# The Analysis Facility at GSI

Sören Fleischer     Raffaele Grosso

GSI

September 27, 2022

- The ALICE T2 at GSI
  - the cluster
  - software solutions
  - current issues
- Prospects for the ALICE Analysis Facility at GSI
  - AF requirements
  - current state

Resources reserved on a shared cluster:

▶ ~13k logical CPU cores (hyperthreading):

| PARTITION | CPUS | S:C:T | MEMORY | NODES |
|-----------|------|-------|--------|-------|
| grid | 96 | 2:24:2 | 191762 | 134 |

under a SLURM workload manager

Resources reserved on a shared cluster:

- ∼13k logical CPU cores (hyperthreading):

  | PARTITION | CPUS | S:C:T | MEMORY | NODES |
  |-----------|------|-------|--------|-------|
  | grid | 96 | 2:24:2 | 191762 | 134 |

  under a SLURM workload manager
- 5.2 PB (in 2022 $\Rightarrow$ 6.1 PB in 2023) disk storage under a Lustre distributed file system

Resources reserved on a shared cluster:

- $\sim$13k logical CPU cores (hyperthreading):

  | PARTITION | CPUS | S:C:T | MEMORY | NODES |
  |-----------|------|-------|--------|-------|
  | grid | 96 | 2:24:2 | 191762 | 134 |

  under a SLURM workload manager

- 5.2 PB (in 2022 $\Rightarrow$ 6.1 PB in 2023) disk storage under a Lustre distributed file system

- Network connection
  - internally 100 Gb/s EDR Infiniband
  - 10 Gb/s LHCONE, 2 Gb/s DFN

Resources reserved on a shared cluster:

▶ ∼13k logical CPU cores (hyperthreading):

| PARTITION | CPUS | S:C:T | MEMORY | NODES |
|-----------|------|-------|--------|-------|
| grid | 96 | 2:24:2 | 191762 | 134 |

under a SLURM workload manager

▶ 5.2 PB (in 2022 ⇒ 6.1 PB in 2023) disk storage under a Lustre distributed file system

▶ Network connection

▶ internally 100 Gb/s EDR Infiniband

▶ 10 Gb/s LHCONE, 2 Gb/s DFN

▶ Memory limits imposed by SLURM via cgroups

▶ limit is set on PSS, thus correctly taking shared memory into account (4.4 GB per physical core)

▶ we don't set virtual memory limits

Jobs run within Singularity containers.

▶ Host: minimal Red Hat Enterprise Linux compatible installation

Jobs run within Singularity containers.

- Host: minimal Red Hat Enterprise Linux compatible installation
- Container:
    - used to be our own image
    - since the move to JAliEn-VOBoxes this is handled by JAliEn, using the Singularity runtime engine and image from `/cvmfs/alice.cern.ch/`.

Jobs run within Singularity containers.

▶ Host: minimal Red Hat Enterprise Linux compatible installation
▶ Container:
  ▶ used to be our own image
  ▶ since the move to JAliEn-VOBoxes this is handled by JAliEn, using the Singularity runtime engine and image from `/cvmfs/alice.cern.ch/`.
▶ In the last months all jobs have been going through the GSI-8Core queue
  ▶ JobAgents are run as SLURM jobs with 8 cores, which they fill with single- and multi-core AliEn jobs

Average: ∼4.3k jobs from Hyperloop trains, max 5600 running jobs

XRootD compiled with following self-developed plug-ins:
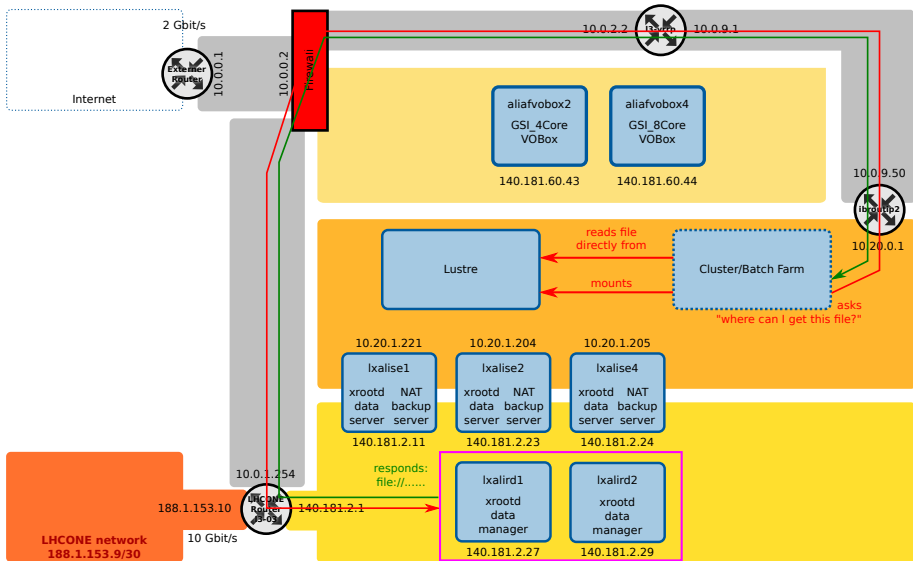
▶ symlink: locally creates an LFN-symlink pointing to the PFN ⇒ allows local access to files via LFN

▶ quota: calls `lfs quota` ⇒ allows MonaLisa to get the correct values for available storage from a shared distributed file system

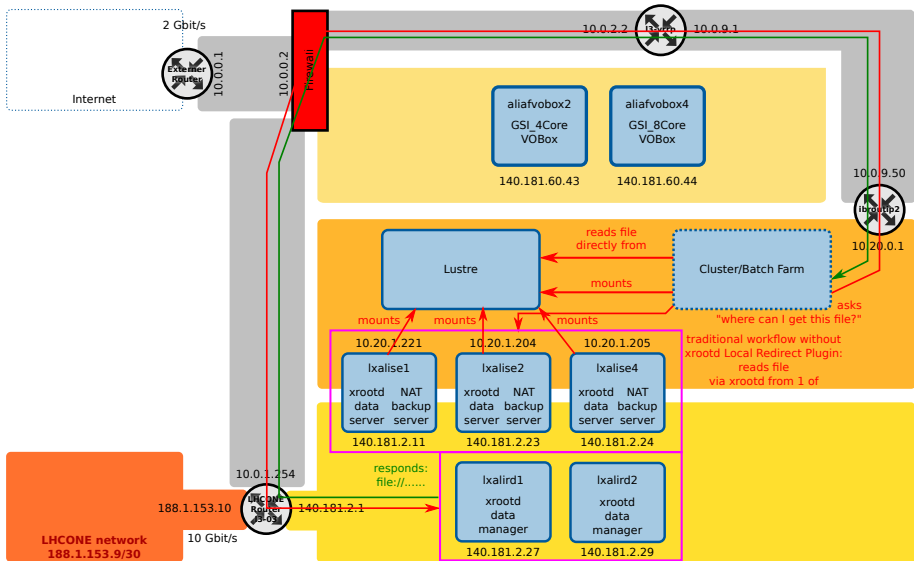▶ local redirect: see description in next slides ⇒ optimizes the I/O throughput of the analysis jobs

alinat

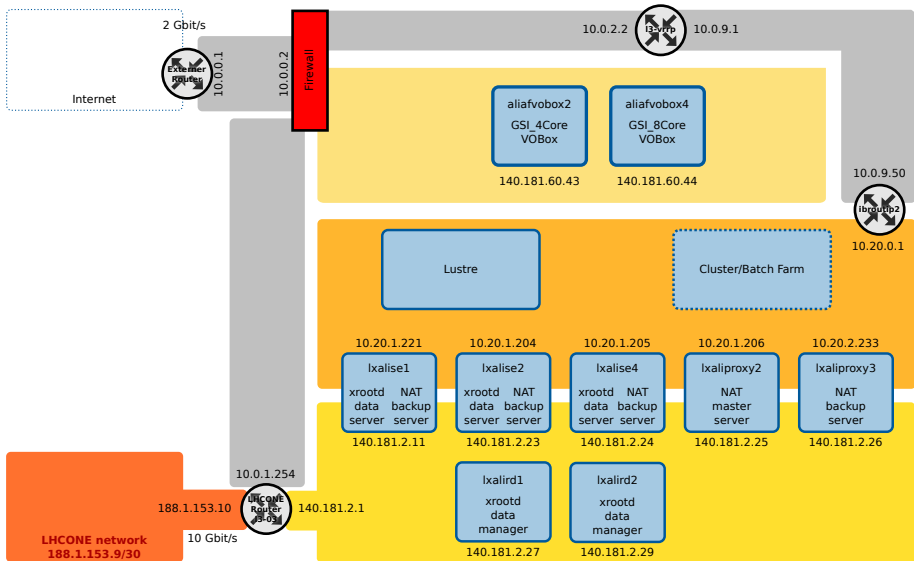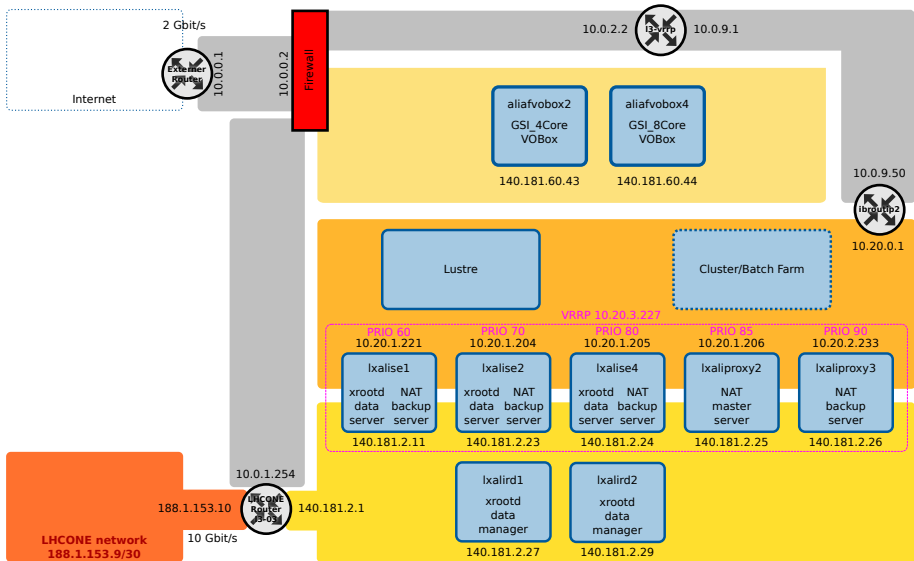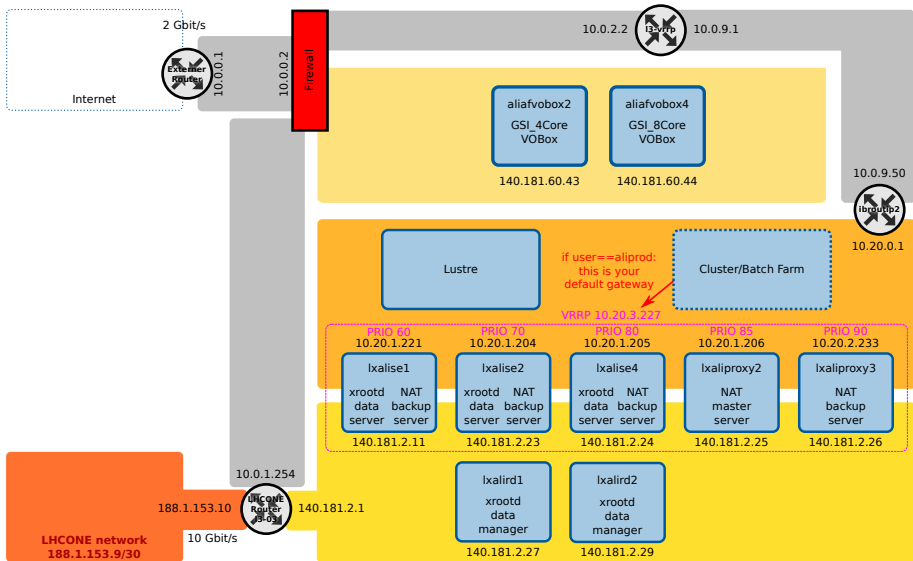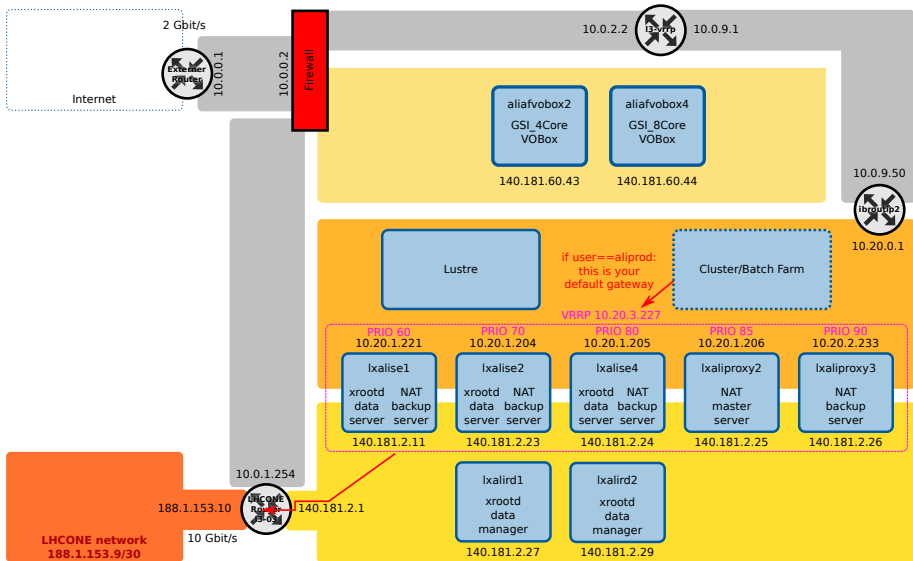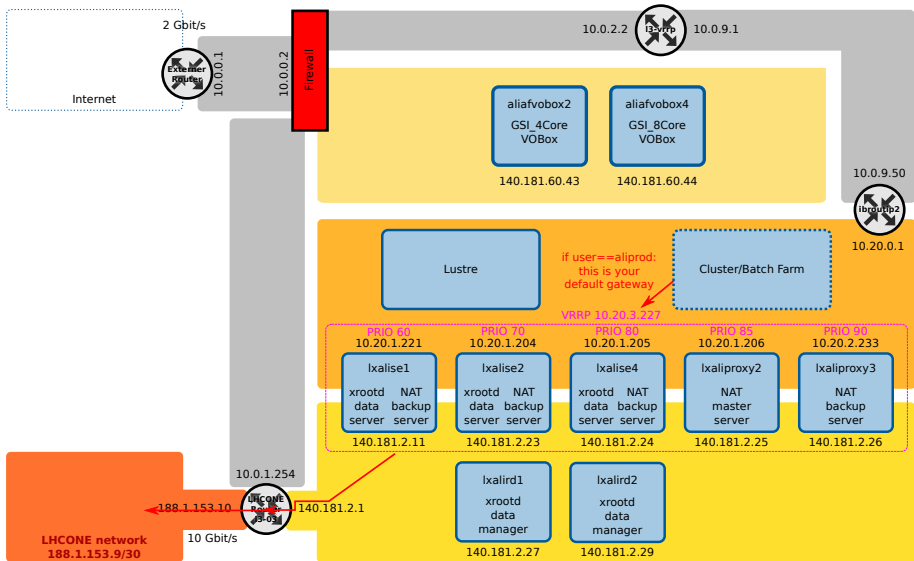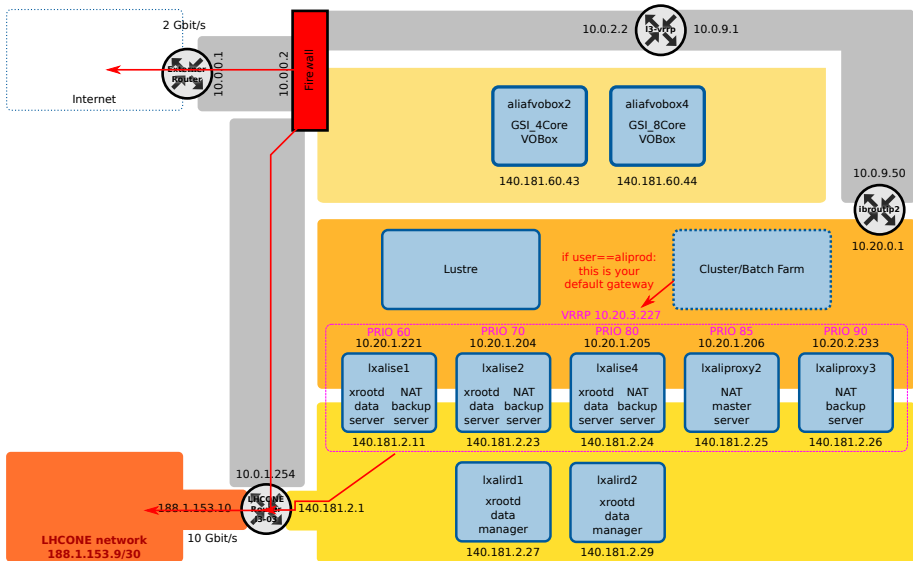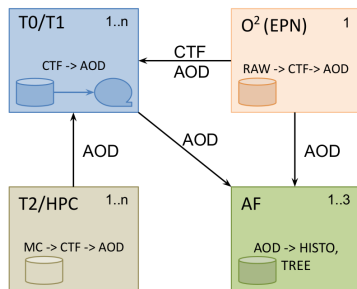▶ CE "freezes": is running but unable to spawn threads, in particular
  squeue commands ⇒ the VOBox silently stops spawning new
  JobAgents

  ▶ temporary solution: monitor log file to promptly restart ce
  ▶ long term solutions: fixes in JAliEn:
    ▶ CE must not fail silently
    ▶ understand why the critical condition is reached

▶ permissions: getting 0600 instead of 0644 with TPC transferred files
  Solution: being discussed with XRootD developers

▶ hard-coded list of Singularity bind mounted directories
  Solution: being discussed with JAliEn developers

- ▶ AFs supposed to provide 50% of CPU share for analysis
  - ▶ receive AODs from $O^2$ farm and T1/T2s
  - ▶ produce histograms and trees
- ▶ 10% of sampled AODs for quick analysis and cut tuning



Requirements:

- ▶ serve 6-8k job slots with ∼15 MB/s/core[1]
- ▶ aggregate throughput of 100 GB/s
- ▶ be able to digest more than 5 PB of AODs in a 12-hour period

---

[1]S.Piano, ALICE Week June 6th

- Already close to the required number of cores.
- Benchmark showed linear increase of throughput vs number of concurrent jobs up to 30 GB/s for 2000 jobs. Challenges to rely on linear increase up to 100 GB/s still to be faced.
    - improve utilization of the shared Lustre file system
    - improve and enforce algorithms for distribution of files over OSTs

Thanks for your attention!

Questions?