



Role of GPU in ALICE Online-offline reconstruction at LHC Run 3

M. Concas, D. Elia, <u>F. Noferini</u>, S. Piano, M. Puccio

Workshop sul Calcolo nell'I.N.F.N.Paestum23 - 27 maggio 2022



ALICE @ Run 3 and 4



ALICE @ Run 3 and 4



Hardware commissioning/magnet training



Continuous Readout

Pb-Pb@50 kHz pp@1 MHz

Run 3/4: LHC will deliver 50 kHz in Pb-Pb collisions

 ALICE aims to record >10 nb⁻¹ integrated luminosity, x50 times more minimum bias data wrt Run 2



ALICE Raw Data Flow in Run 3



ALICE O² computing model



2/3s of CTFs processed by $O^2 + T0$

1/3 of CTFs exported, archived and

and archived at T0;

processed on T1s;



Processing plan for Pb-Pb





ALICE O² activities

- ALICE upgrade to continuous readout required a new Online-Offline (O²) framework
 - message-queues based processing by separate device(processes)
 - Data Processing Layer (DPL)
 - Workflows are built for group of Devices by automatic matching of their Inputs and Outputs

G. Eulisse, R. Shaoyan (CHEP 2019)





Using GPU in Run 3 (reconstruction) with ALICE O²

Processing on dedicated farm at experimental site

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

\rightarrow 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



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
 - \circ ~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)



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



Reconstruction steps for GPU-offload







Reconstruction time covered by GPUs

room of improvements in async reconstruction

Synchronous processing		Asynchronous processing		
Processing step	% of time	Processing step		% of time
TPC Processing	99.37 %	TPC Processing		72.01 %
EMCAL Processing	0.20 %	TRD Tracking		12.69 %
ITS Processing	0.10 %	TOF-TPC Matching		9.94 %
TPC Entropy Coder	0.10 %	MFT Tracking		1.69 %
ITS-TPC Matching	0.09 %	ITS Tracking		0.78 %
MFT Processing	0.02 %	TPC Entropy Decoder		0.73 %
TOF Processing	0.01 %	Secondary Vertexing		0.69 %
TOF Global Matching	0.01 %	ITS-TPC Matching		0.56 %
PHOS / CPV Entropy Coder	0.01 %	Primary Vertexing		0.14 %
ITS Entropy Coder	0.01 %	TOF Global Matching		0.11 %
FIT Entropy Coder	0.01 %	PHOS / CPV Entropy Decoder		0.10 %
TOF Entropy Coder	0.01 %	FIT Entropy Decoder		0.10 %
MFT Entropy Coder	0.01 %	ITS Entropy Decoder		0.06 %
TPC Calibration residual extraction	0.01 %	TOF Entropy Decoder		0.05 %
TOF Processing	0.01 %	MFT Entropy Decoder		0.05 %
Bunning on GPU in baseline scenario	Running on GPU in baseline scenario Running on GPU in optimistic scenario Preliminar			ome algorithms

not yet complete or not optimized!



Reconstruction time of O2 GPU Synchronous Reconstruction on





Speedup of O2 GPU synchronous reconstruction versus CPU





Number of GPUs required to run O2 GPU synchronous reconstruction at 50 kHz Pb-Pb





HPC and cloud resources (simulation)

- Thanks to the new O² simulation and reconstruction code (Run 3) possible to fully exploit the multi process features
- Significant progress has been made to incorporate HPC and cloud resources in the standard ALICE Grid workflows
 - Multicore queues at CERN used to test and benchmark the O² MC code
 - Intel based HPCs (Marconi @ CINECA, Cori and Lawrencium @ LBNL) were used for the O² MC challenge
 - Cloud resources delivery at CERN direct integration of Azure cloud as a Grid node
- Next steps:
 - porting the O² code to Power 9 and ARM platform



O² simulation

We definitely entered in the Run 3 phase: O² simulations replaced Run 2 ones in the GRID.

We recently simulate 1 billion pp@13.6 TeV events. Even if we are still collecting metrics the improvement with respect to Run 2 data was confirmed to be better than a factor $3 \rightarrow$ impact on resources per simulated event is lower than Run 2 case both for Wall Time and disk usage.

Additional optimizations in place to speed up simulations: e.g. embedding in digitization \rightarrow <u>S. Wenzel</u> (CHEP 2018)

Usage of GPUs not feasible \rightarrow Geant4 transportation code is still the dominant component



Summary

ALICE will record 50 kHz Pb-Pb minimum bias data in Run 3 without trigger.

• Continuous TPC readout, time frames of ~10 (or ~20) ms instead of events.

Full online data processing on GPUs.

- Computing farm consists of 250 servers, with 8 AMD MI50 GPUs, 2 32-core Rome CPUs, and 512 GB RAM each.
- Currently 230 servers are sufficient for processing 50 kHz Pb-Pb (peak load).
- MI50 GPU replaces ~80 CPU cores (sync reco) or ~55 CPU cores (async reco)

All GPU software written in generic way, can run on different GPUs and on the CPU.

Processing farm used for synchronous (online) and asynchronous processing (periods without beam).

- Full baseline scenario with synchronous GPU processing ready.
- Planning to use GPUs as much as possible also in asynchronous processing.
 - In the optimistic scenario, we will be able to offload ~95% of the workload to the GPU.

Usage of GPUs in GRID nodes not yet explored (reco is performed at Tier-1s), is it something foreseen in INFN plans?

backup



Data buffer for Online-Offline (O²) facility

- 100 PB raw capacity, RS 10+2 erasure coded (level of security to be defined, 84 PB usable space)
- Based on cheap JBODs, SATA drives, EOS managed



20



Generic software and GPU Benchmarks



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)



Reconstruction time of O2 GPU Synchronous Reconstruction on

GPU and CPU



Runs only EPN servers:

- 1x Supermicro server, 512GB of memory
- 2x 32cores Rome CPUs
- 8x AMD MI50 GPUs

Standalone tests for GPU performance:

- Time linearly scales with the size of processed data
- GPUs are much faster than CPUs