# XCache studies and Logfile analytics

G. Duckeck, <u>N. Hartmann</u>, C. A. Mitterer, R. Walker, E. Schanet

LMU Munich

October 1, 2019

# XCache

# What is XCache?

- Disk caching proxy using xrootd (`libXrdFileCache.so`)

- Data is cached in blocks

- Simply prepend xcache server url - e.g.
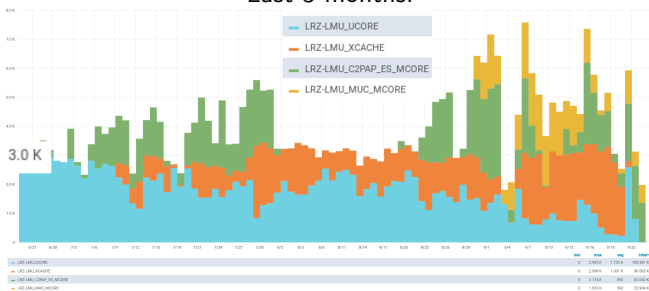  `TFile::Open("root:[xcache-server]:[port]//[xrootd-path]")`

- Optionally use rucio DIDs via N2N plugin:
  https://github.com/wyang007/rucioN2N-for-Xcache
  $\rightarrow$ allows usage of rucio DIDs instead of xrootd path
  $\rightarrow$ tracks identical files distributed at different locations
  (internal symlink `.../scope/XX/YY/filename`)

# Setup

- Hardware: Old dCache pool node (from 2012):
  - Dell R710, 2x6 core Xeon L5640, 32 GB RAM, 10 Gb Ethernet
  - 60 TB Raid-6 (2x12x3TB HDD)
- Xrootd version 4.10.0
- Setup w/ singularity SL6 image. Full configuration: https://gitlab.physik.uni-muenchen.de/Nikolai. Hartmann/xcache-singularity-lrz/
- XCache settings:
  ```
  pfc.ram 14g
  pfc.blocksize 1M
  pfc.prefetch 10
  ```
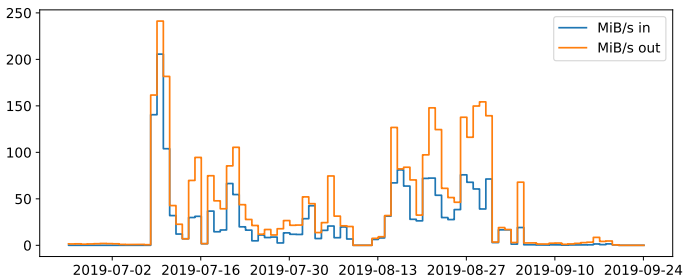
# Test XCache in ATLAS production queue

Last 3 months:



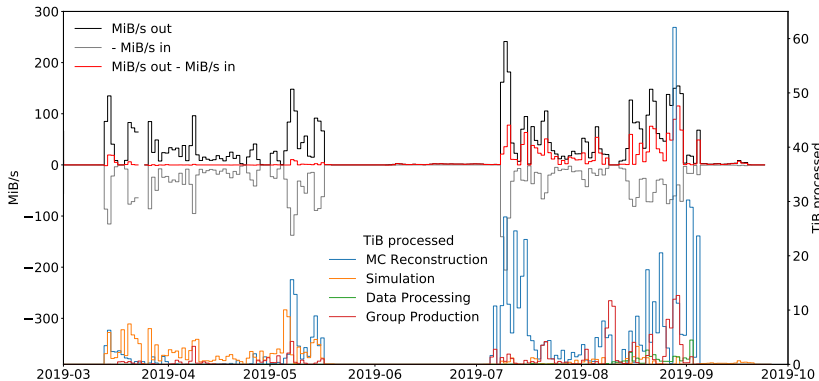ATLAS production queue in Munich that retrieves all files via XCache

- Remote destination is nearby MPP Munich storage
- Can take a quite significant fraction of the jobs
- Works surprisingly well, given that all trafic goes throuh a single server

# Caching works



$\rightarrow$ Output volume already larger than input volume ($\approx 1.8$)
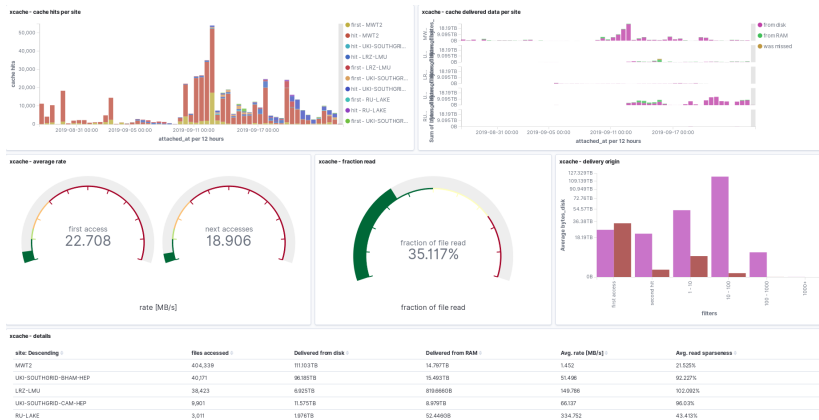
# But hit rate depends on type of job



$\rightarrow$ largest hit rate for MC Reconstruction (here mainly pileup overlay)

# Central monitoring for ATLAS XCaches

Since a few weeks we are (together with other ATLAS XCaches)
monitoring file access statistics to an ElasticSearch instance in Chicago

# Bugs/Issues

Found 2 Problems when XCache is under high load:

- Number of open files increasing until system limit is hit
  (https://github.com/xrootd/xrootd/issues/975) → fix in
  work
  → partially mitigated by settings: `pss.ciosync 60 900`
- Segfaults/Crashes
  (https://github.com/xrootd/xrootd/issues/1026)
  → mostly fixed in xrootd 4.10, but occasionally still seen for very
  high load (pileup jobs)

Lead to corrupted files: wrong checksum for file in cache, $\approx$ 90 out of
200k files
→ not observed any more after fixes/mitigations
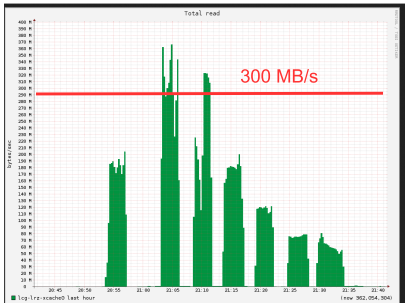→ still, we want to have a check for corrupted files in the future

# Performance for parallel reads - Raid6 vs single disks

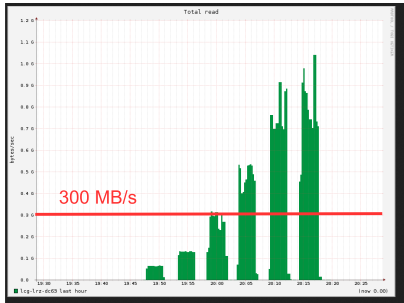Feedback from xrootd developers: Use multidisk-mode instead of Raid

Performance test at LRZ:



Raid 6

300 MB/s

1   2   5 10  20 40
simultaneous reads

individual disks

300 MB/s

1   2   5 10  20 40
simultaneous reads

→ multi-disk mode might perform better than Raid for caching system
→ similar test with additional 50% writes gives the same picture

# Outlook

- Test Multi-disk mode instead of Raid system
  $\rightarrow$ claimed to perform better, tests seem to confirm

- Investigate more use cases for caching:

  - Analysis jobs/Direct read instead of copy-to-scratch
    $\rightarrow$ continue tests, saw issues with long running jobs

  - Test XCache in columnar data analysis (e.g. with Pandas/Dask)