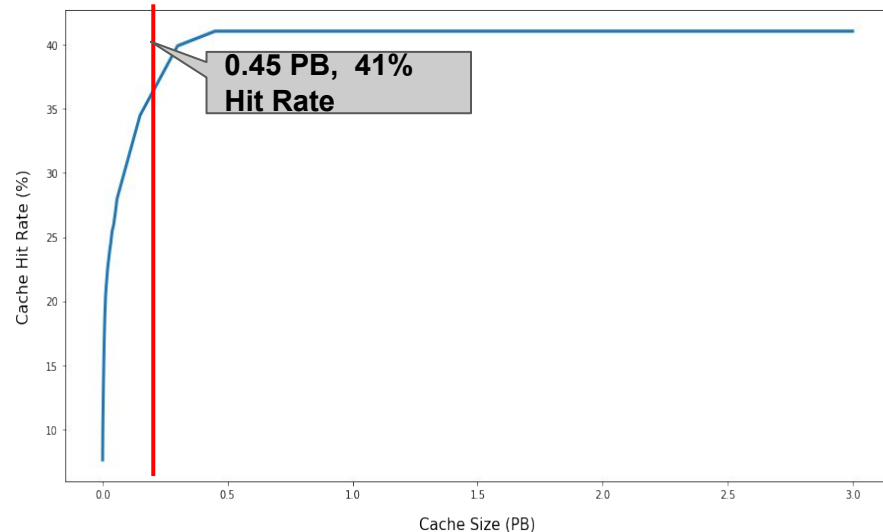


# Example 3: ATLAS data popularity studies

- ATLAS studies started last year
  - Based on half a year of RUCIO access data
    - These traces do not cover the access to final analysis product (ntuples)
  - Started to look also at staging traces
    - In different (better) format
    - Combining will be difficult
- All work based on the current analysis model
  - Which is quite different from the future model
- ATLAS is discussing a new model
  - <https://indico.cern.ch/event/769501/>
  - In this model new DAOD formats are introduced (PHYS/PHYSLITE)
    - Smaller, less versions, DAOD production from TapeCarousel

# Follow Up on Cache Studies with ATLAS data

- Based on 1 month of logs
- Picked PragueLCG2 as a “typical” T2
  - 32k cores, 6 PB
- Simulated cache
  - Hit rate/ cache size
  - Repeated later with data-served-from-cache/total-data-read (small difference)
- AOD+DAOD+HITS = 87% of data
  - AOD+DAOD 56%
- Simulated caches for different data types
  - AOD + DAOD: 61 % @ 256 TByte
  - HITS : 90% @16 TByte

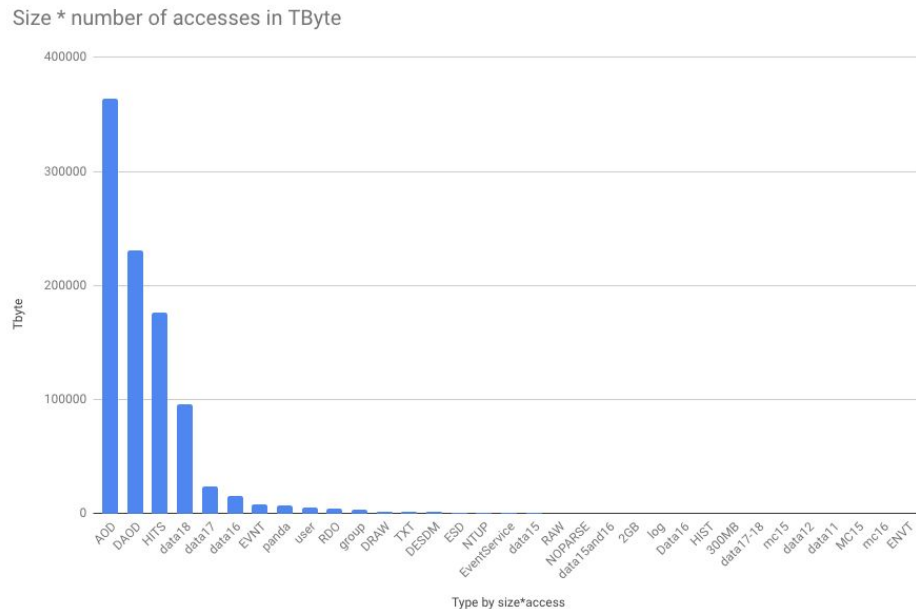


# Additional Work:

- David Smith started work on a **stress test system for caches**
  - Using access records from logs
  - Using profiles from measured workloads
- Goal is to understand **what performance is needed by a cache node** to handle realistic site loads
  - From this **cost for caches can be derived**
    - Based on site cost models
    - Human effort still to be evaluated
      - Feedback from sites using caches needed ( setup/ops)
      - Since Xcache has been containerised this shouldn't be too hard...

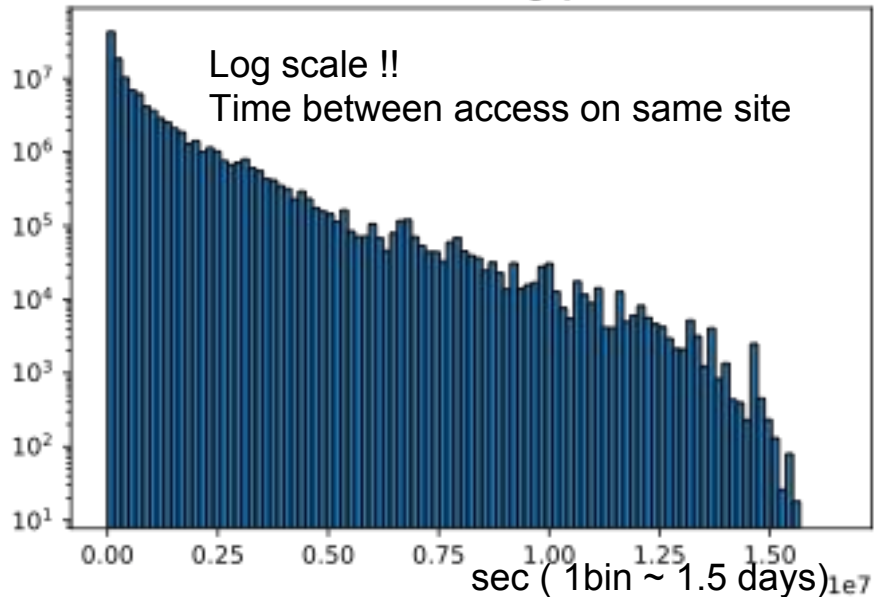
# Some work on global data access patterns

- ATLAS data access logs
- Data type (AOD/DAOD/HITS...)
  - as expected
- Looked at “impact” = number of accesses \* size
- Looked at many different aspects
  - Time between access
  - Number of sites
  - Time between first and last access
  - Number of accesses
  - Correlations....

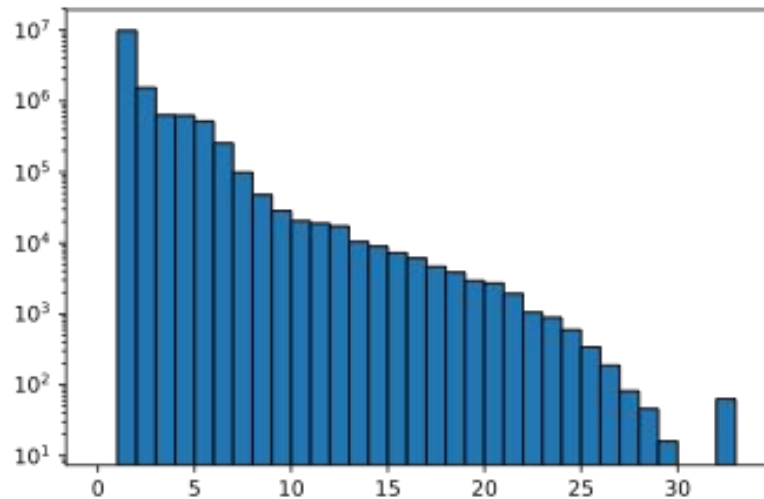


# Some examples: AODs

AOD : time gap



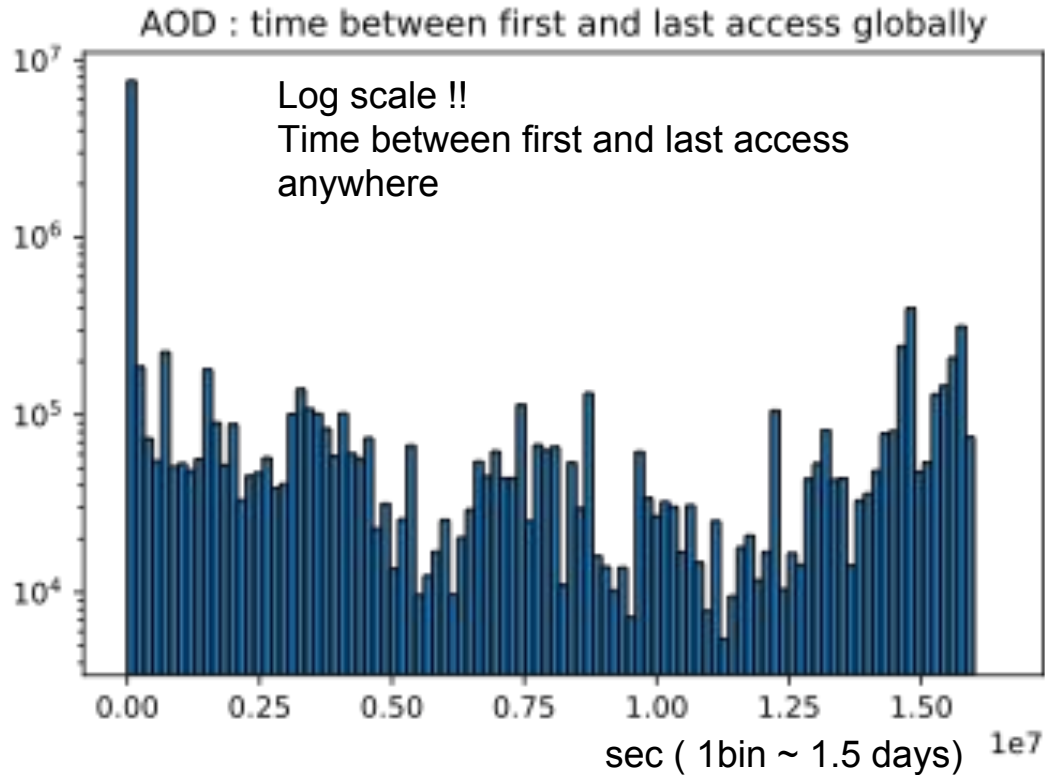
AOD : number of sites with the same files



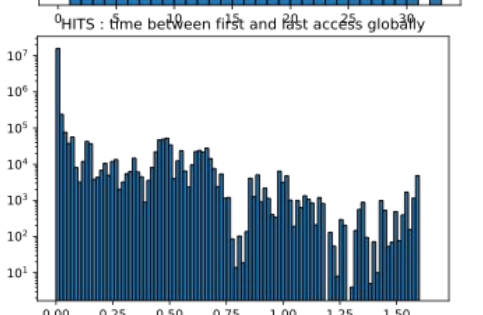
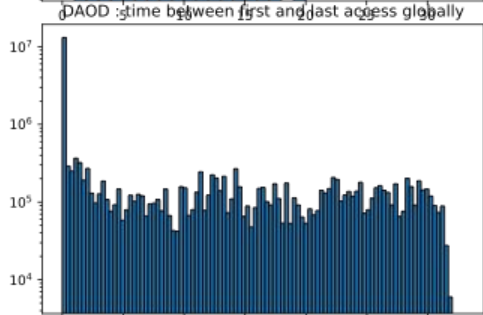
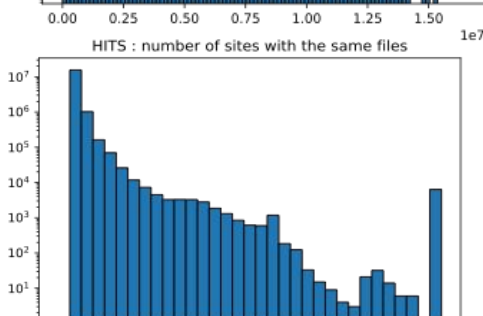
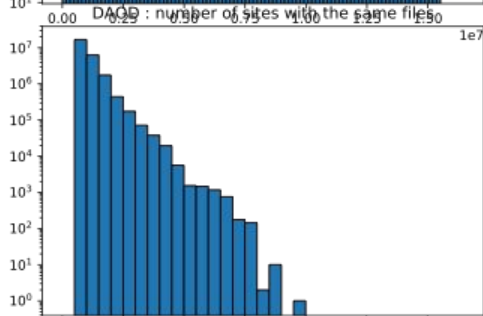
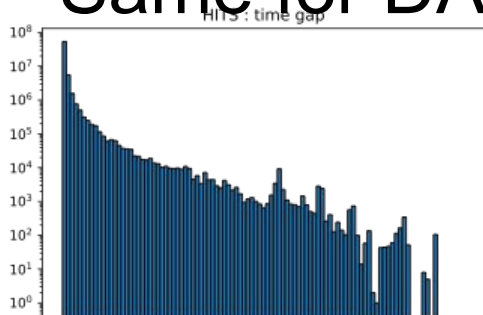
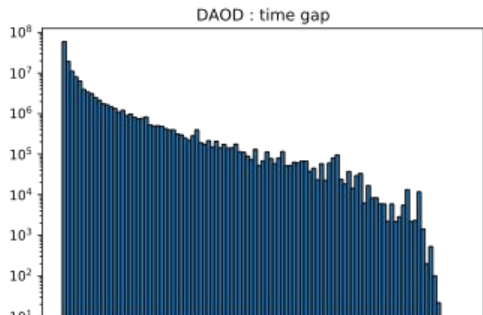
Number of sites on which the same file is read  
Log scale



# AODs



# Same for DAODs, HITS....



- And many more....
- What could be learned?
  - Data isn't accessed very often
  - Most likely to be re-read within days
  - Only on log scale structures become visible
- What is missing
  - A lot .....
  - 6 months isn't long enough !
  - Need to add staging and deletion information
    - To measure "active" vs "passive" time
  - Looking for access rate/absolute time → seasons etc.

# Similar studies have been done for CMS Data

- Andrea Sciaba
- CMS data contains all accesses
- To be discussed

# Data Access and popularity study at PIC

- PIC Tier-1 is doing an analysis of the CMS data access and popularity based on dCache billingDB
- Looking into **file accesses**:
  - Accesses from **remote or local** IPs
  - **Data type** (MC, Data, and the type of data accessed: RAW, RECO, AOD)
  - Time since creation to first access - number of accesses - time from last access to deletion
  - **Bytes transferred** from accesses
  - **Usage of the disk space** (files accessed as compared to total files stored, as a function of a sliding window)
- Millions of files accessed per month - complex analysis
- Once the procedures are setup, there is the plan to compare to a Tier-2 (CIEMAT, Madrid)

# Summary

- Ours and DOMA access studies indicate that caches can have a huge impact on how storage is organized
- Have to **derive the cost impact** from the measurements
  - In terms of storage and compute **resources**, this is straightforward
  - **Network** cost is more complicated, due to the step function when current bandwidth limits are reached
  - **Operation** costs differences between managed storage and caches are difficult to estimate at the moment
    - With more and more “hands-on” experience it will become feasible
- Data formats and analysis strategies are currently in flux
  - Focus on developing **flexible approaches** rather than very detailed analyses