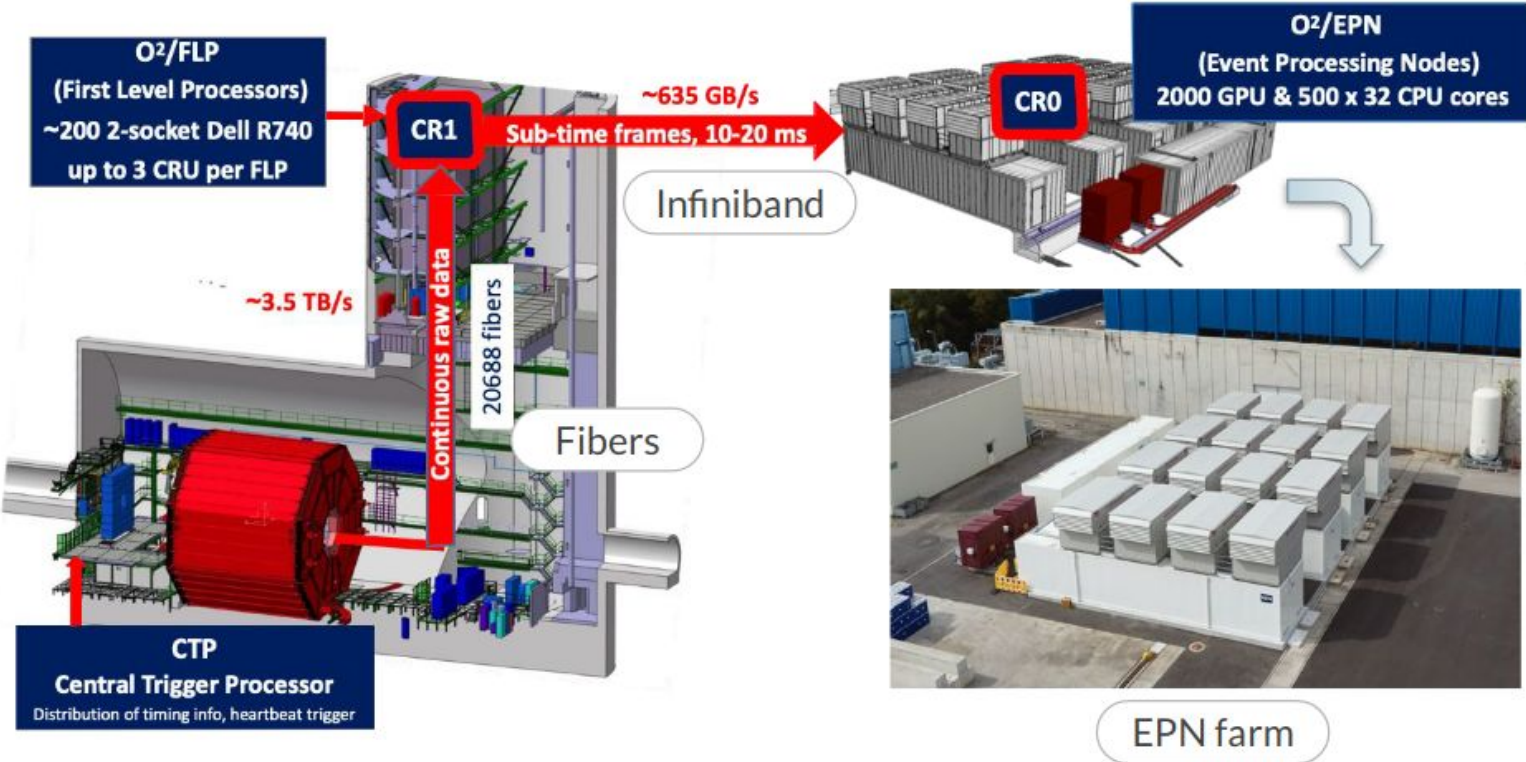# Calcolo ALICE in Run 3 su GPU e architetture eterogenee: validazione e attività in produzione

F. Noferini
INFN sez. Bologna

# ALICE DATA flow in Run-3

# ALICE on heterogeneous architectures

- ALICE underwent several upgrades in view of Run 3 to deal with a much higher luminosity wrt Run 1/2
- One of the major upgrades is $O^2$ (offline-online) software/framework with the goal to support continuous readout (triggerless) in data taking and synchronous processing/reconstruction
- Software was completely revised to allow for
  - a high level of paralyzation in all steps (data processing, MC, analysis)
  - allowing to exploit GPU resources (strictly retired for synchronous processing)
  - providing a common framework to cover all sync and syncs needs (data processing layer, DPL)
  - I/O and CPU optimization in analysis (RDataFrame)
  - …

# ALICE on heterogeneous architectures

- ALICE underwent several upgrades in view of Run 3 to deal with a much higher luminosity wrt Run 1/2
- One of the major upgrades is $O^2$ (offline-online) software/framework with the goal to support continuous readout (triggerless) in data taking and synchronous processing/reconstruction
- Software was completely revised to allow for
    - a high level of paralyzation in all steps (data processing, MC, analysis)
    - allowing to exploit GPU resources (strictly retired for synchronous processing)
    - Providing a common framework to cover all sync and syncs needs (data processing layer, DPL)
    - I/O and CPU optimization in analysis (RDataFrame)
    - …

$O^2$ software is flexible enough to support also other king of resources, e.g. HPC and ARM

4

# Potential applications

There are 3 potential areas where ALICE activity can be expanded (in **bold** what will be covered in this talk)

- **GPU**
  - **The role of GPUs in ALICE already presented in many discussion (e.g. <u>link</u>, <u>link2</u>). In this talk we will focus on its effectiveness in production activity**

- HPC
  - Some results on MC simulation were presented in the past →Marconi/Cineca (<u>link</u>) and other (<u>link</u>)

- **ARM**
  - **Some preliminary studies were performed and reported in this talk**

# Using GPU in Run 3 (reconstruction) with ALICE $O^2$

Processing on dedicated farm at experimental site

- 250x Event Processing Nodes (EPNs) with 2x32 core CPUs (to be expanded soon)
- 8x AMD Graphic Processing Units (GPUs)
- ~1600 GPUs required to process 50 kHz Pb-Pb collisions

→**GPU usage is mandatory for sync reconstruction and calibration**

➢ All GPU software written in a generic way
➢ Same software runs on GPUs of different vendors and on the CPU

# ALICE reconstruction

## SYNC

- Rough corrections/calibrations for all detectors
- Full reconstruction of TPC (data reduction on GPU + space distortion corrections)
- TPC – ITS tracks matching (for a small subsample)
- Tracks propagation to outer detectors (TRD, TOF)
- Global track fits
- Primary and secondary vertices
- PID hypothesis
- **CTF and calibrations as output**

When the EPN farm is not (fully) used for synch. processing, it will be used for asynch. processing of the raw data stored on the disk buffer
**EPN will perform ~1/3 of the Pb-Pb asynchronous processing**

## ASYNC

- Full correction of TPC distortions (nominal resolution), full calibration for all detectors
- TPC – ITS tracks matching
- Tracks propagation to outer detectors (TRD, TOF)
- Global track fits
- Primary and secondary vertices
- PID hypothesis
- **Calibration/QC and AOD as output**

- ➢ Different relative importance of GPU / CPU algorithms compared to synchronous processing
- ➢ TPC part faster than in synchronous processing (less hits, no clustering, no compression

7

# Reconstruction time covered by GPUs

| Synchronous processing (50 kHz Pb-Pb, MC data, processing only) | |
| --- | --- |
| Processing step | % of time |
| TPC Processing (Tracking, Clustering, Compression) | 99.37 % |
| EMCAL Processing | 0.20 % |
| ITS Processing (Clustering + Tracking) | 0.10 % |
| TPC Entropy Encoder | 0.10 % |
| ITS-TPC Matching | 0.09 % |
| MFT Processing | 0.02 % |
| TOF Processing | 0.01 % |
| TOF Global Matching | 0.01 % |
| PHOS / CPV Entropy Coder | 0.01 % |
| ITS Entropy Coder | 0.01 % |
| Rest | 0.08 % |

| Asynchronous processing (650 kHz pp, real data, calorimeters not in run) | |
| --- | --- |
| Processing step | % of time |
| TPC Processing (Tracking) | 61.41 % |
| ITS TPC Matching | 6.13 % |
| MCH Clusterization | 6.13 % |
| TPC Entropy Decoder | 4.65 % |
| ITS Tracking | 4.16 % |
| TOF Matching | 4.12 % |
| TRD Tracking | 3.95 % |
| MCH Tracking | 2.02 % |
| AOD Production | 0.88 % |
| Quality Control | 4.00 % |
| Rest | 2.32 % |

**Running on GPU in baseline scenario** **Running on GPU in optimistic scenario**

Credits D. Rohr, CHEP 2023

# Reconstruction time covered by GPUs

Max ~6.5x speedup, since 85% of the compute power is in the GPU

- **Today**, offloading the ~60% of the async to the GPU should yield a speedup around **2.5x**.
- We remove 60% of the CPU time, while we are still CPU-bound, but we have some overhead CPU resources for driving the 8 GPUs.
- In the o**ptimistic scenario**, by offloading 80% we might get close to **5x**.

**Asynchronous processing**
**(650 kHz pp, real data, calorimeters not in run)**

| Processing step | % of time |
|---|---|
| TPC Processing (Tracking) | 61.41 % |
| ITS TPC Matching | 6.13 % |
| MCH Clusterization | 6.13 % |
| TPC Entropy Decoder | 4.65 % |
| ITS Tracking | 4.16 % |
| TOF Matching | 4.12 % |
| TRD Tracking | 3.95 % |
| MCH Tracking | 2.02 % |
| AOD Production | 0.88 % |
| Quality Control | 4.00 % |
| Rest | 2.32 % |

**Running on GPU in baseline scenario**    **Running on GPU in optimistic scenario**

Credits D. Rohr, CHEP 2023

# GPU workflow topology



Credits G. Eulisse and D. Rohr, CHEP 2023

# Production activity on EPN (GPU)

Full reprocessing of 2022 HIR pp@13.6 TeV data profited of EPN node availability

1/3 of the processing on EPN (conf 2 NUMA domains but only 1 GPU+64 cores per domain).

Next processing on EPN is now configured to run with 2 FULL NUMA domains (4 GPU, 64 cores per domain): 1 EPN job/domain expected to replace (today) ~16-20 GRID/CPU 8-core jobs



apass3
2022 data

In parallel with
~3300 8-core
GRID jobs

11

# Test on ARM architecture

We joined tests on an ARM node provided by E4 (one test node with 256 cores, 227 GB ram) in the context of Spoke-2

Requirements:

- cvmfs and local software installation
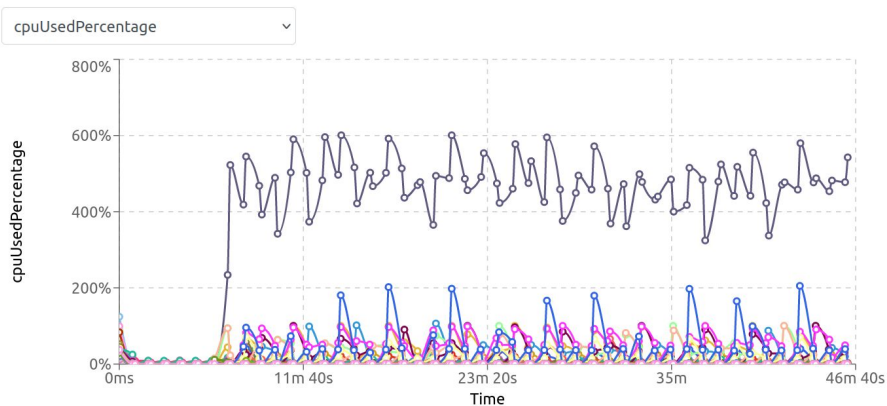- Network access to external ALICE CCDB (ALICE token auth. → grid cert)
- CentosOS

Test full workflows for

- Data reconstruction (GRID configuration)
- Full simulation

NB|| build on cvmfs for arm were sporadically provided (now built by default for all official production tag)
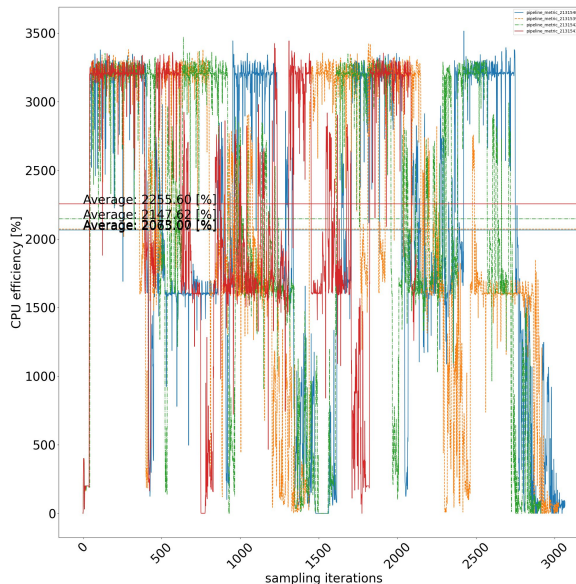
# Test on ARM architecture (II)

Data reconstruction



Simulation



4 concurrent 16 workers MC job (actually we fill half of the node in this test)

- GRID setting tuned for 8 cores per node
- CPU efficiency consistent with what observed in the GRID node
- Physics results consistent with the one of GRID jobs (unfortunately we cannot test with exactly the same software tag)
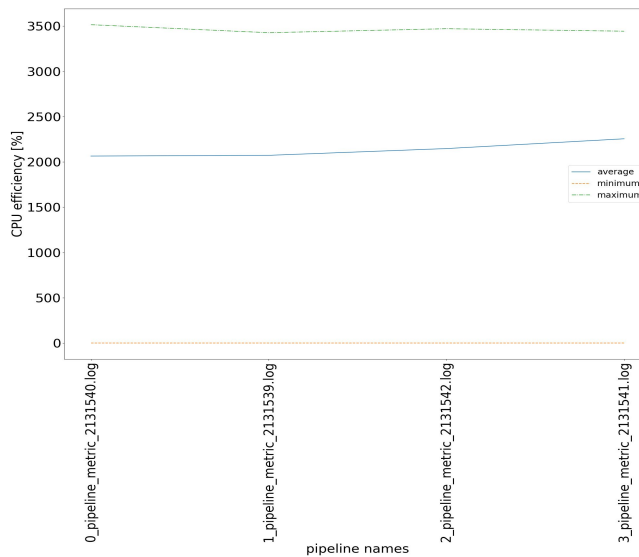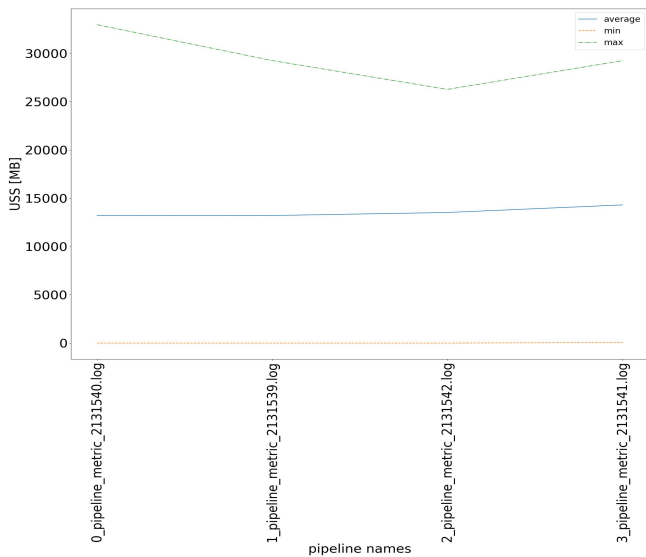
NB|| CPU usage improved

MC workflow is more aggressive than data reconstruction (cpu overbooking). If more CPU are available it will use them (up to a factor 2 more). RAM usage is the hard limit (2GB/core)

13

# Test on ARM architecture (III)

PbPb 10TFx50ev=500ev x 4job = 2000ev
export NWORKERS=16
export CPULIMIT=16
export MEMLIMIT=32000
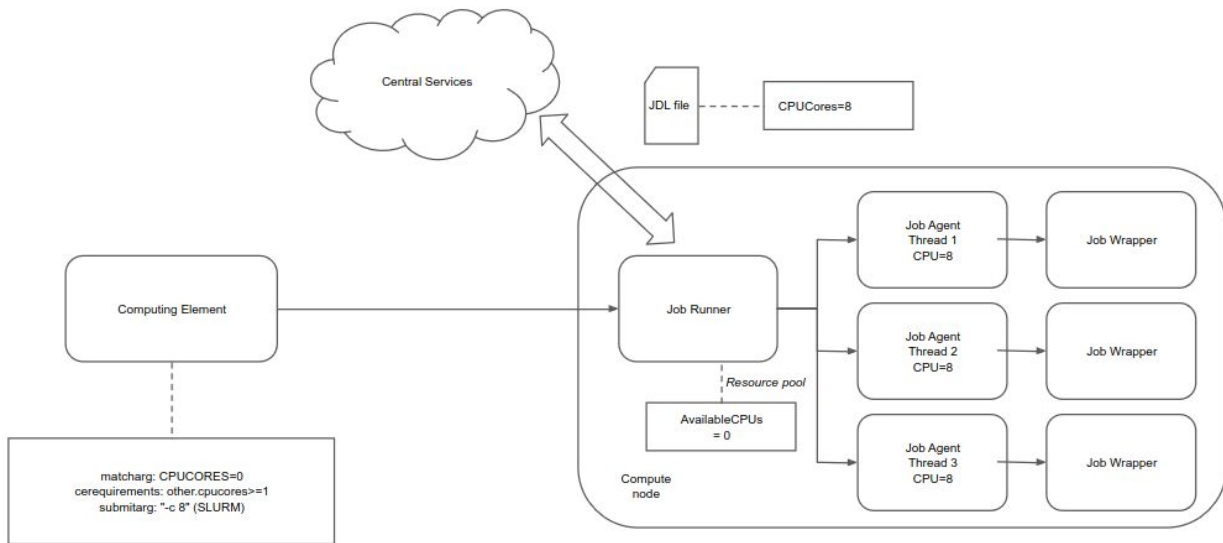→ Tot 64 cpu, 128 GB mem

Simulation



4 concurrent 16 workers MC job (actually we fill half of the node in this test)

MC workflow is more aggressive than data reconstruction (cpu overbooking). If more CPU are available it will use them (up to a factor 2 more). RAM usage is the hard limit (2GB/core)

14

# Whole node configuration



**Whole-node Scheduling Run Through**

Credits S. Weisz , ALICE Tiers workshop 2022

A new design of the Job runner in the worker node allows to instantiate more Job Agents (tasks) in a single slots and assign to them part of the resources available (e.g. in a 8-core slot we can run 8 single-core tasks)

This level of flexibility allows us to deal with a whole node configuration (full node reserved to an ALICE job)

- ideally for HPC resources
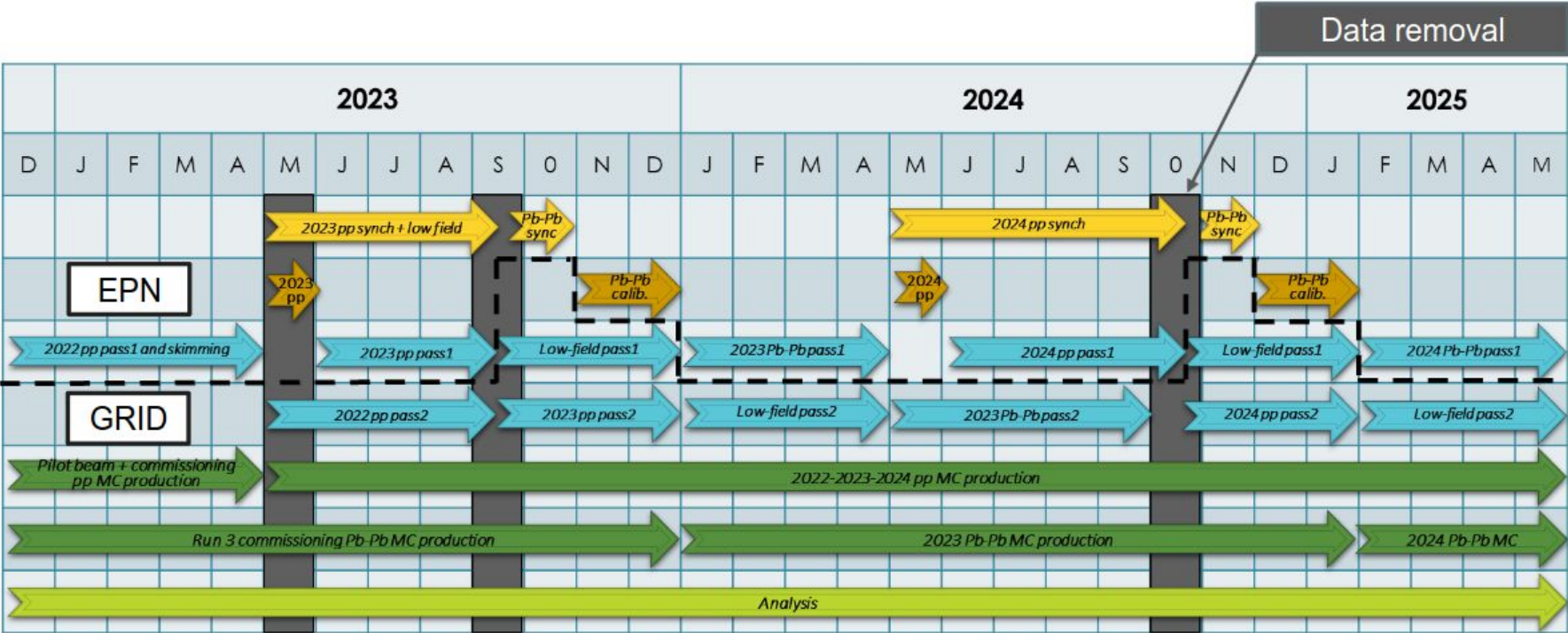- potentially also in other scenarios

Thanks for your attention!

Backup

# 2023-2024 processing timeline

# Generic software and GPU Benchmarks

Vendor/architecture-independent software:
- All algorithms are written in generic C++, and can be dispatched to HIP, CUDA, OpenCL on GPUs or OpenMP on CPUs using small wrappers → good code maintainability
- GPU libraries linked dynamically on demand → can distribute same binary software to CPU and GPU nodes

Benchmarking of the synchronous software completed in August 2020:
- GPU performance @ 50kHz Pb-Pb
  - ~1600 AMD MI50 and ~1100 NVIDIA Quadro RTX 6000
  - Compatible with our previous estimates <2000 GPU including 20% margin
- GPU Memory optimization
  - 128 orbit TF (~ 11 ms) needs 24 GB
- EPN Full System Tests performed with 70 orbit TF
  - Validated processing rate of 1/230 of assumed rate at 50 kHz Pb-Pb (nominal 1/250)
  - Max. server memory consumption 280 GB and CPU load 44 cores (+20% in the final setup)