



Energy Efficient Computing with ALICE O²

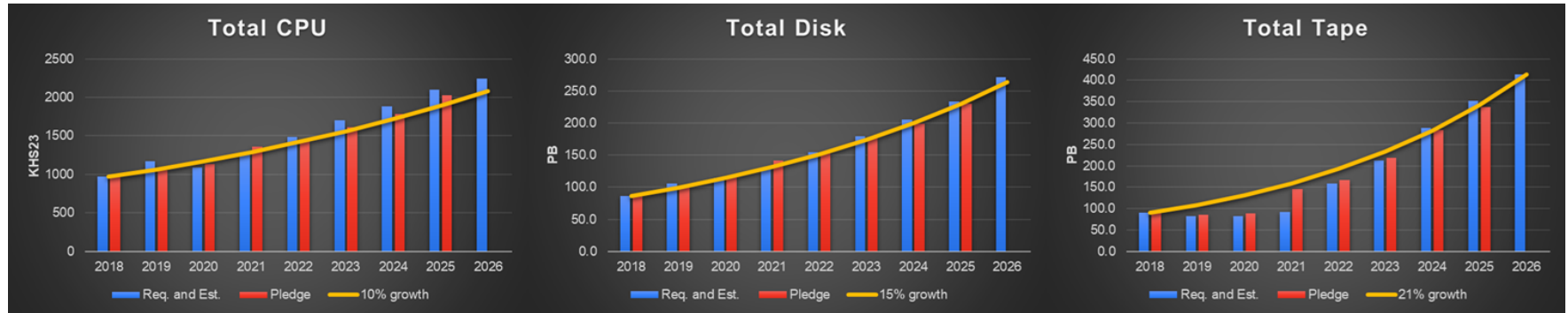
Latchezar Betev (CERN), Irakli Chakaberia (LBNL),
Federico Ronchetti (CERN – INFN LNF),
Stefano Piano (INFN TS)

The ALICE computational challenge in Run 3/4+

- Run 1/2 - triggered acquisition @ 7-10 kHz IR, down to 1 kHz for central part of the detector
- Run 3/4 - continuous readout acquisition @ 50 kHz for Pb-Pb, 650 - 1000 kHz for p-p
 - Enhanced study of probes with small S/B for which triggering is not possible, as charm and beauty
 - Lossless data compression
 - Data skimming for p-p to 4% of original volume
- **The Challenge** - cope with the data rates without dramatically increasing the computational requirements
 - Entirely new software
 - Special compression facility with extended use of accelerators

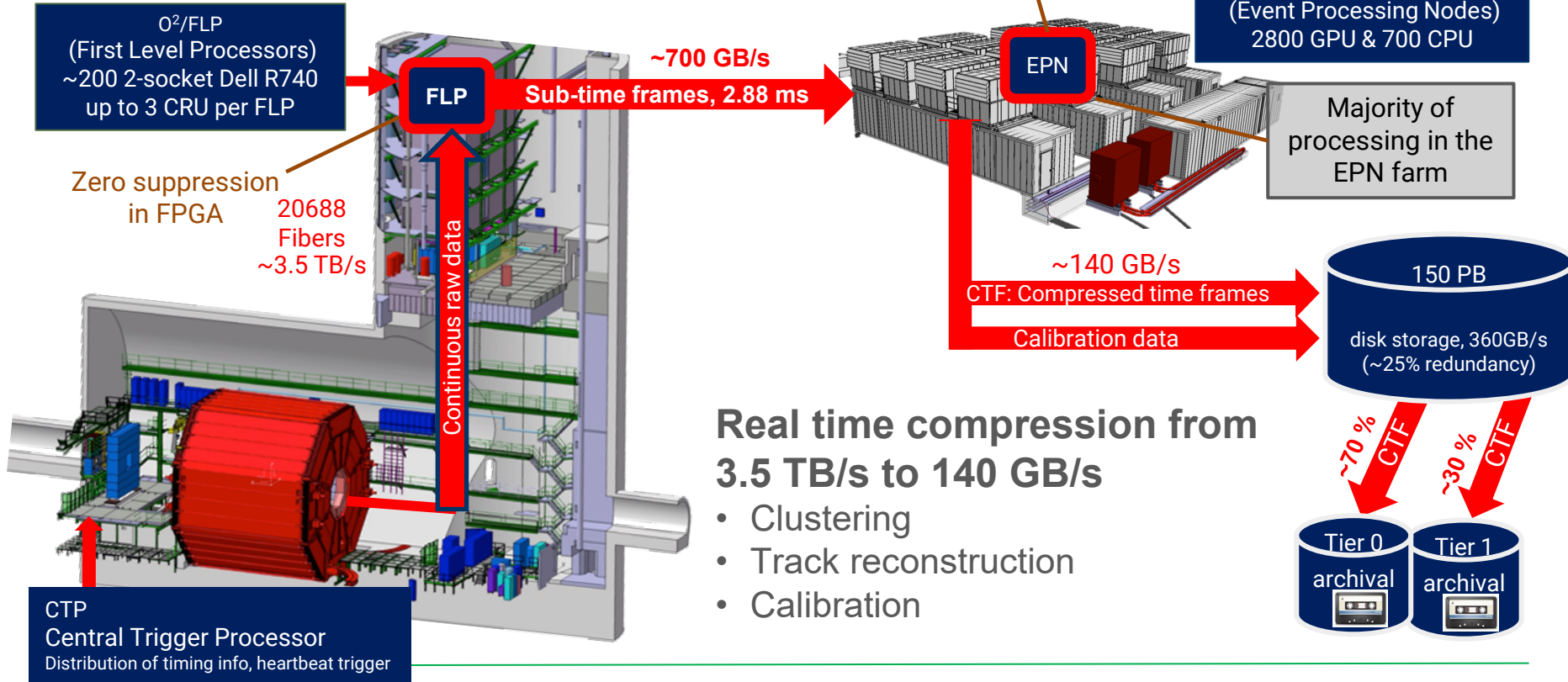
The Offline-Online framework O²

- Online and offline processing with a single framework (O²) in two stages
 - **Synchronous** - data compression, calibration, and QC on dedicated farm close to the detector
 - **Asynchronous** - final reconstruction output (multiple passes)
- Collect 100x more events (compared to Run 1/2), with unchanged resources profile growth @~15% year





Data compression



The Event Processing Node (EPN) farm

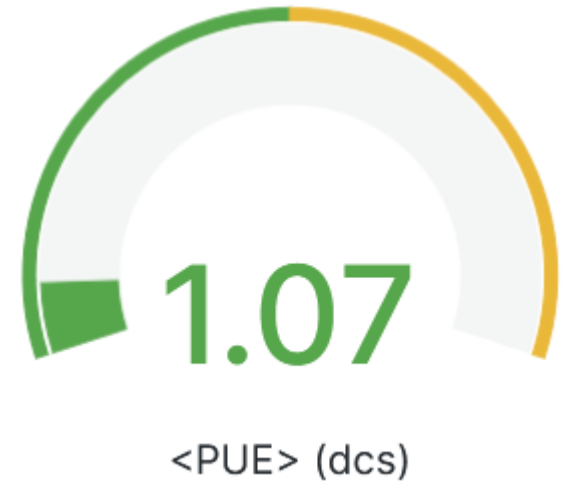


O²/EPN: 350 nodes with
2800 GPU & 700 CPU

- **Modular containers** → allow easy extensions of the farm
- **Adiabatic cooling** when temperature is higher than the setpoint, otherwise air-air cooling
- Build and operated by the EPN team
 - Preventive and second level maintenance from CERN and ALICE Technical Coordination

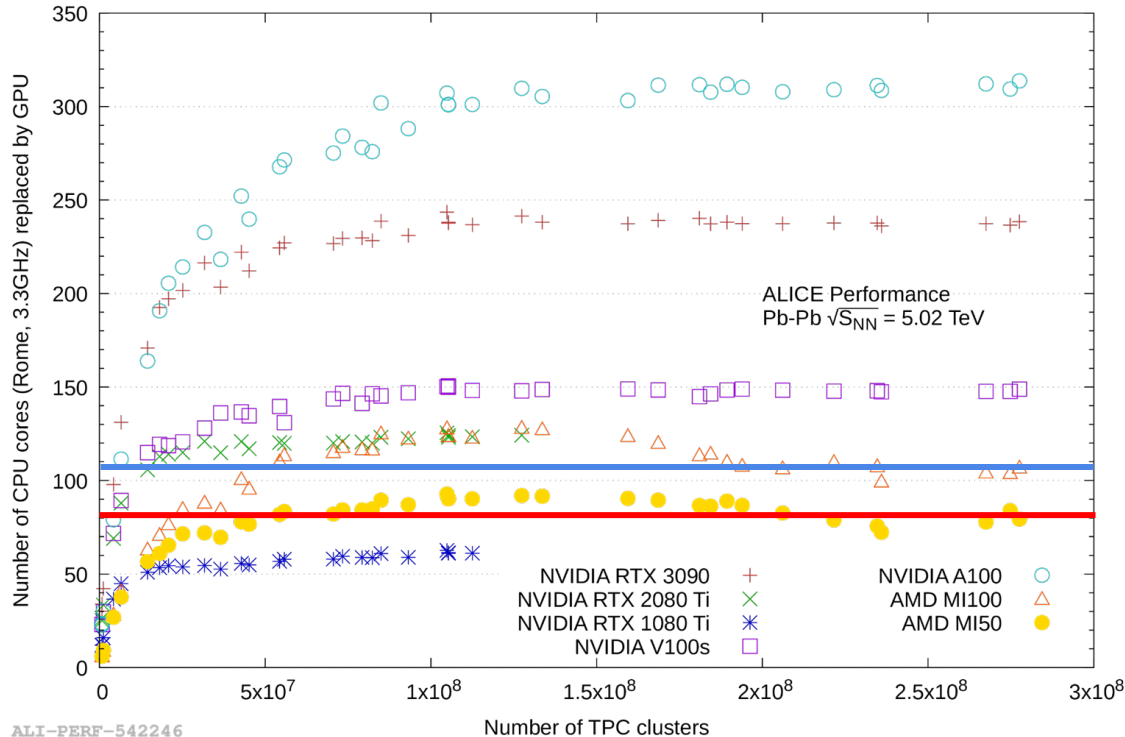
Power Usage Efficiency (PUE) of EPN infrastructure

- Cooling and power design => low PUE < 1.1





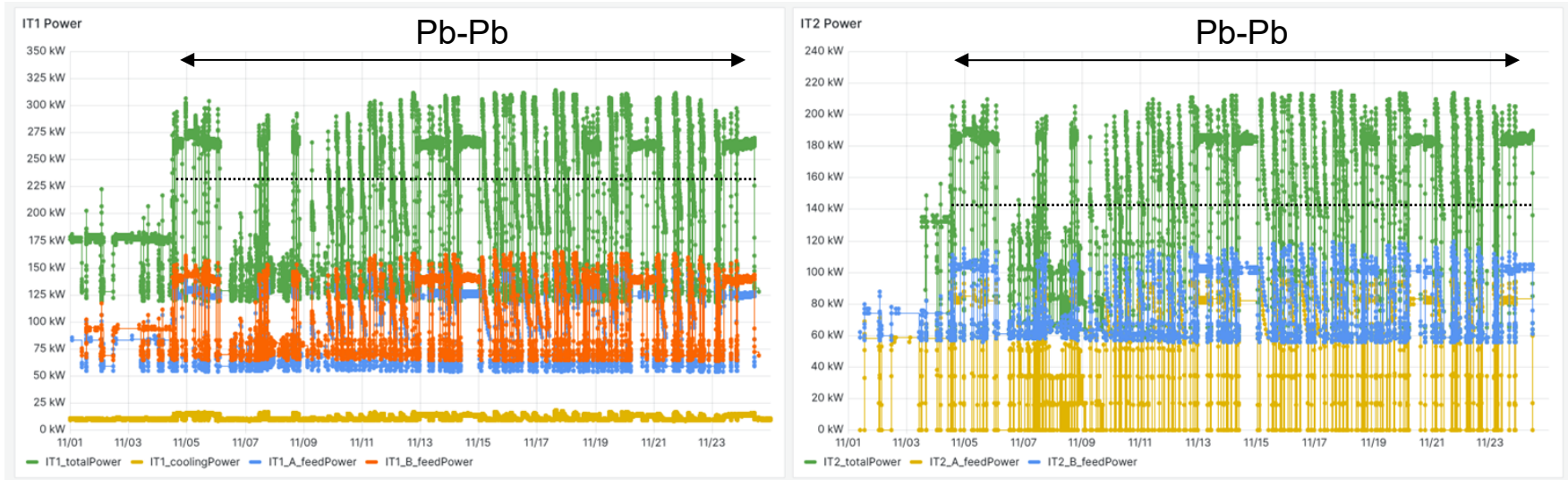
Lower cost and energy consumption - replacing CPUs with GPUs



ALI-PERF-542246

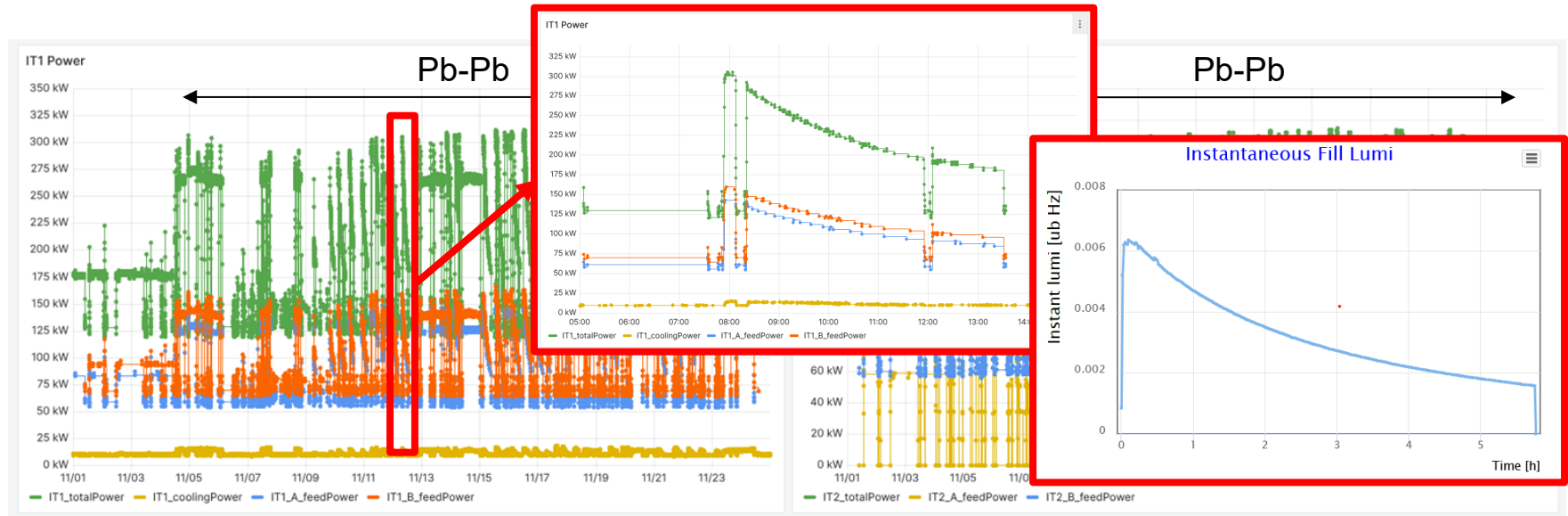
- Cost and energy efficient replacement
 - 1 GPU can replace 50-300 CPU cores @ 3.3 GHz for ALICE TPC online processing
 - One **MI50** (**MI100**) equivalent to ~**80** (**107**) CPU cores
- Compact form factor - 8 GPUs per chassis, instead of 2 CPUs

EPN power consumption during Pb-Pb data taking



- Average power consumption of **380 kW** during Pb-Pb at ~50 kHz

EPN power consumption during Pb-Pb data taking



- Average power consumption of **380 kW** during Pb-Pb at ~ 50 kHz
- Power profile aligns well with the luminosity delivered by LHC

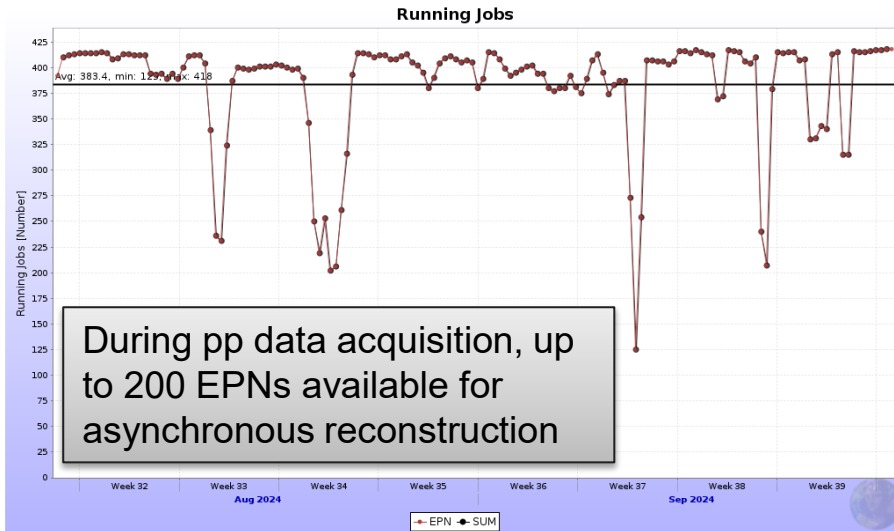
Power reduction with GPUs

	EPN nodes	CPU Cores	GPUs	Equivalent CPU cores per GPU	Equivalent CPU cores	Total needed cores	Equivalent Nodes 2 CPUs
EPN MI50	280	17'920	2'240	80	179'200	197'120	3'080
EPN MI100	70	6'720	560	107	59'733	66'453	692
TOTAL	350	24'640	2'800		238'933	263'573	3'772

- **380 kW** average power consumption of 350 node EPN farm
- More than 3'700 nodes needed if only CPU used
- Assuming average of 200 W/node => **740 kW** for CPU-only farm (*)
 - Almost 100% more power

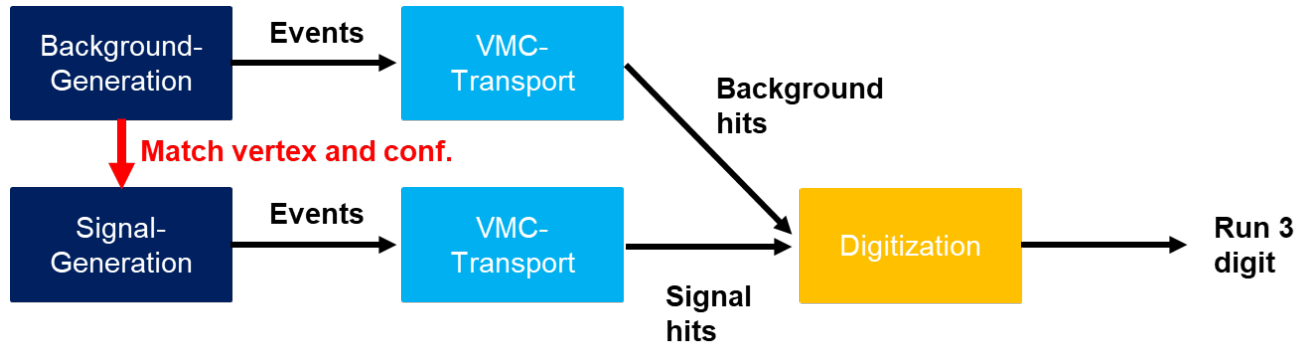
(*) detailed measurements will be published soon

Asynchronous processing on GPUs



- EPN farm use is balanced between synchronous and asynchronous processing according to online load
 - Allocation of the partitioning - managed by Run Coordinator
- 60% of asynchronous reconstruction is done on GPU: current **speedup of 2.5x**
- Ongoing efforts to offload up to 80% of CPU time to GPU, aiming to achieve a speedup of 5x

O² simulation ecosystem: embedding in digitization



- ALICE can leverage embedding techniques, i.e. generating TPC pile up events, to accelerate simulation processes, early benchmark on test machine:

Description	Collision System	IR	Number of Events	Without embedding		With embedding		
				HS06 s / event	Speedup wrt Run 2	Expected gain	HS06 s / event	Speedup wrt Run 2
min bias	pp	500kHz	10000	184	5.2	1.5	123	7.8
min bias	Pb-Pb	50kHz	250	10778	2.2	3	3593	6.7

- The total improvement considering the embedding gain is up to 7-8 times the Run 2 simulation



O² Analysis Model

Number of input files
Input size
Output size
PSS Memory
Private Memory
Timing



- ALICE has a strong tradition of organized analysis
 - Saves resources, as users run together over the same data
- Novel O² analysis framework **10x** faster analyses wrt Run 2
 - New AO2D data formats **9x (Pb-Pb) to 25x (pp)** less storage space than Run 2 ESD and AOD formats
 - Smaller formats for derived data samples for specific analyses in Run 3
- **Extensive use of Analysis Facilities with lower energy footprints: Green Cube @ GSI, HPC systems @ LBNL, Wigner Scientific Computing Laboratory**

Estimated emissions 🌱

This is an estimate of the 'CO2 equivalent' which would be produced by running. The estimated 1053 kg produced would be greater than flying from Geneva to New York .

To find out more about these estimates, take a look at the Hyperloop documentation, accessible from the ? button in the page header.

The analysis train test includes all relevant parameters required for validation, along with an average carbon estimate to enhance user awareness

Emerging CPU architectures

- Alternative and more energy efficient architectures are gaining significant traction
- ARM CPU @Glasgow provided temporarily for ALICE
 - Dual Ampere Altra Q80-30 CPU's with ~160 cores, 3.2GB RAM/core, 16k cores total, 11HS23/core
 - Performance suggests different memory management, beneficial to our code
 - Further tests ongoing @ CNAF aiming at full physics validation
- A very good option for Grid CPU, if cost-effective
- Aarch64 fully supported
 - Automatic matching of binaries and containers for ARM WNs
- Changes kept as generic as possible
 - Allows us to easily slot-in support for more architectures in future (e.g. RISC-V)



Green Summary

- The ALICE computing facilities and software model were upgraded to accommodate enhanced detector performance and higher data rates in Run 3
- Significant energy and CO₂ savings were achieved by ALICE through
 - Algorithmic optimization, parallelizing and vectorizing CPU code
 - Maximizing GPU usage in purpose-built facilities with low PUE - the EPN
- Targeting efficient PUE centers allows for optimization of energy consumption and reduction of the carbon footprint of the analysis
- Reusing generated events with embedding techniques accelerates simulations and reduces computing cost
- More energy-efficient architectures are now fully supported by ALICE



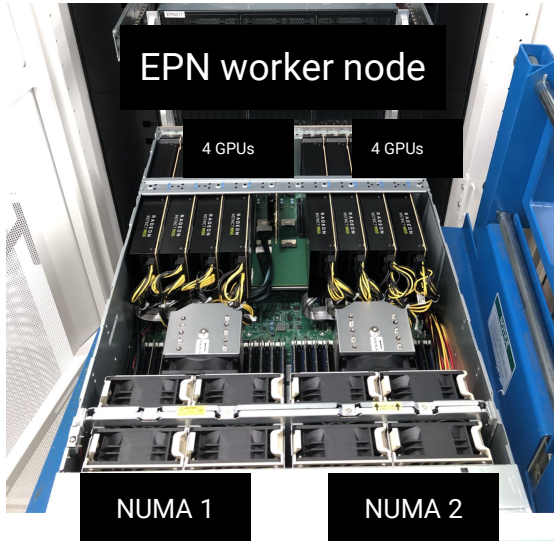
Backup



The ALICE EPN FARM

- 280 nodes equipped with 8 AMD MI50 32GB GPUs
- Additional 70 nodes (installed in 2023) equipped MI100 32GB
- Grand total of 350 (280 MI50 + 70 MI100) nodes and 2800 GPUs (equivalent to ~373 MI50 nodes at MI100 = 4/3 MI50)

	70 MI100 EPNs	280 MI50 EPNs	4 Calib Nodes
GPU	8 AMD Instinct™ MI100 32 GB	8 AMD Instinct™ MI50 32 GB	
CPU	2 AMD EPYC™ 7552, 48 cores	2 AMD EPYC™ 7452, 32 cores	2 AMD EPYC™ 7452, 32 cores
MEMORY	1TB DDR4 3200 MHz	512GB DDR4 3200 MHz	512GB DDR4 3200 MHz
Networks	IB 100 Gb/s, ETH 1 Gb/s		



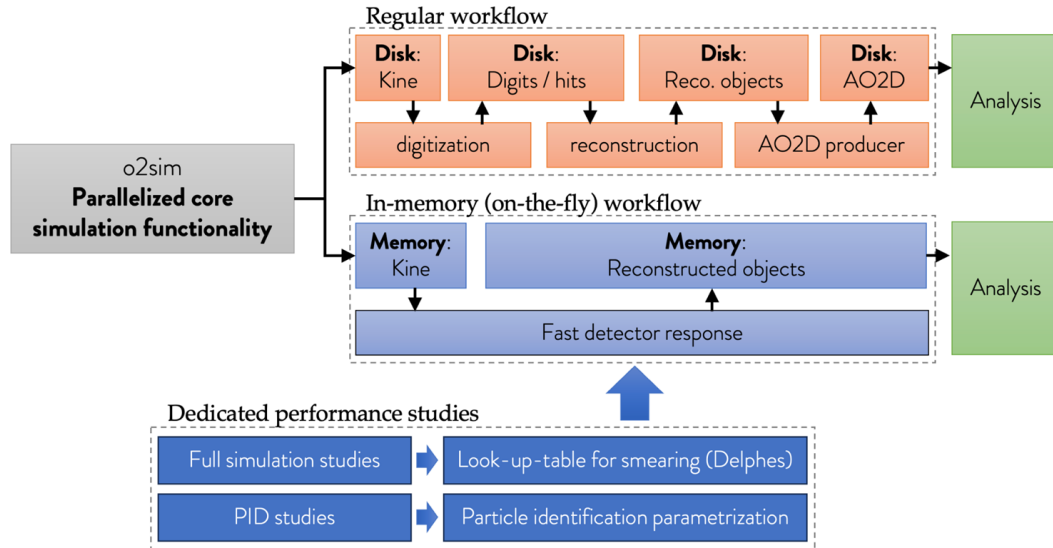
PERFORMANCE	MI-100
Compute Units	120
Stream Processors	7,680
Peak BFLOAT16	Up to 92.3 TFLOPS
Peak INT4 INT8	Up to 184.6 TOPS
Peak FP16	Up to 184.6 TFLOPS
Peak FP32 Matrix	Up to 46.1 TFLOPS
Peak FP32	Up to 23.1 TFLOPS
Peak FP64	Up to 11.5 TFLOPS
Bus Interface	PCIe® Gen 3 and Gen 4 Support ³

NODES	GPU FP32	TFLOPS (FP32)
280	13.3	3724
70	23.1	1617
350		5341

PERFORMANCE	MI-50
Compute Units	60
Stream Processors	3,840
Peak INT8	Up to 53.6 TOPS
Peak FP16	Up to 26.5 TFLOPS
Peak FP32	Up to 13.3 TFLOPS
Peak FP64	Up to 6.6 TFLOPS
Bus Interface	PCIe® Gen 3 and Gen 4 Supported ²



Fast, modular and flexible: 'on-the-fly simulations'



Regular simulation chain:

- Heavy to run
- Frequent use of scratch space on disk

Fast simulation chain:

- Parametrised: >>100x faster
- No use of disk at all
- Relies on separate **detector studies**
- Ideal for upgrade studies, since detector can be exchanged ~trivially
- Fully integrated into O2 framework
- Other use cases to be studied

- ✓ “FastTracker” tool now operational with “on-the-fly” MC generation: no content saved to disk
- ✓ State-of-the-art tracking: strangeness tracking included in fast tracking tool
- ✓ Maximum CPU efficiency, zero intermediate storage, enormous flexibility