



ALICE continuous readout and data reduction strategy for Run3 and 4

R.Shahoyan for the ALICE collaboration,
CHEP 2019, Adelaide, 05/11/2019



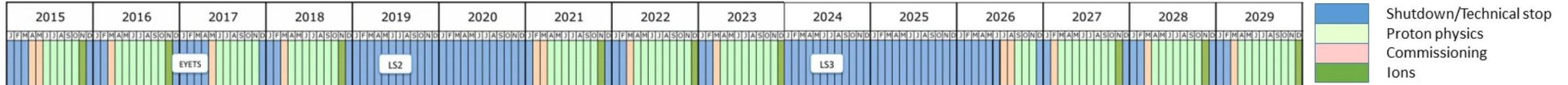
Executive Summary



- Aim of the hardware upgrade in Run3/4 is to boost amount of collected collisions by factor ~ 100
- Operating ALICE will require cardinal change in approach to data processing $\Rightarrow O^2$ framework
- Main challenges for handling the data to come:
 - The volume and the rate of data to store
 - \Rightarrow requires combination of data reduction and compression by factor >35
(~ 3.5 T/s \rightarrow <100 GB/s) \Rightarrow ~ 50 PB/y from Pb-Pb and reference pp data
 - Computing power to process these data
 - \Rightarrow requires drastic improvements in reconstruction code (factor ~ 10 already achieved wrt Run2)
and adapting heaviest operations to GPU processing
- Lot of other challenges (not covered in this talk), most important being coping with Space-Charge Distortions in the TPC (up to 15 cm for space-points with target precision of ~ 500 μm)



Run2 → Run3 and 4



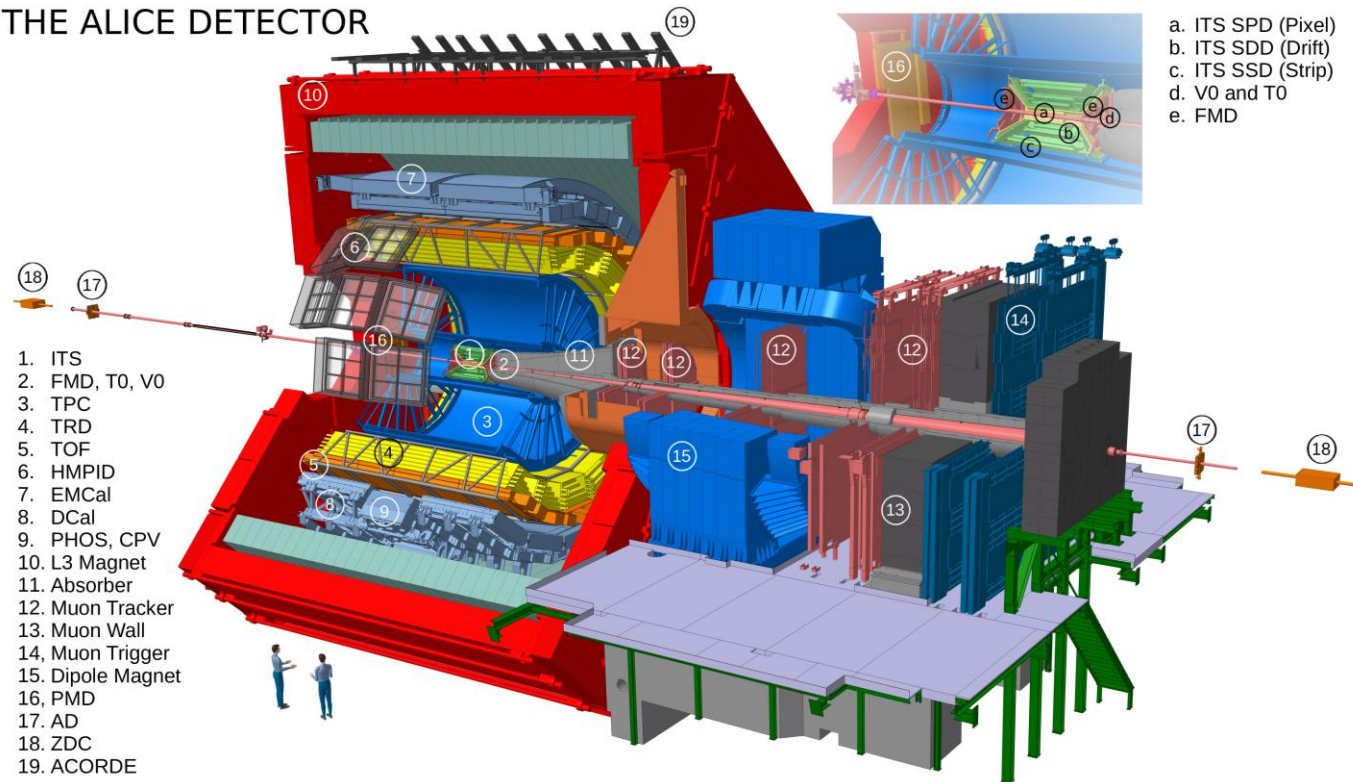
Run2: $\mathcal{L}_{int}^{Pb-Pb} = 1.0 \text{ nb}^{-1}$

Run3: $\mathcal{L}_{int}^{Pb-Pb} = 6.0 \text{ nb}^{-1}$

Run4: $\mathcal{L}_{int}^{Pb-Pb} = 7.0 \text{ nb}^{-1}$

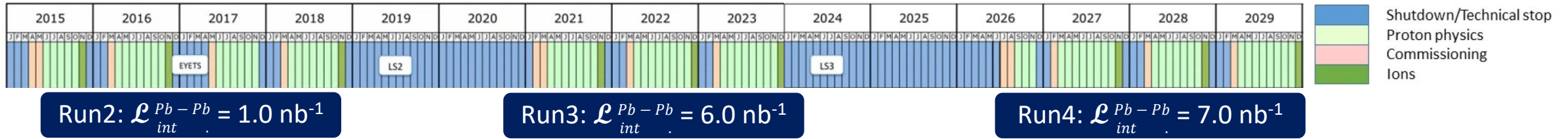
- In Run2 ALICE operated at Pb-Pb interaction rates $\sim 7\text{-}10 \text{ kHz}$ (inspected $\sim 1 \text{ nb}^{-1}$) with trigger rate $< 1 \text{ kHz}$

THE ALICE DETECTOR

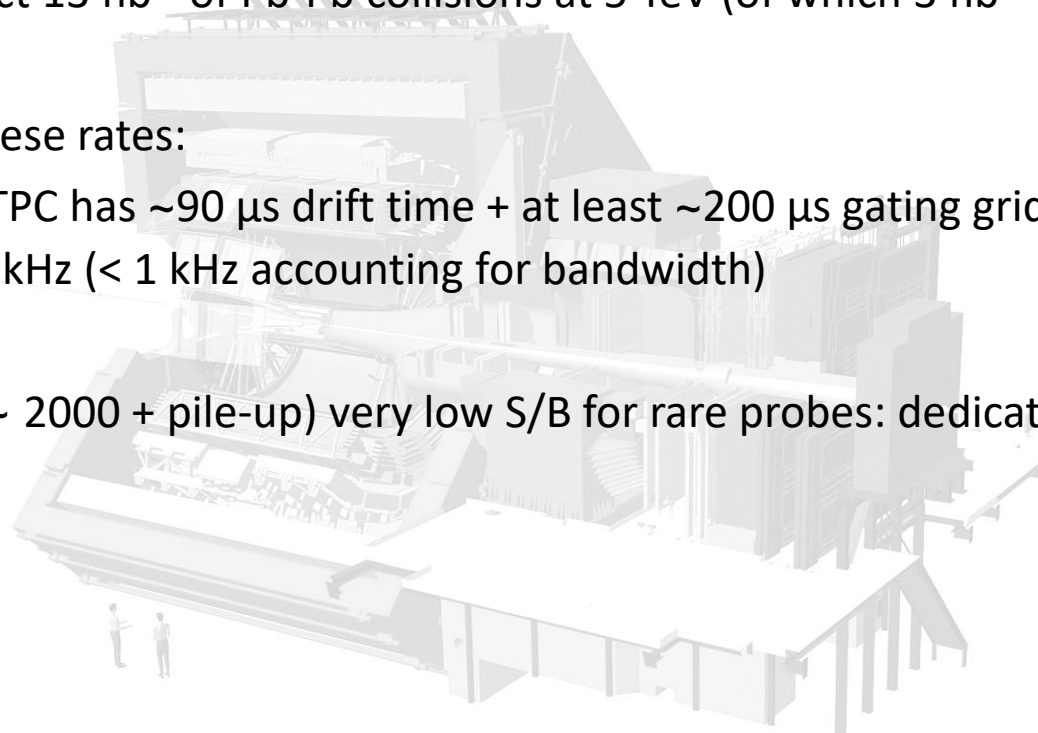




Run2 → Run3 and 4

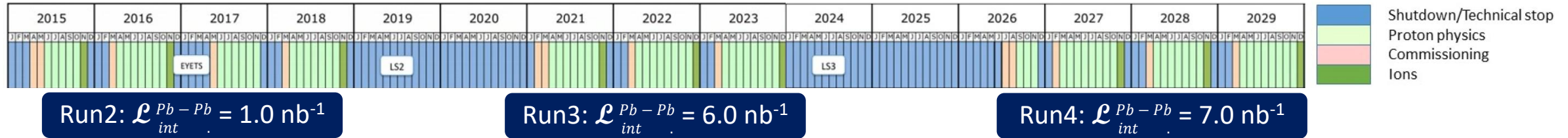


- In Run2 ALICE operated at Pb-Pb interaction rates $\sim 7\text{-}10 \text{ kHz}$ (inspected $\sim 1 \text{ nb}^{-1}$) with trigger rate $< 1 \text{ kHz}$
- LHC plans to deliver 50 kHz Pb-Pb interaction rate after LS2
- ALICE plans for Run3&4: collect 13 nb^{-1} of Pb-Pb collisions at 5 TeV (of which 3 nb^{-1} with reduced field)
- Main limitations to work at these rates:
 - Principal tracking detector, TPC has $\sim 90 \mu\text{s}$ drift time + at least $\sim 200 \mu\text{s}$ gating grid to collect the ion backflow → trigger rate limited to $\sim 3 \text{ kHz}$ ($< 1 \text{ kHz}$ accounting for bandwidth)
- At high multiplicities ($dN/d\eta \sim 2000 + \text{pile-up}$) very low S/B for rare probes: dedicated (HLT) trigger is not realistic





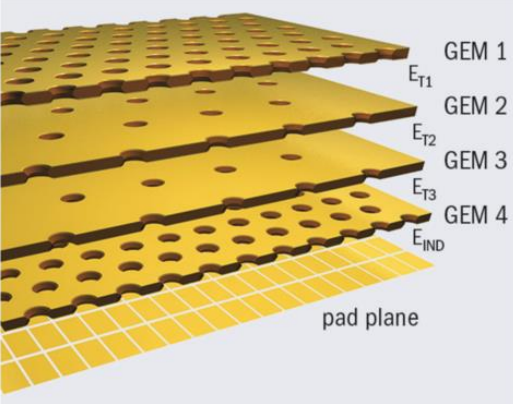
Run2 → Run3 and 4



- In Run2 ALICE operated at Pb-Pb interaction rates $\sim 7\text{-}10 \text{ kHz}$ (inspected $\sim 1 \text{ nb}^{-1}$) with trigger rate $< 1 \text{ kHz}$
 - LHC plans to deliver 50 kHz Pb-Pb interaction rate after LS2
 - ALICE plans for Run3&4: collect 13 nb^{-1} of Pb-Pb collisions at 5 TeV (of which 3 nb^{-1} with reduced field)
 - Main limitations to work at these rates:
 - Principal tracking detector, TPC has $\sim 90 \mu\text{s}$ drift time + at least $\sim 200 \mu\text{s}$ gating grid to collect the ion backflow → trigger rate limited to $\sim 3 \text{ kHz}$ ($< 1 \text{ kHz}$ accounting for bandwidth)
 - At high multiplicities ($dN/d\eta \sim 2000 + \text{pile-up}$) very low S/B for rare probes: dedicated (HLT) trigger is not realistic
- ➔ Use continuous readout at least for TPC (no gating grid), increase bandwidth
Read out all events, store compressed data and inspect all events offline



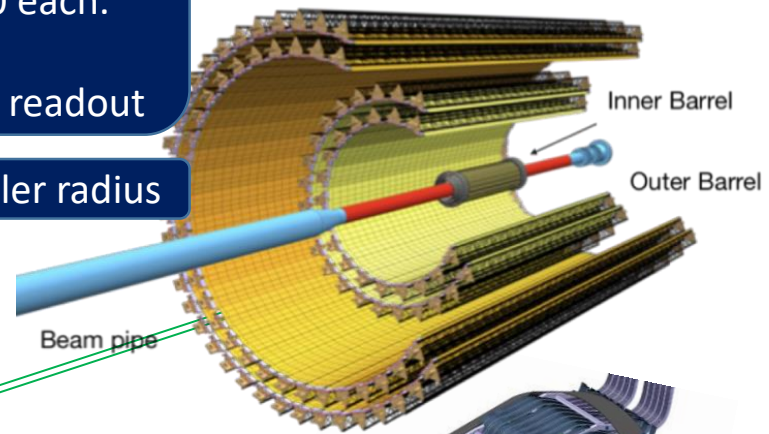
ALICE HW upgrades



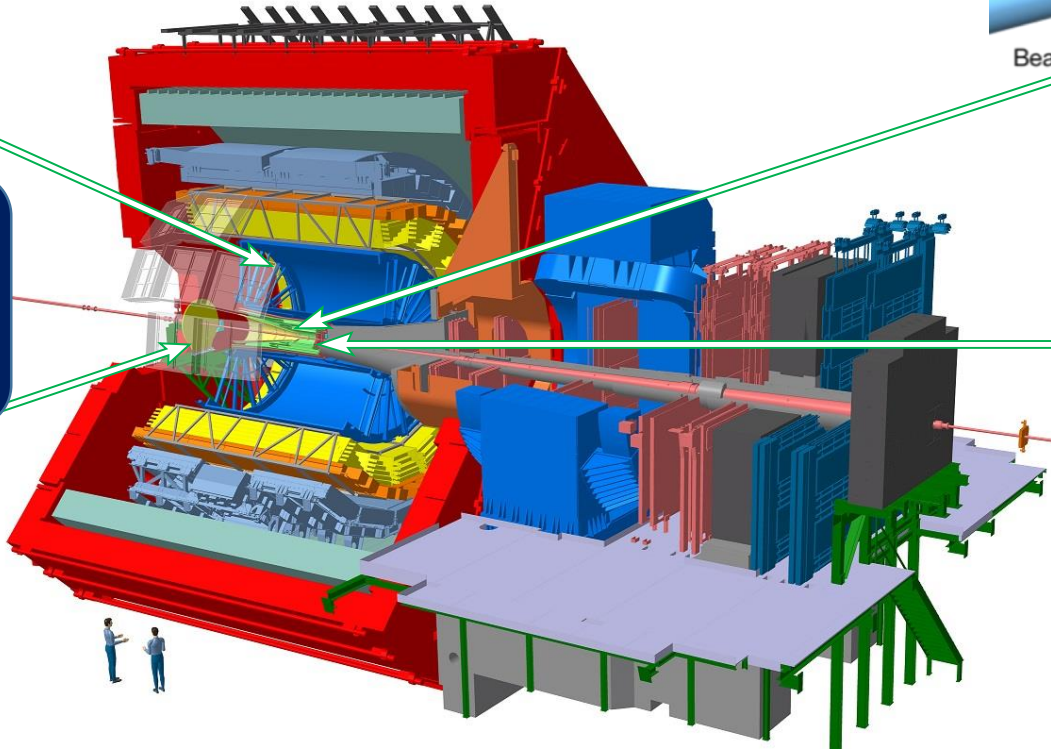
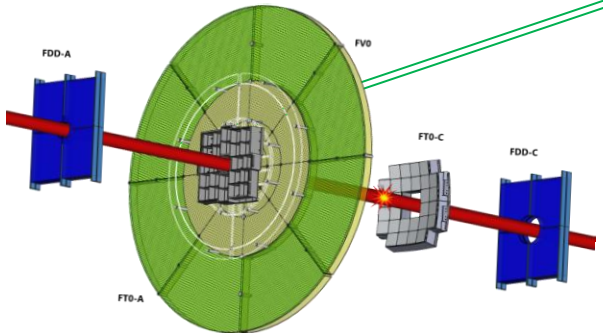
TPC MWPC readout → 4 layer GEM
(Intrinsic ion backflow ~99% blocking)
5MHz continuous sampling

New Si Inner Tracker: 10 m² of
MAPS with 29x27μm² pixel size
3 inner layers ~0.3% X0 each.
Closer to the beam
50-500 kHz continuous readout

New beam pipe of smaller radius



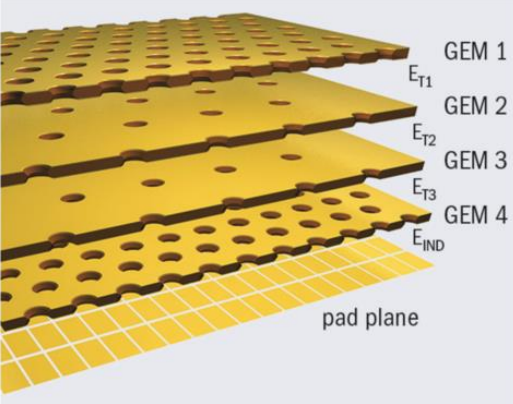
Fast Interaction Trigger (FIT) detector
Scintillator (FV0, FDD) + Cerenkov (FT0)
detectors to provide Min.Bias trigger
for detectors with triggered R/O



Muon Forward Tracker
to match muons before
and after the absorber.
Same Si chips as new ITS



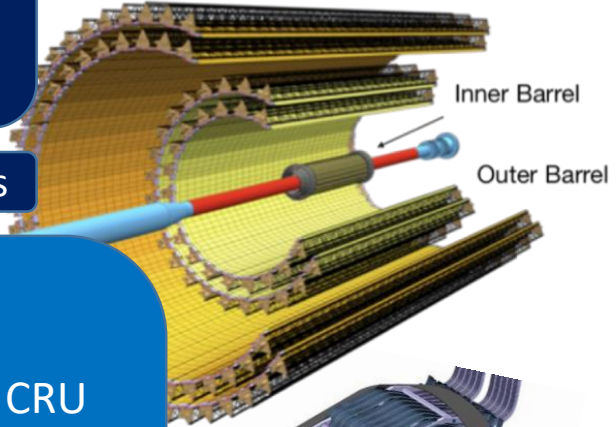
ALICE HW upgrades



TPC MWPC readout → 4 layer GEM
(Intrinsic ion backflow ~99% blocking)
5MHz continuous sampling

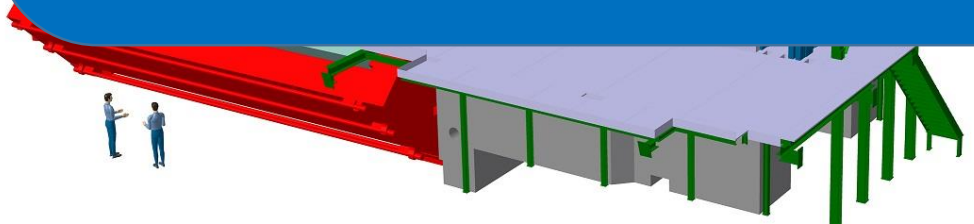
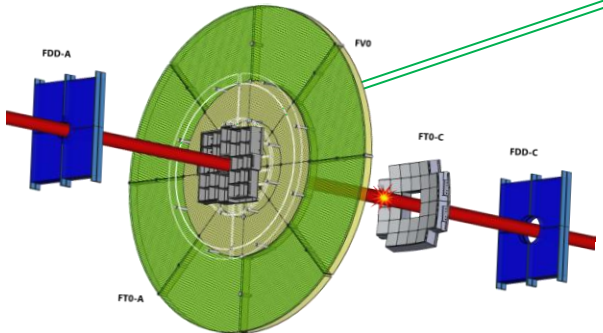
New Si Inner Tracker: 10 m² of
MAPS with 29x27μm² pixel size
3 inner layers ~0.3% X0 each.
Closer to the beam
50-500 kHz continuous readout

New beam pipe of smaller radius



Fast Interaction Trigger (FIT) detector
Scintillator (FV0, FDD) + Cerenkov (FT0)
detectors to provide Min.Bias trigger
for detectors with triggered R/O

New readout (all except EMCal, PHOS and HMPID) via CRU
(Common Readout Unit, PCIe40 /Arria10 FPGA/, developed by LHCb)
Detectors can be read out in continuous or triggered modes,
except triggered-only EMCal, PHOS/CPV, TRD (~40kHz)
and HMPID (2.5 kHz)



Muon Forward Tracker
to match muons before
and after the absorber.
Same Si chips as new ITS

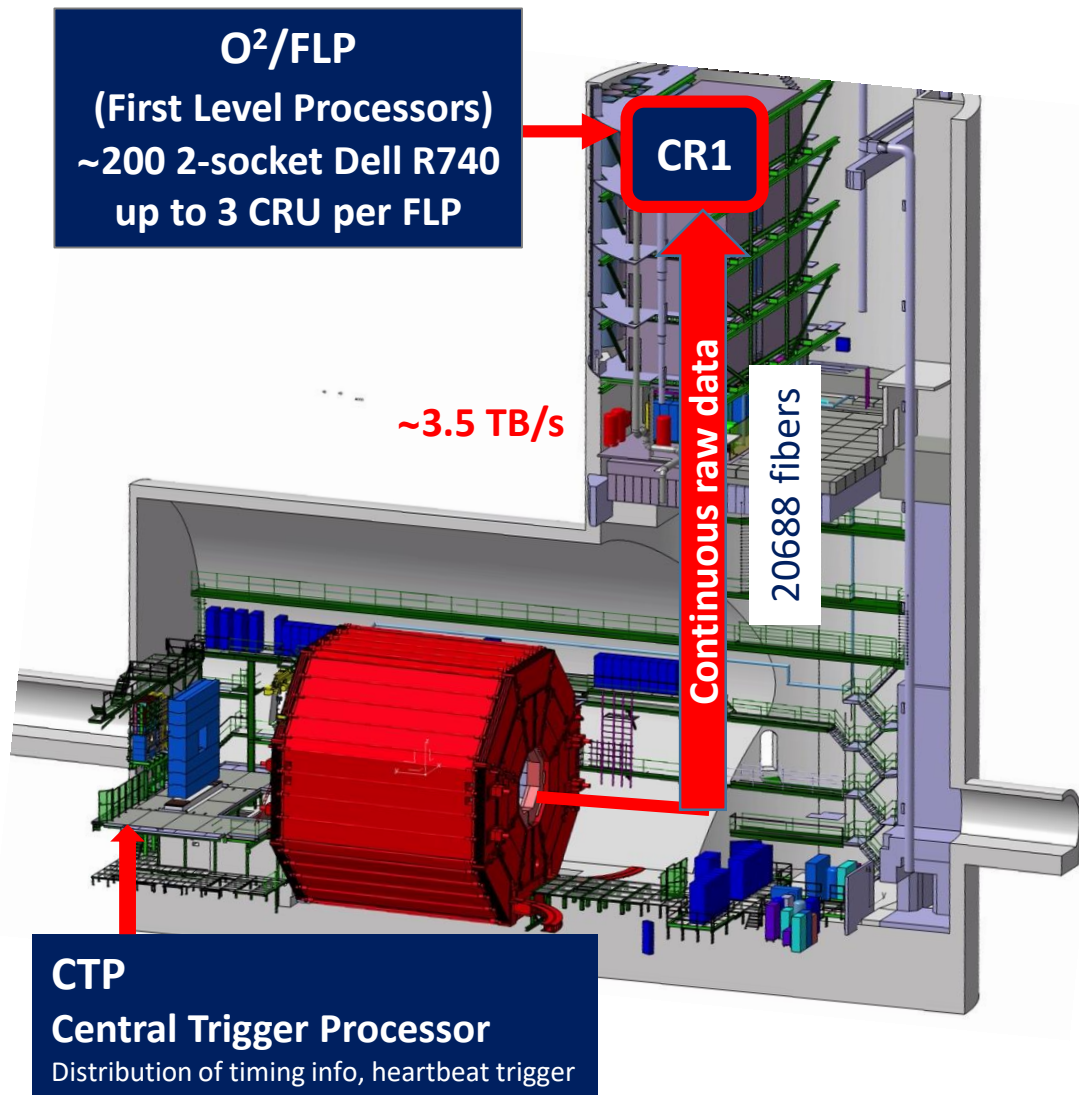


ALICE raw data flow in Run3

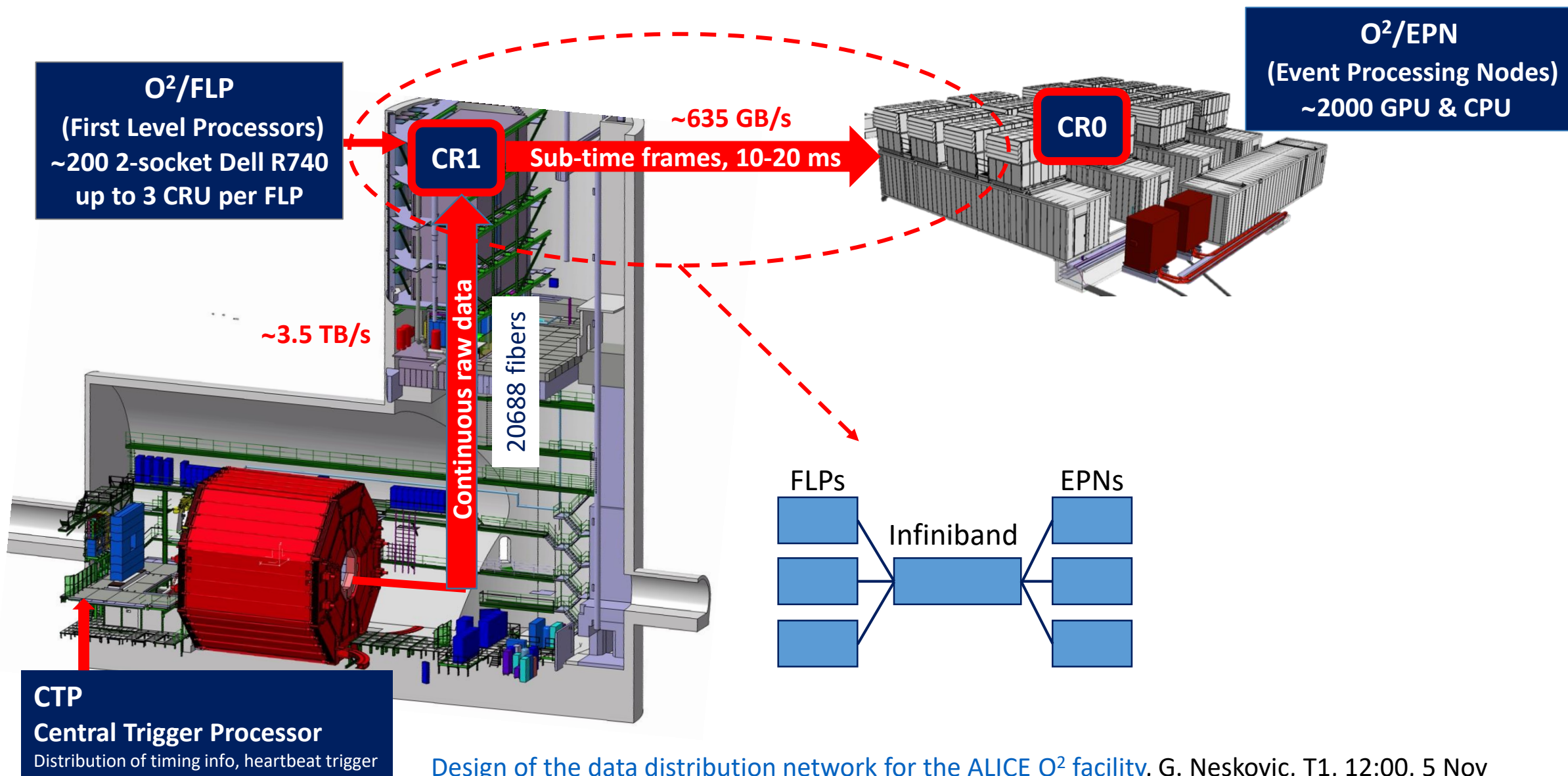


ALICE

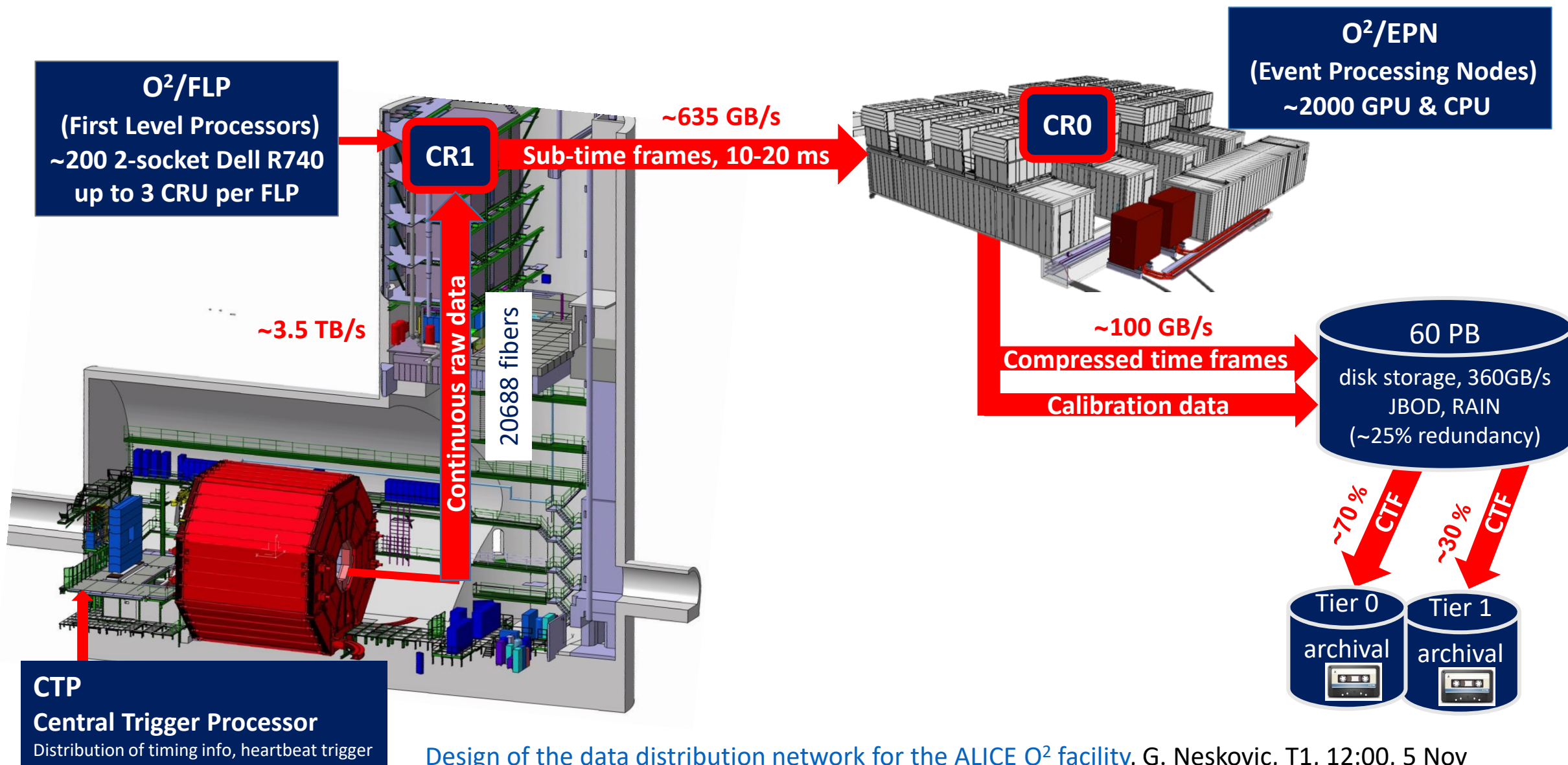
Assessment of the ALICE O² readout servers, F. Costa, T1, 11:00, 4 Nov



ALICE raw data flow in Run3



ALICE raw data flow in Run3





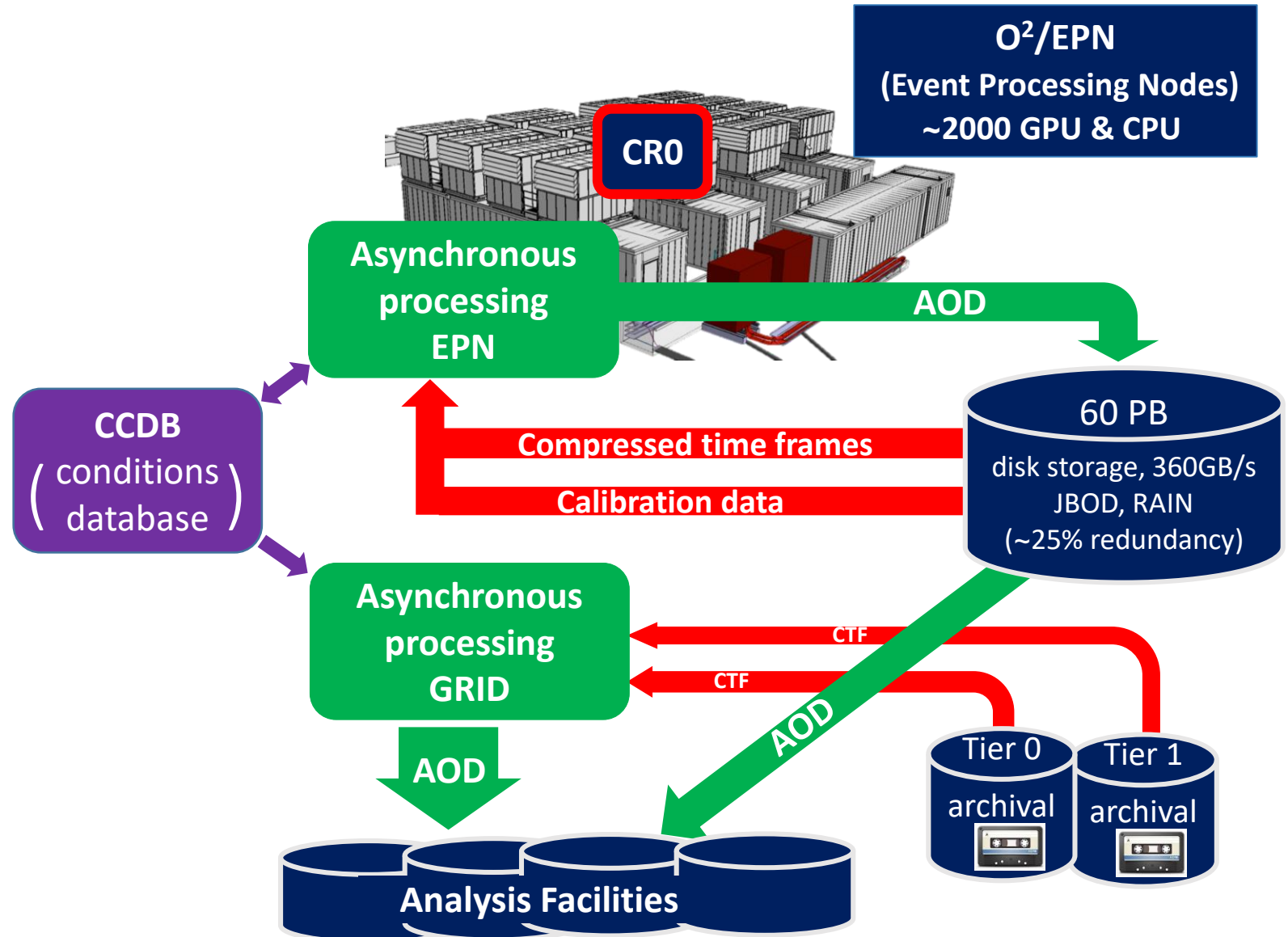
ALICE reconstructed data flow in Run3



Asynchronous phase:

EPN will be used to perform calibrations from data accumulated at the synchronous stage and populate the CCDB (Condition and Calibration Data Base)

Together with GRID will participate in the final reconstruction and distribute Analysis Object Data (AOD) over Analysis Facilities





Heart Beat and Time-Frame



Heart Beat (HB)

issued in continuous & triggered modes to all detectors

Physics trigger

can be sent to upgraded detectors
will be sent to non-upgraded detectors

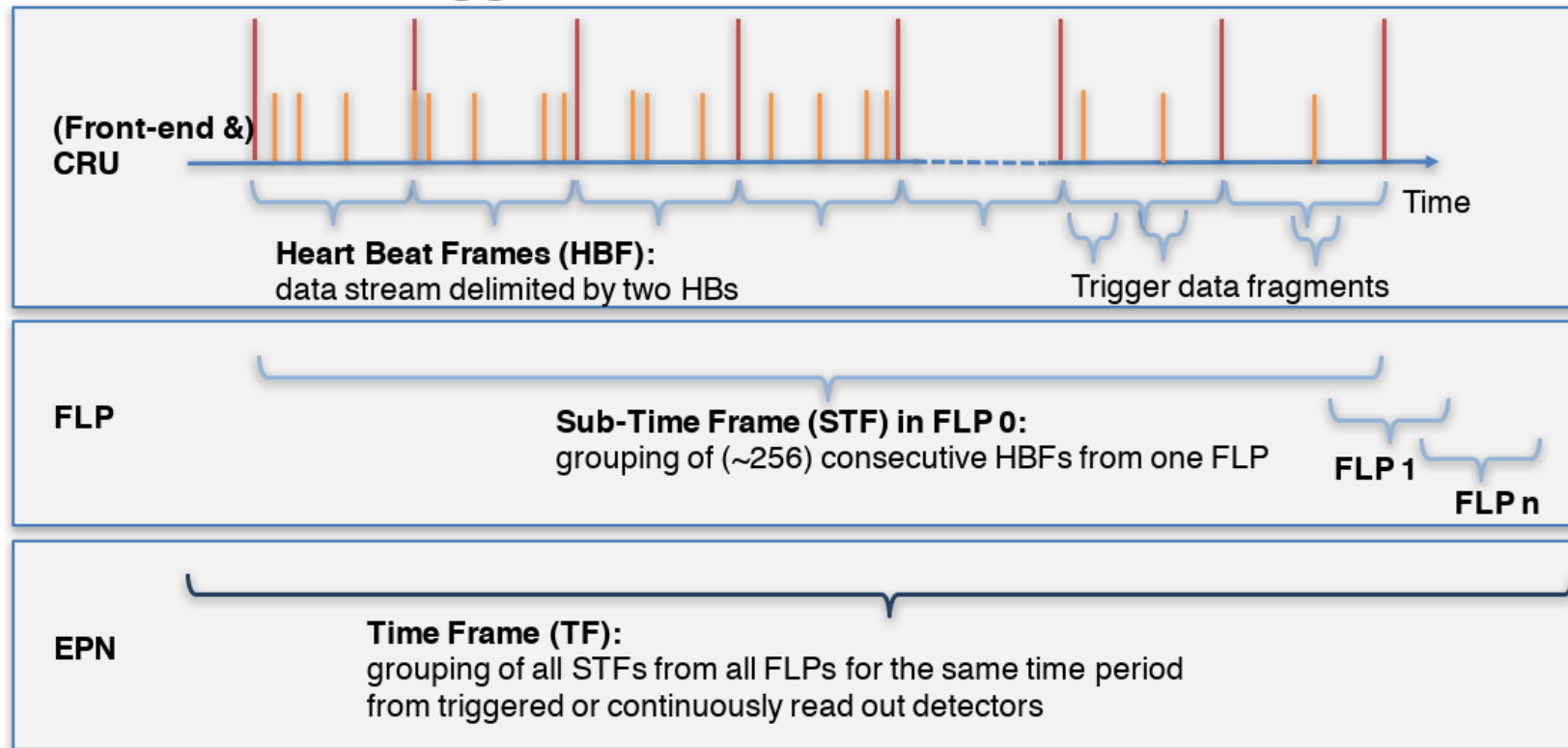
HBF and TF rates programmable

Typical values:

- HB: 1 per orbit, $89.4 \mu\text{s}$: $\sim 10 \text{ kHz}$
- TF: 1 every $\sim 20 \text{ ms}$: $\sim 50 \text{ Hz}$
- $\rightarrow 1 \text{ TF} = \sim 256 \text{ HBF}$

- HB allows synchronization and TF sampling from detectors with continuous and triggered readouts
- Synchronized with LHC clock

Continuous & Triggered read-out



- HB Frame is the smallest chunk of data which is inspected by CTP and can be dropped if the quality is bad.
- FLP sees all data for part of the detector (except small detectors)
- Set of same HBFs (STF) sent to the next EPN for aggregation to TF
- Single EPN sees non-consecutive TFs
 \Rightarrow collisions happening in last HBF may have their TPC clusters (drift \sim HBF) in next TF on other EPN
 \Rightarrow $< 0.5\%$ data loss

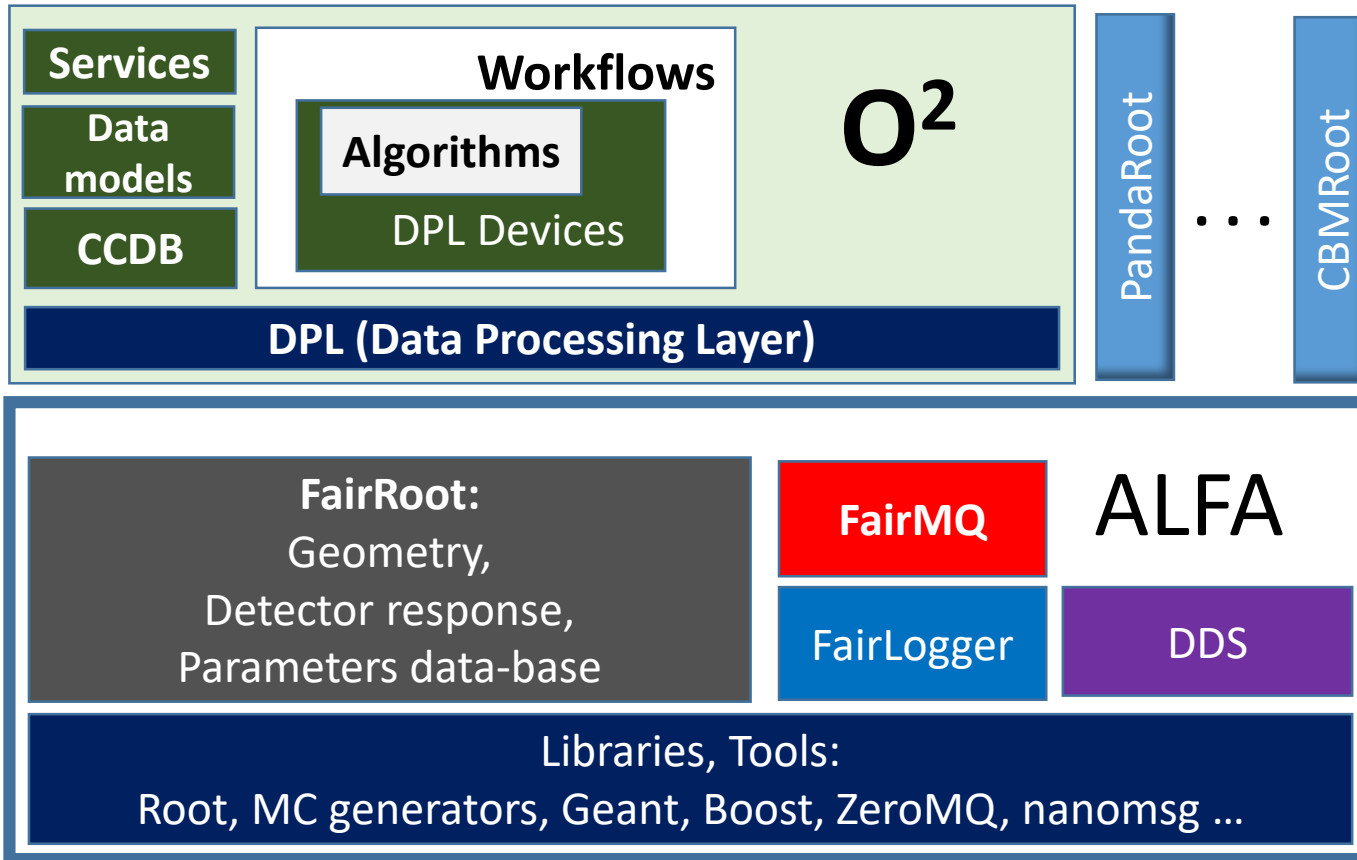


O² software framework



Based on ALFA platform: common project of ALICE and FAIR, derived from FairRoot

See [ALFA: A framework for building distributed applications](#), M. Al-Turany, T5, 11:30, 4 Nov.



- Key feature: **message-queues based parallel processing** done by separate Devices (processes)
- FairMQ supports multiple transport engines (ZMQ, nanomsg) over different protocols (Ethernet, Infiniband, shared memory access)
- Technicalities of messages exchange are hidden in the Data Processing Layer (DPL)
- DPL allows to wrap algorithms to Devices with particular Input and Output message type specifications.
- Workflows are built for group of Devices by automatic matching of their Inputs and Outputs

[Data Analysis using ALICE Run3 Framework](#), G.Eulisse, T6, 11:45, 5 Nov.

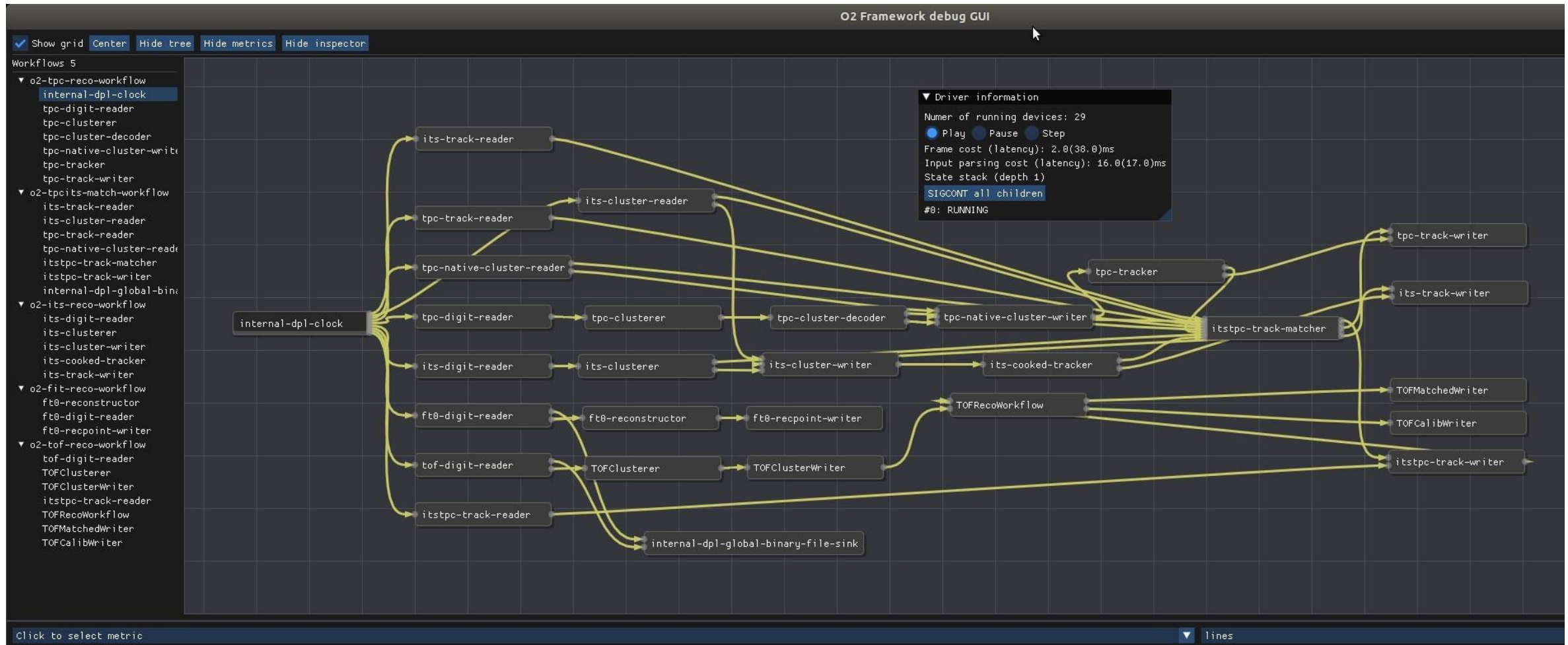
[Running synchronous detector reconstruction in ALICE using declarative workflows](#), M. Richter, TX, 16:30, 5 Nov.



Example of DPL workflow



ITS,TPC,TOF clustering → ITS, TPC tracking → ITS-TPC matching → TOF matching



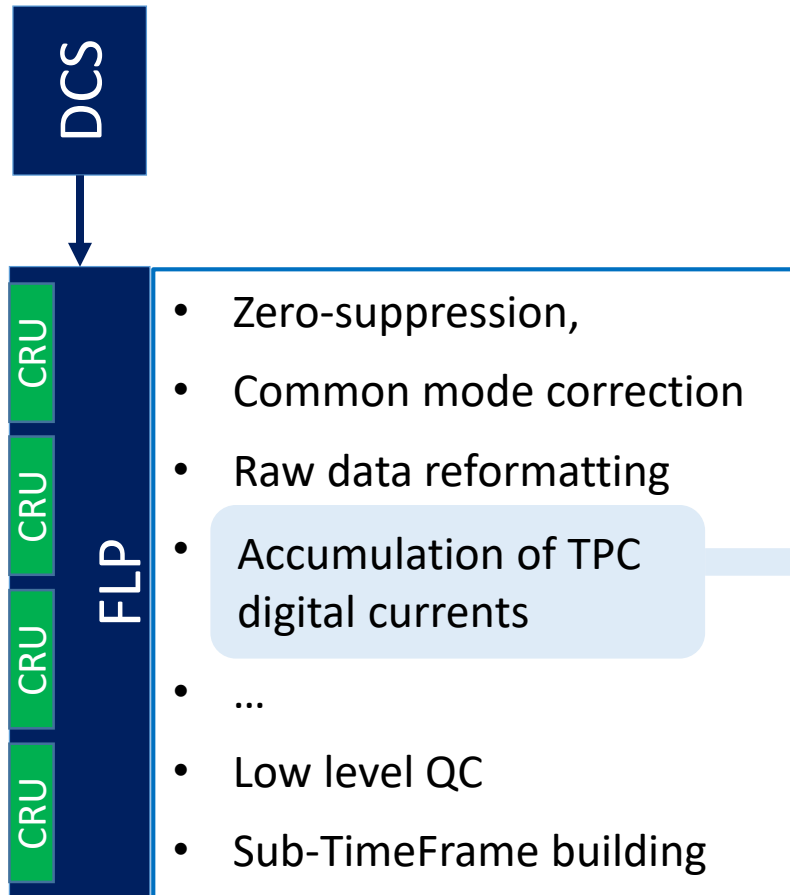
Workflow topology is built at run time by DPL



Outline of the TF synchronous stage processing



Aims: Raw data reduction/compression to Compressed Time Frame (CTF), calibration data accumulation, QC



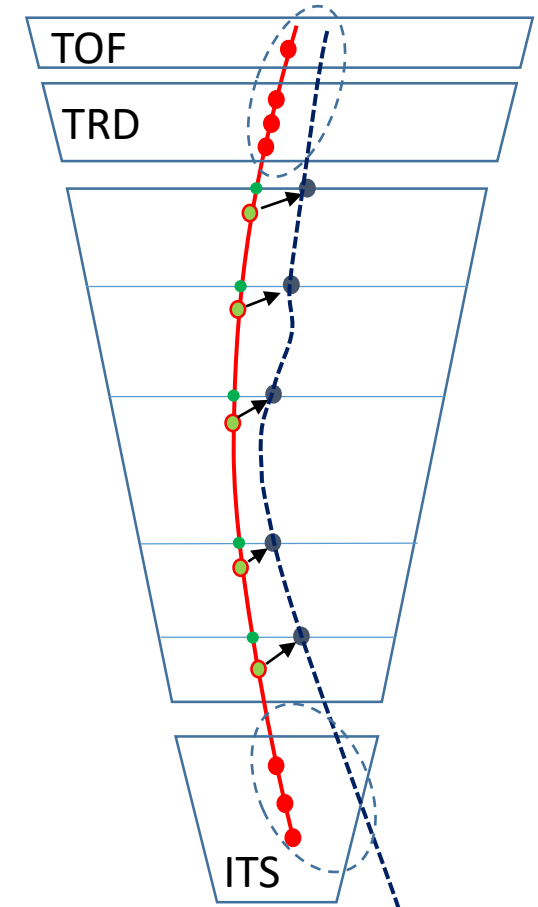
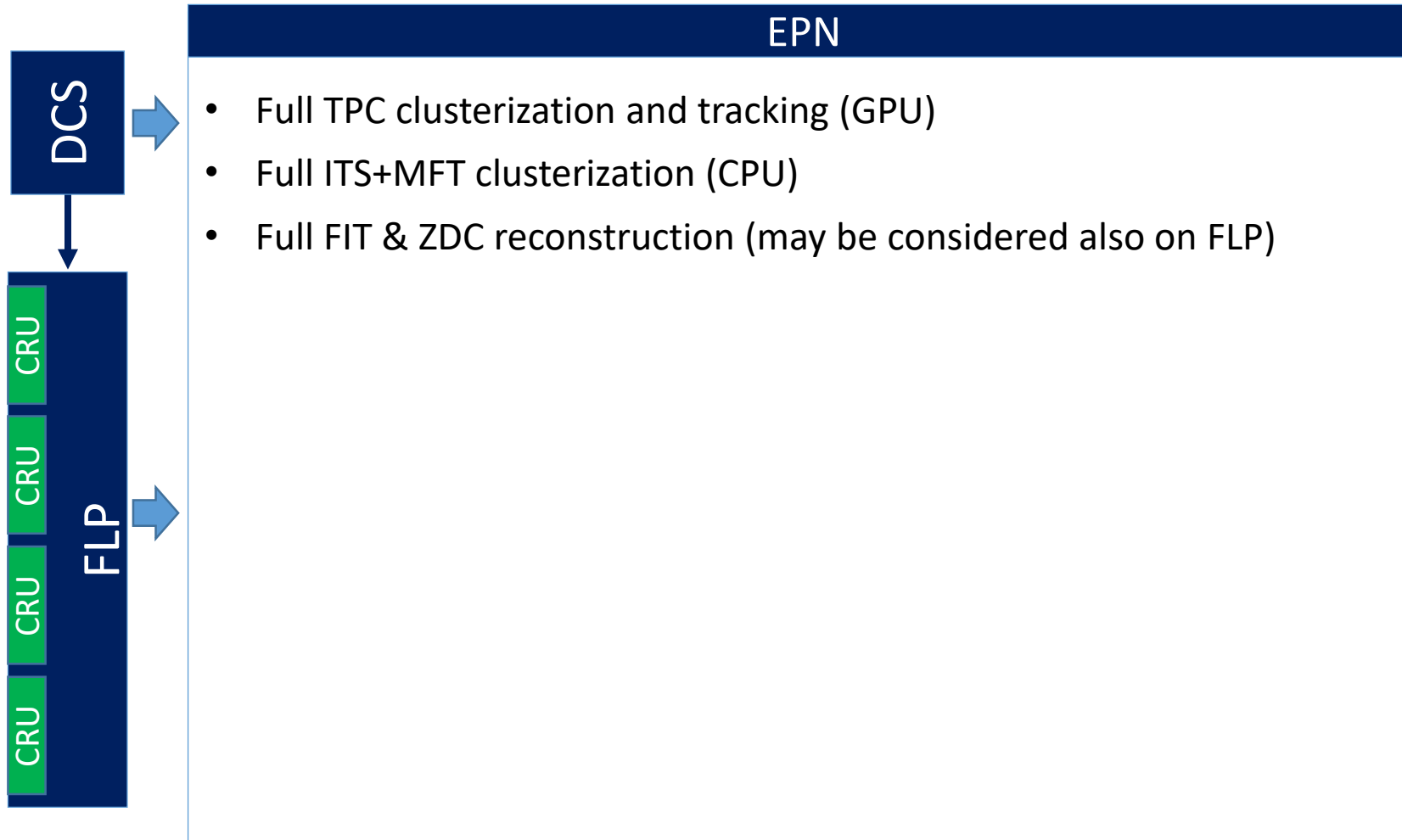
Charge integrated by every TPC pad (or their group) in ~ 1 ms intervals.
Used in TPC calibration to map fluctuations of
Collected charge @ RO planes \rightarrow Space charge \rightarrow Drift line distortions
Due to the slow ion drift, the history for ~ 160 ms (~ 8000 Pb-Pb collisions)
is needed to calibrate ~ 5 ms interval \Rightarrow the only cross-TF data stream



Outline of the TF synchronous stage processing



Aims: Raw data reduction/compression to Compressed Time Frame (CTF), calibration data accumulation, QC

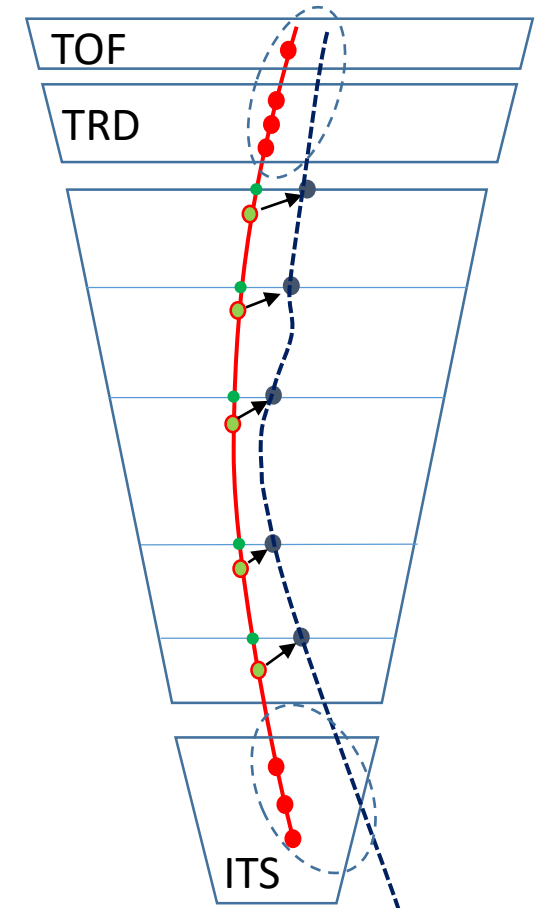
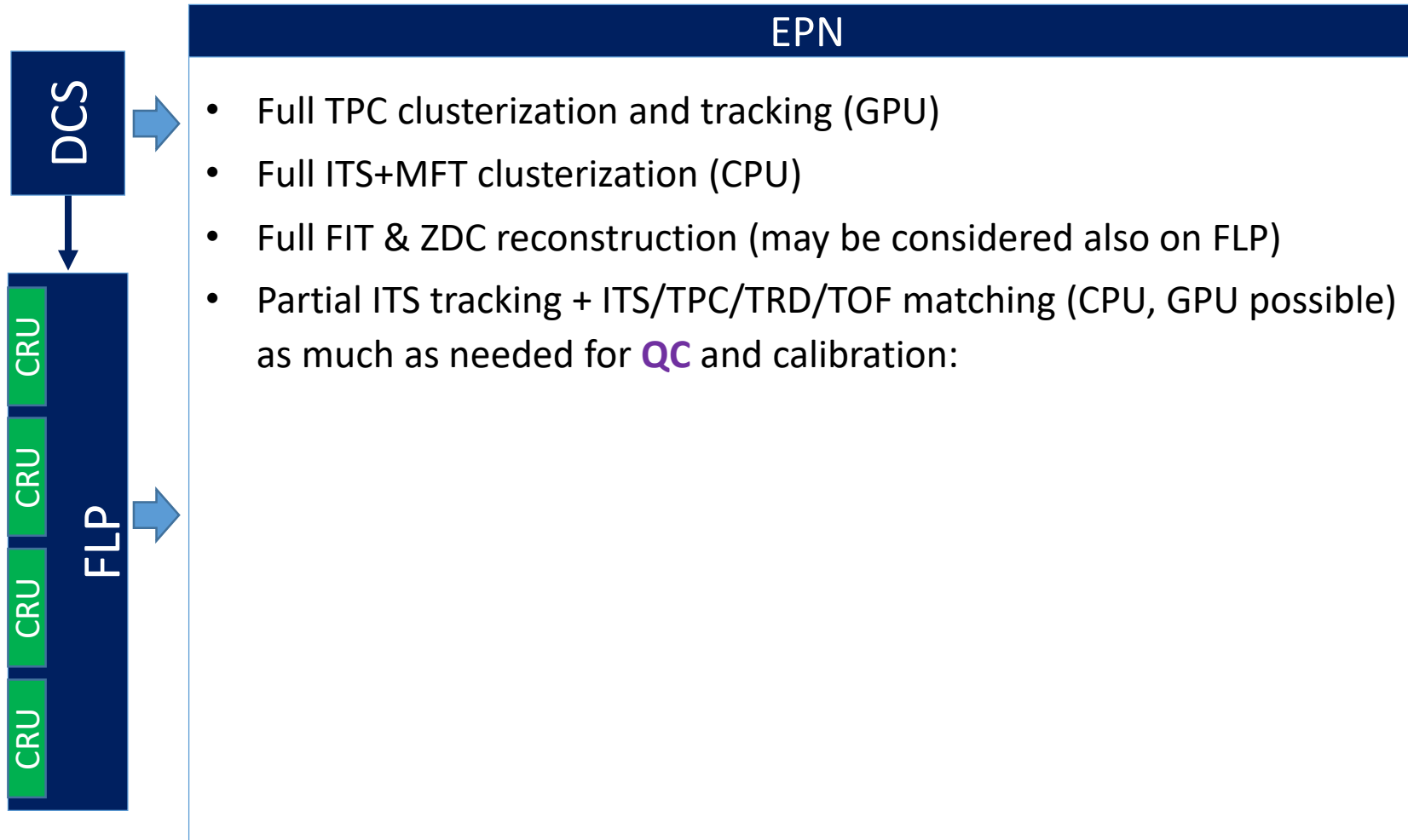




Outline of the TF synchronous stage processing



Aims: Raw data reduction/compression to Compressed Time Frame (CTF), calibration data accumulation, QC

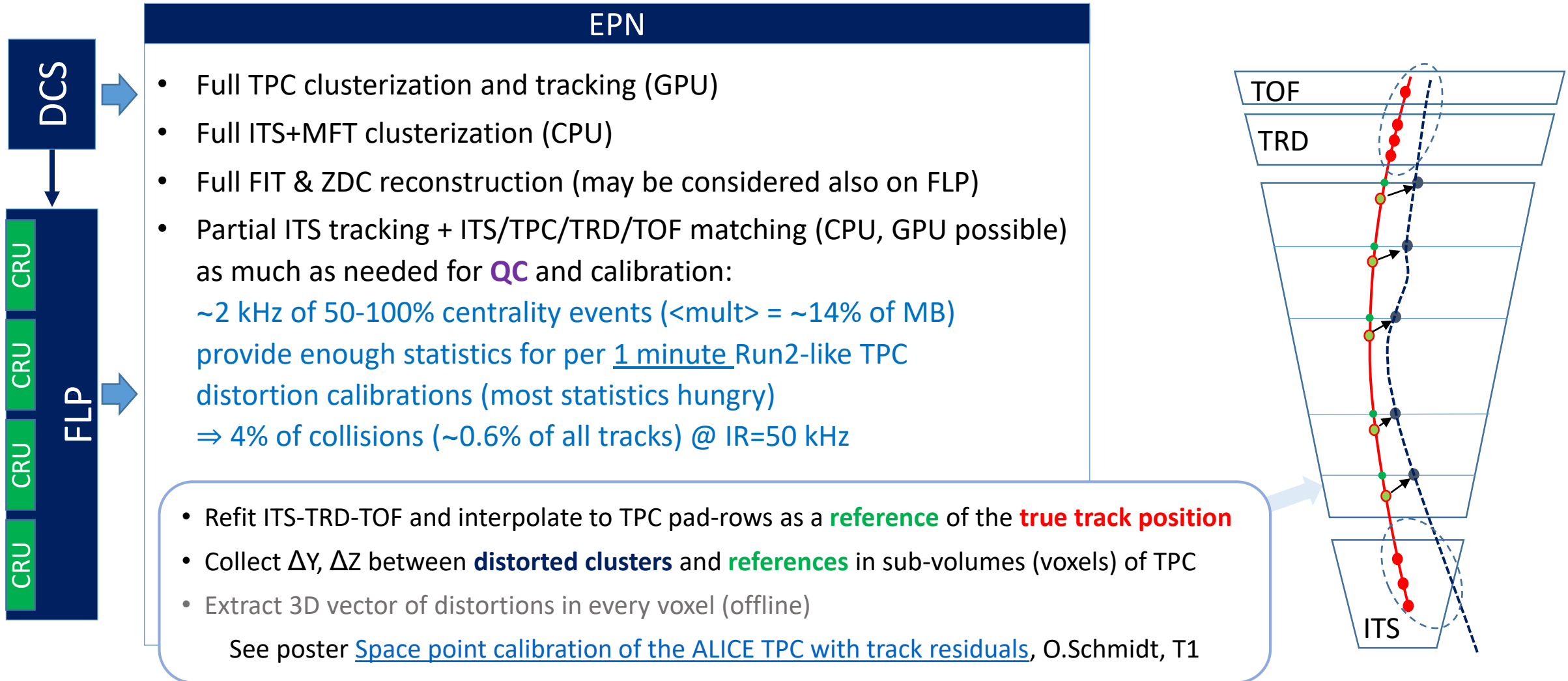




Outline of the TF synchronous stage processing



Aims: Raw data reduction/compression to Compressed Time Frame (CTF), calibration data accumulation, QC



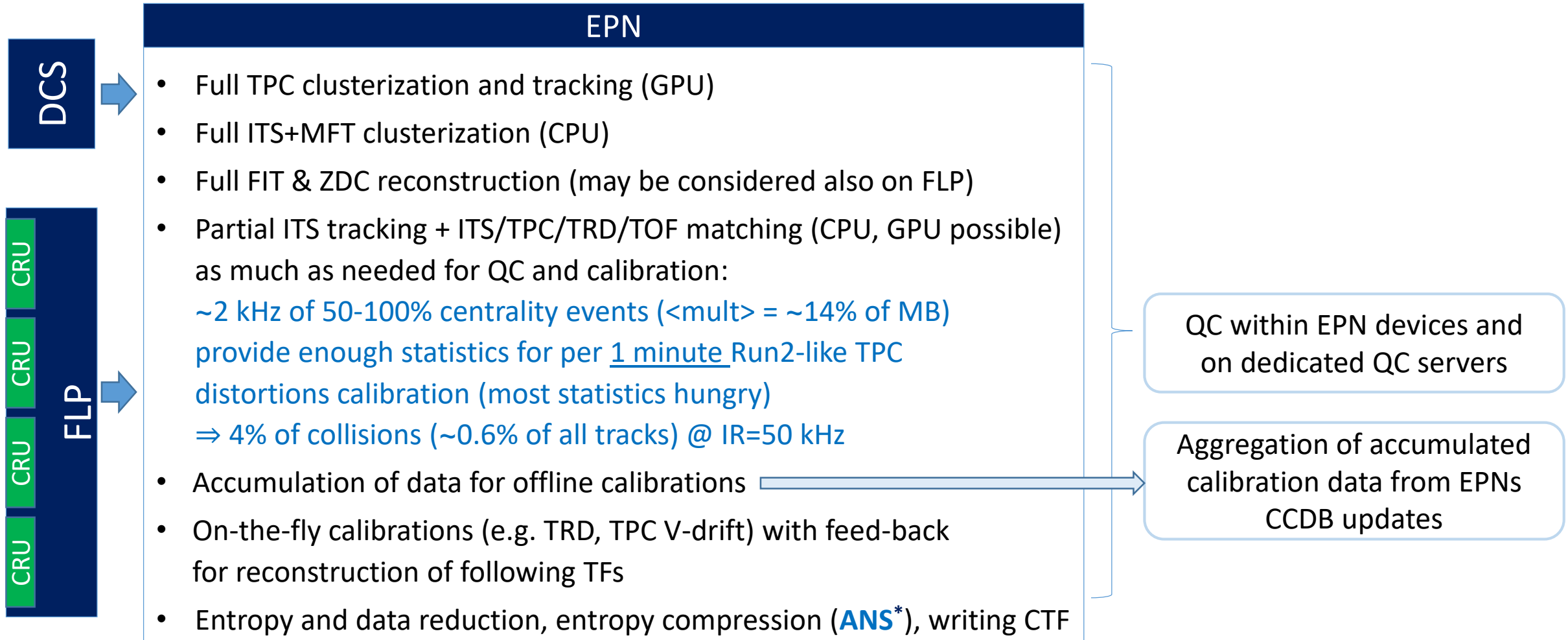
For the **QC** see [The ALICE data quality control system](#), P. Konopka, T1, 15:15



Outline of the TF synchronous stage processing



Aims: Raw data reduction/compression to Compressed Time Frame (CTF), calibration data accumulation, QC



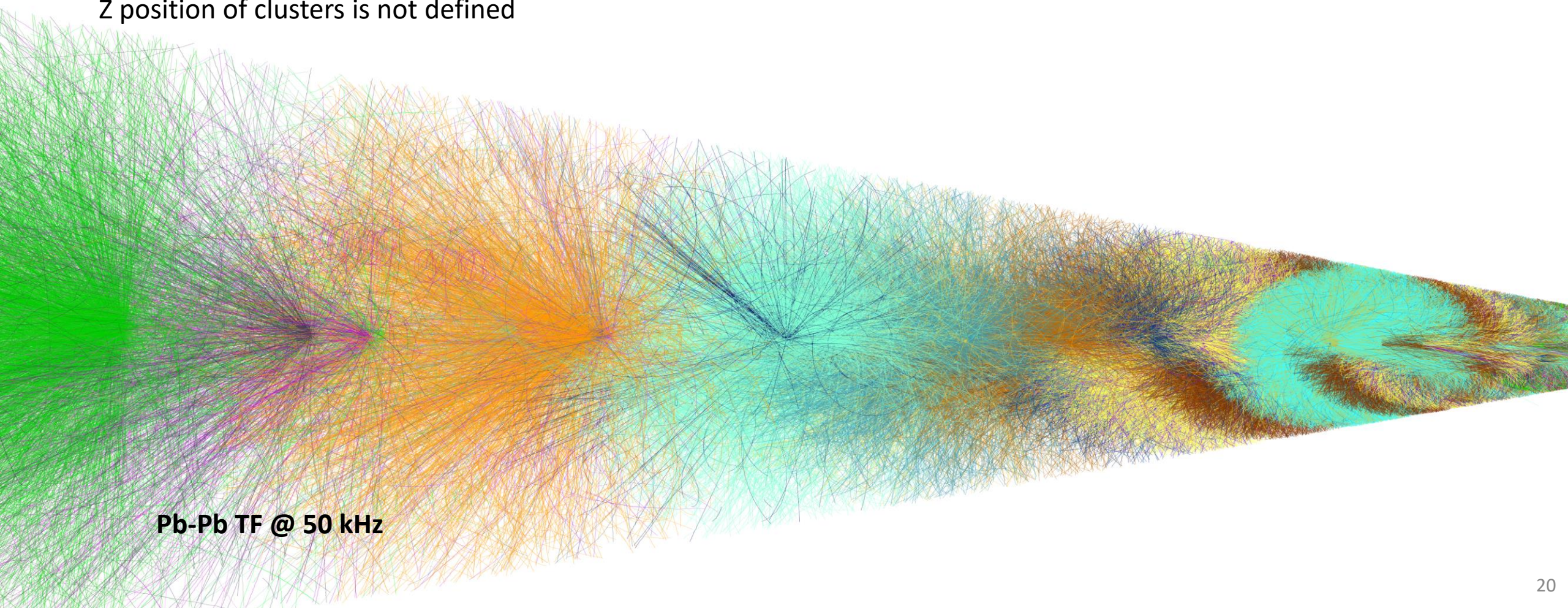


TPC tracking



ALICE

- Principal difference wrt Run1,2:
 - Whole TF (~1000 Pb-Pb) reconstructed in one shot
 - In absence of triggers (reference for drift time estimate)
Z position of clusters is not defined



Pb-Pb TF @ 50 kHz



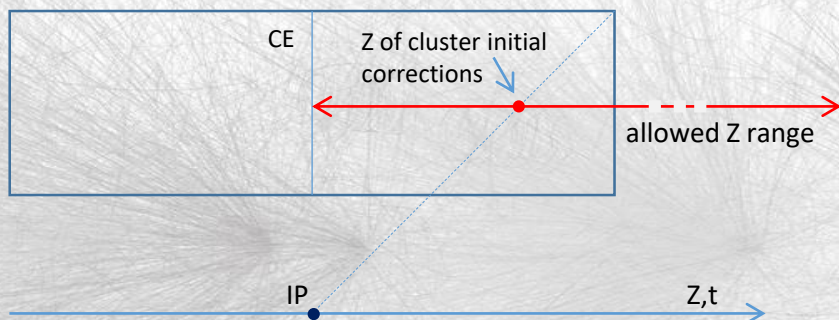
TPC tracking



- Principal difference wrt Run1,2:
 - Whole TF (~1000 Pb-Pb) reconstructed in one shot
 - In absence of triggers (reference for drift time estimate)
Z position of clusters is not defined

→ Cluster correction depends on their Z-position

perform initial calibration assuming cluster belongs to $|\eta| = 0.45$ track from the IP



After matching TPC track to ITS (fixes TPC clusters Z)
recalibrate cluster positions before the refits.

In the synchronous stage will be running on the GPU, for details see:

[GPU-based reconstruction and data compression at ALICE during LHC Run3](#), D.Rohr, TX, today, 14:15

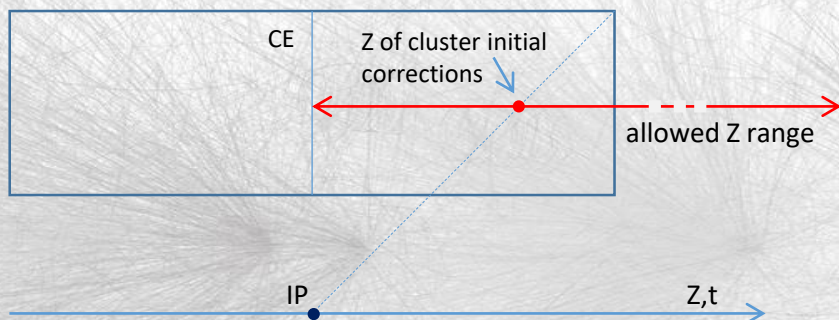


TPC tracking



- Principal difference wrt Run1,2:
 - Whole TF (~1000 Pb-Pb) reconstructed in one shot
 - In absence of triggers (reference for drift time estimate) Z position of clusters is not defined

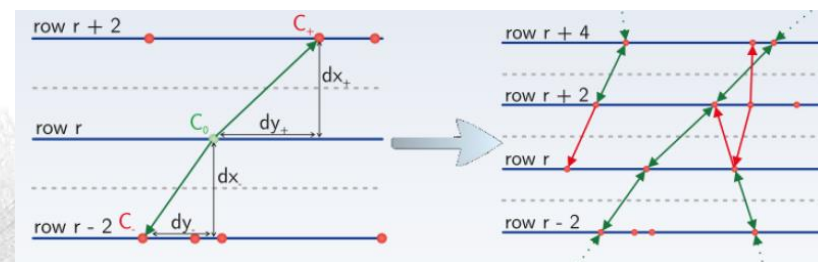
→ Cluster correction depends on their Z-position
perform initial calibration assuming cluster belongs to $|\eta| = 0.45$ track from the IP



After matching TPC track to ITS (fixes TPC clusters Z)
recalibrate cluster positions before the refits.

Improved version of Run1,2 HLT tracking:

- Per sector track finding:
 - Track seeding using Cellular Automaton:
 - Finding straight-line hit triplets on adjacent rows
 - Concatenating compatible triplets to seeds, fitting



○ Track following (Kalman filter)

- Merging between sectors and over central electrode
- Inside/outside refit with cluster reattachment
- Looping track legs merging (found from $p_T > \sim 15$ MeV/c)

In the synchronous stage will be running on the GPU, for details see:

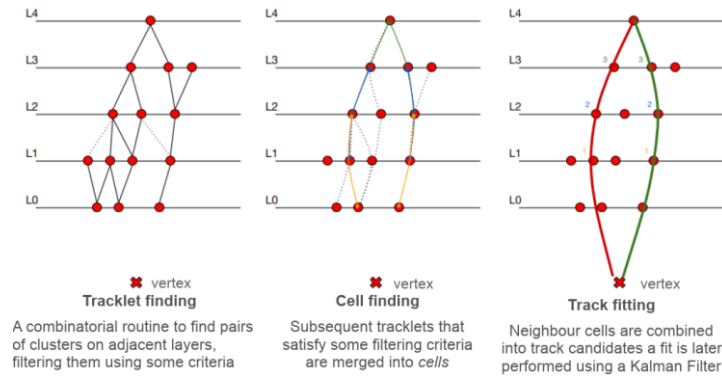
[GPU-based reconstruction and data compression at ALICE during LHC Run3](#), D.Rohr, TX, today, 14:15



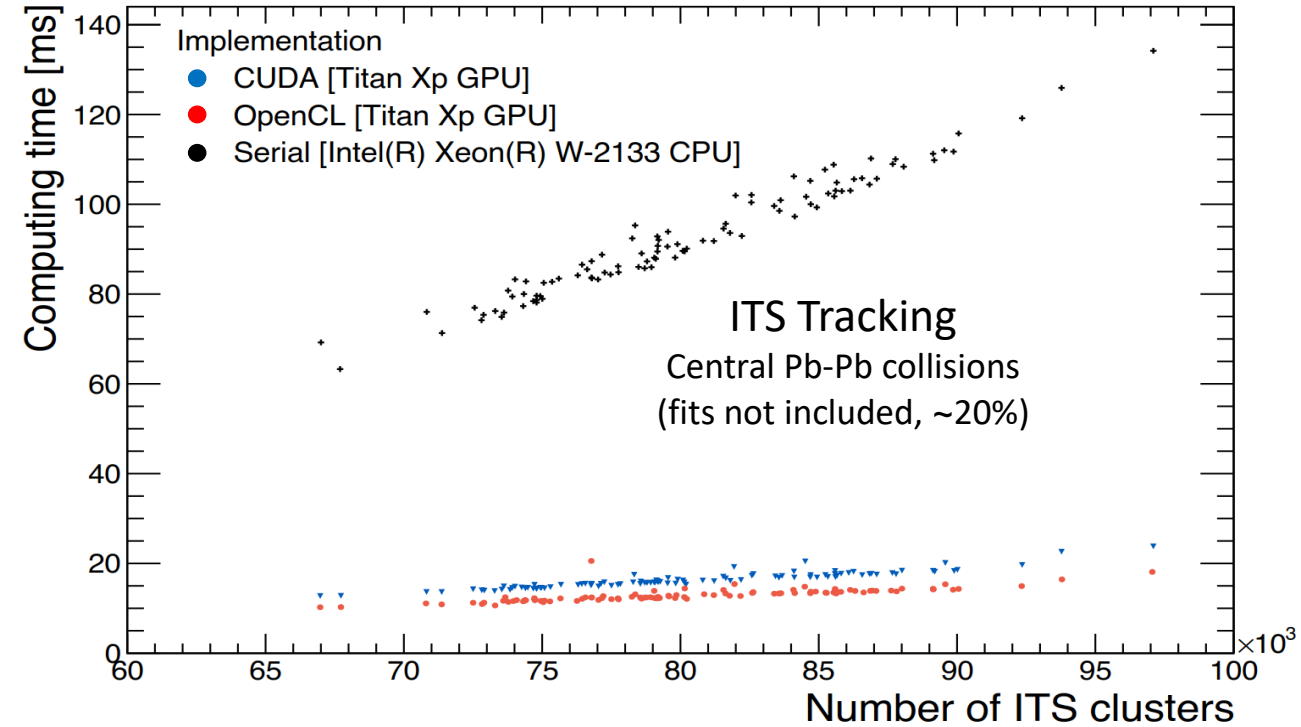
ITS tracking, matching to TRD, TOF



- Fast primary vertex finding with tracklets from 3 inner layers
- Standalone track-finding in ITS using Cellular Automaton



- Multiple passes for prompt (constrained to vertex) long, incomplete and off-vertex tracks (in asynchronous phase)
- Runs both on CPU and GPU

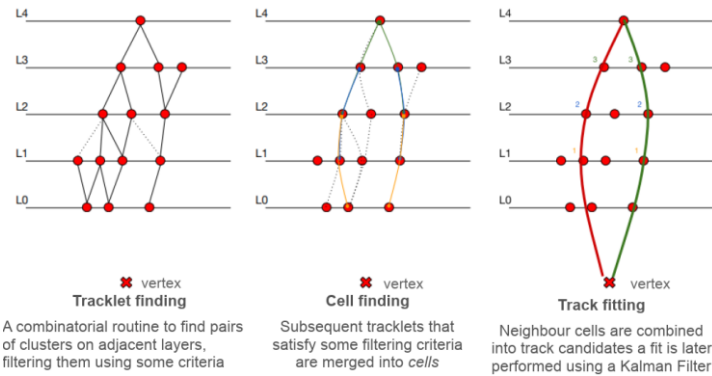




ITS tracking, matching to TRD, TOF

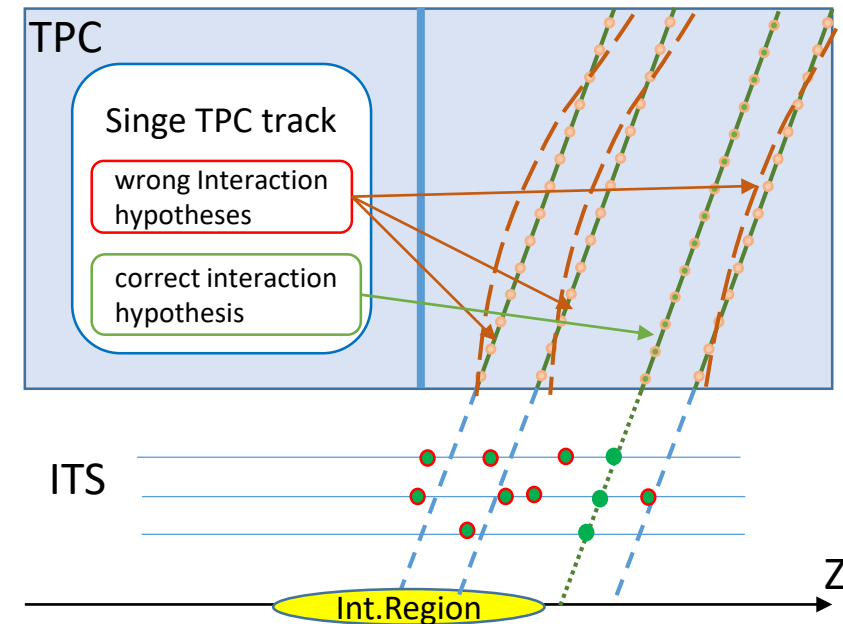


- Fast primary vertex finding with tracklets from 3 inner layers
- Standalone track-finding in ITS using Cellular Automaton



- Multiple passes for prompt (constrained to vertex) long, incomplete and off-vertex tracks (in asynchronous phase)
- Runs both on CPU and GPU

- ITS-TPC track-to-track matching
- ITS in Run3 has long integration time ($\sim 2 - 20 \mu\text{s}$)
 → may see multiple collisions in single frame
 Matching resolves interaction time to $\sim 100 \text{ ns}$
- ITS-TPC refitted outwards and matched to TRD and TOF
- Afterburner (in asynchronous phase only): matching remaining ITS clusters of TPC tracks with Z constrained by different interaction times from FIT





TPC data reduction

- Target in ideal case: reject ~50% of clusters
- Two alternative scenarios:

A: Keep all clusters except those from identified

- 1) (looping) tracks $p_T < 50$ MeV/c (not needed for physics)
- 2) extra legs of loopers $50 < p_T < 200$ MeV/c
- 3) segments of tracks with high inclination to pad-rows ($\varphi > 70^\circ$)

currently only ~13% rejection rate achieved

in principle, ~39% rejection achievable if

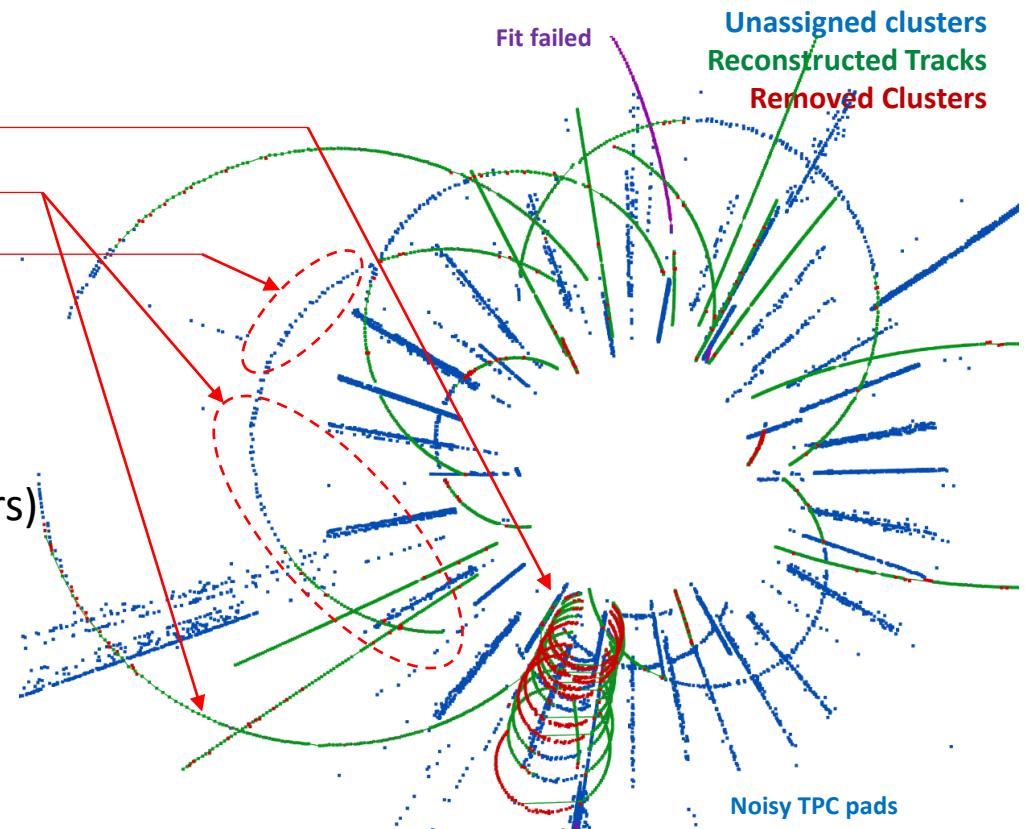
- merging of looper legs improved
- ➔ • looper tagging can be extended to $p_T < 10$ MeV (~15% of clusters)
(track radii < 6cm, Hough transform is tested)

B: Keep only clusters attached or in the vicinity of tracks

interesting for physics ($p_T > 50$ MeV/c, principal leg for loopers)

currently ~37% rejection achieved

in principle, ~52% rejection achievable in case of ideal loopers legs merging

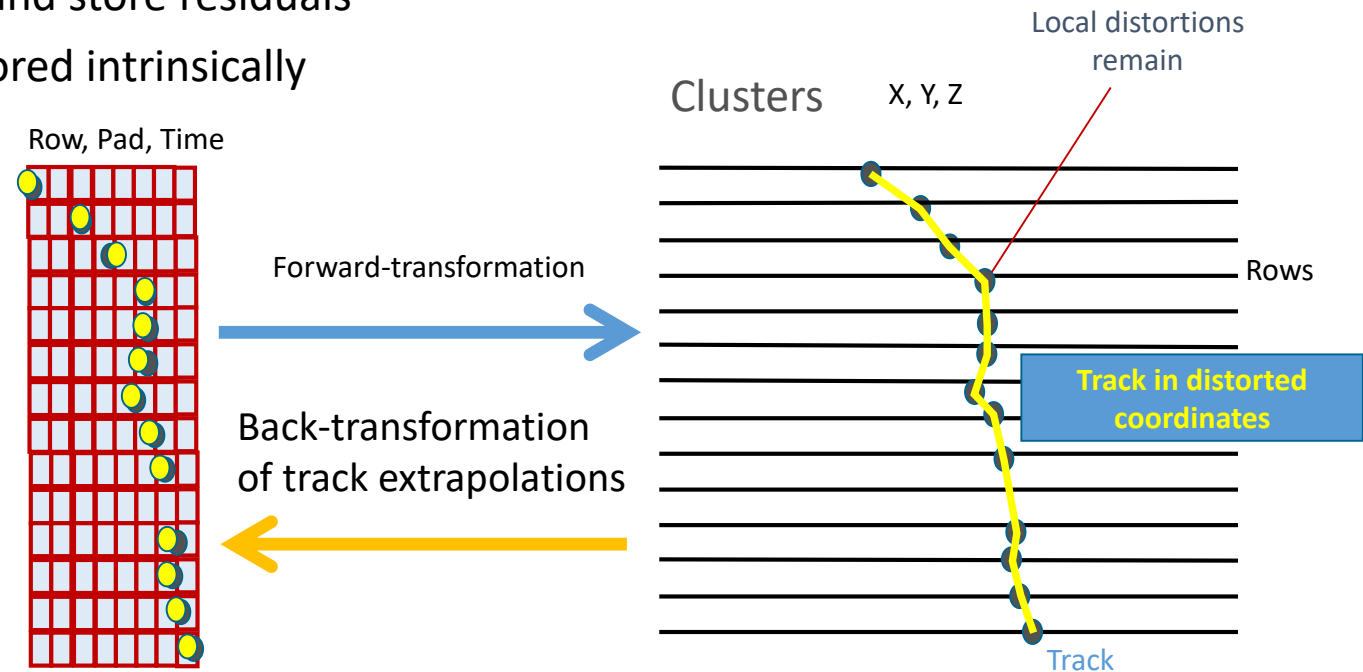




TPC data compression

We aim for efficient storing of raw clusters (pad, row, time) data in view of later reconstruction with best possible calibration.

- 👍 ~ 50% of clusters can be attached to tracks \Rightarrow minimize entropy by exploring correlation between these clusters:
 - Find tracks with clusters transformed and corrected with best available calibration
 - Transform again raw clusters of the track using fast linear transformation and refit tracks in distorted coordinates
 - Starting with distorted track parameters perform Kalman update / extrapolation from row to row
 - Transform extrapolations back to raw coordinates and store residuals
- 👍 Additional benefit: cluster to track association is stored intrinsically
 - Data rounded to relevant (logarithmic) precision and encoded with ANS compression





TPC data compression

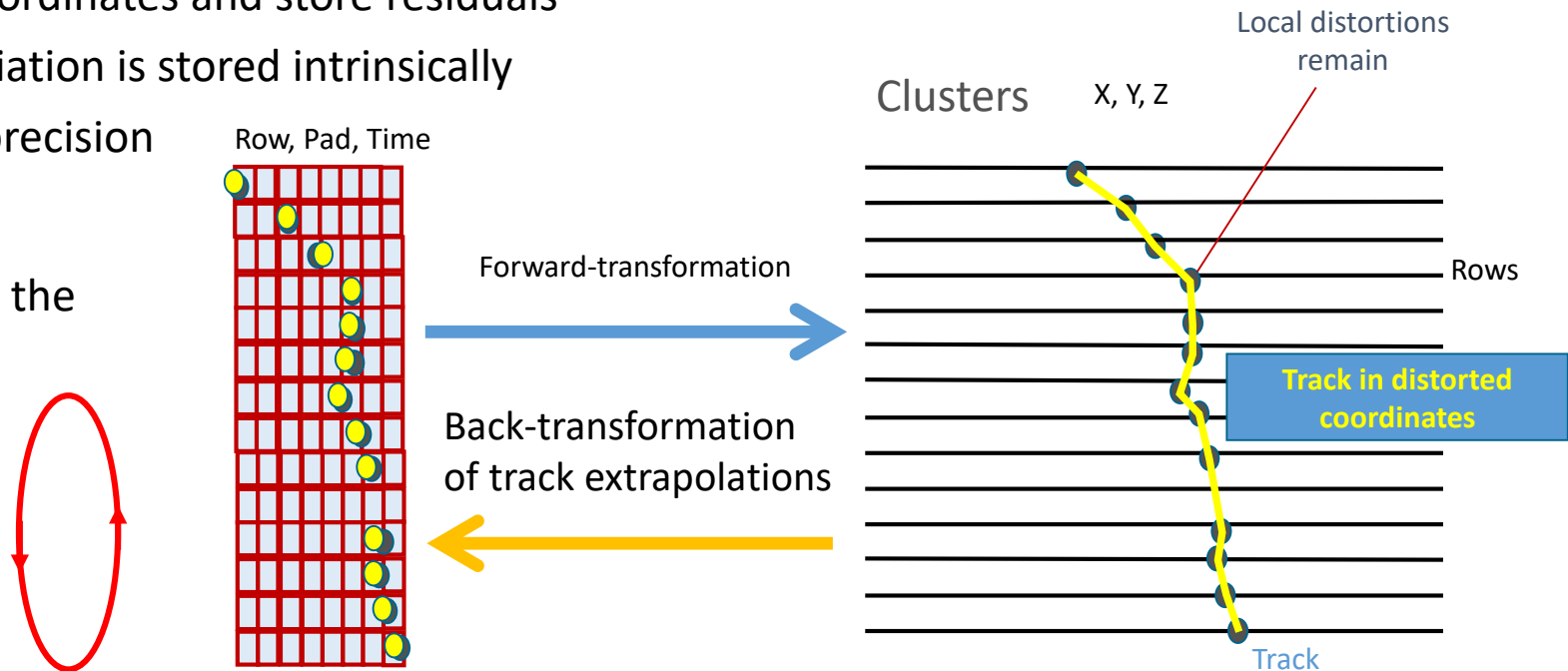


We aim for efficient storing of raw clusters (pad, row, time) data in view of later reconstruction with best possible calibration.

- ~ 50% of clusters can be attached to tracks \Rightarrow minimize entropy by exploring correlation between these clusters:
 - Find tracks with clusters transformed and corrected with best available calibration
 - Transform again raw clusters of the track using fast linear transformation and refit tracks in distorted coordinates
 - Starting with distorted track parameters perform Kalman update / extrapolation from row to row
 - Transform extrapolations back to raw coordinates and store residuals
- Additional benefit: cluster to track association is stored intrinsically
 - Data rounded to relevant (logarithmic) precision and encoded with ANS compression

Decoding: perform inverse procedure using the same linear transformation:

- recover cluster coordinates from stored residuals + current track position
- do Kalman update with cluster
- extrapolate to next layer...





ITS (MFT) clusters compression

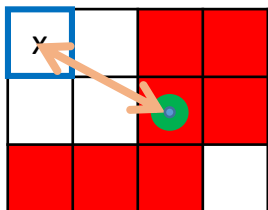


- ITS will ship ~30 GB/s data in continuous readout with ~6 μ s strobe (~10¹⁰ fired pixels/s with noise 10⁻⁷ /channel)
- Data are redundant: ~50% of fired pixels are fired in 2 consecutive strobos (TOT \approx strobe)
 \Rightarrow Mask repeatedly fired pixels and store clusterized data.

- Profit from the highly non-uniform frequency of different cluster patterns (low entropy)

\rightarrow Build table of patterns sorted in frequency, for each pattern pre-calculate its properties:

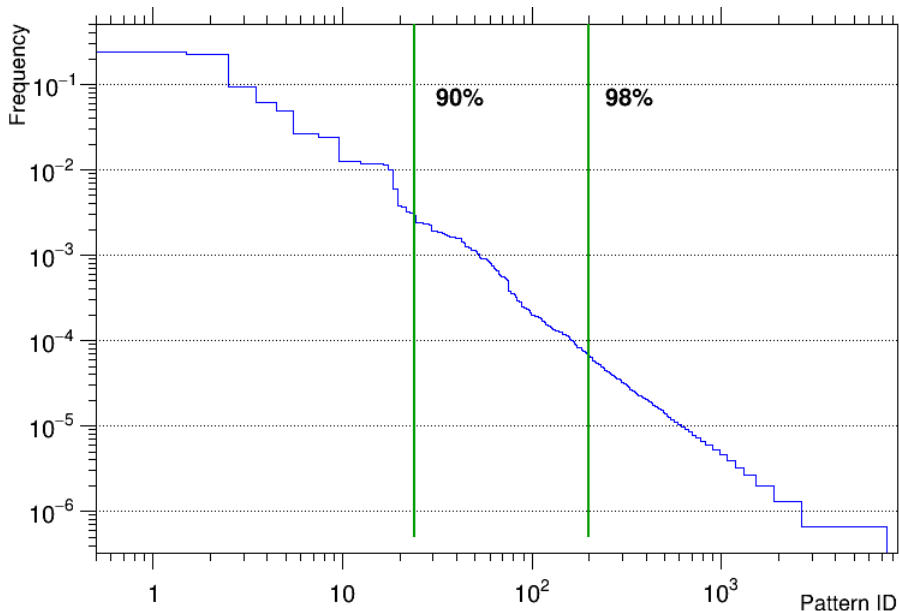
- Offset of **COG** wrt **reference pixel** (e.g. corner of cluster bounding box)



- Mean error between COG and MC hit position
- (in future) error vs impact angles

\rightarrow Encoding: store reference pixel column/row and pattern ID

\rightarrow During decoding directly obtain cluster position w/o actual COG calculation



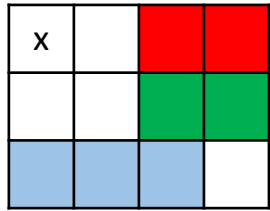
Pattern	Fraction	ID
	24.3%	0
	22.5%	1
	9.2%	2
	6.0%	3
	4.9%	4
	2.6%	5
	2.4%	6
	2.4%	7
	1.2%	8



ITS (MFT) clusters compression



Find cluster, define b-box



Convert to bitmap



calculate its hash



#



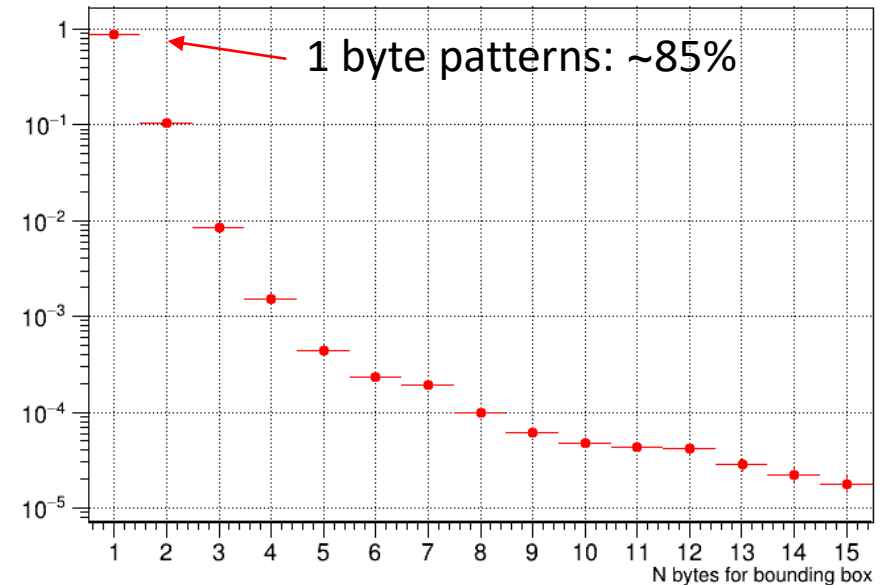
assign ID of this pattern
from pre-calculated list
sorted in frequency



For short patterns ($N_{col} \times N_{row} \leq 8$)
extract ID directly from LUT

Clusters within the chip are sorted in column, chips sorted in ID...

➔ store ΔChipID . Δcolumn , row, patternID

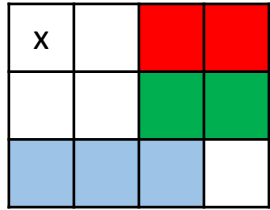




ITS (MFT) clusters compression



Find cluster, define b-box



Convert to bitmap



calculate its hash

#



assign ID of this pattern from pre-calculated list sorted in frequency

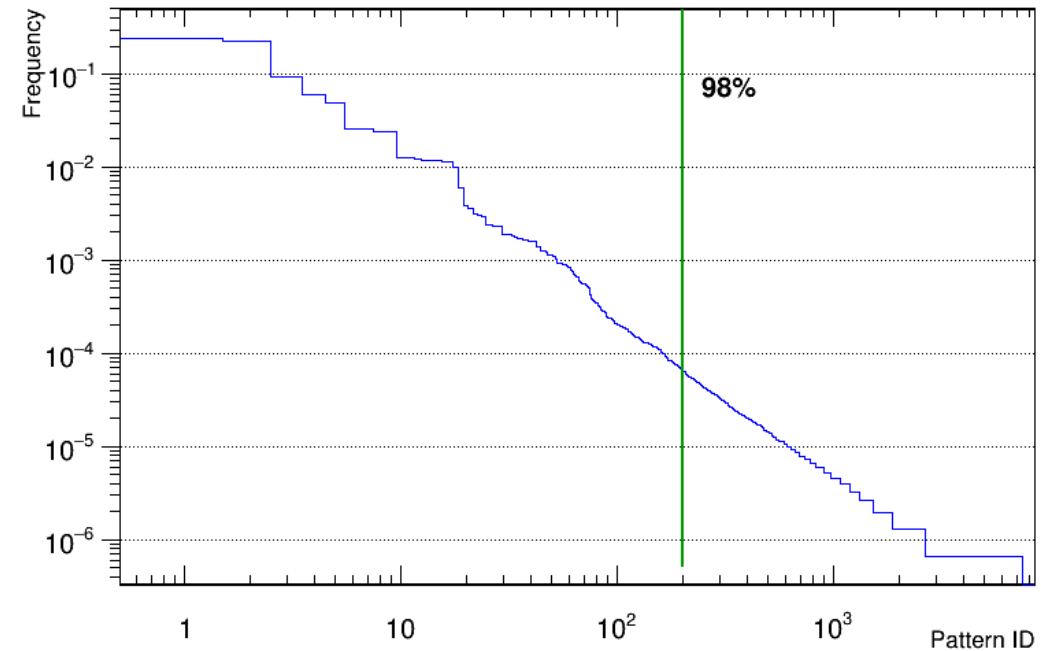


For short patterns ($N_{col} \times N_{row} \leq 8$) extract ID from directly from LUT

Clusters within the chip are sorted in column, chips sorted in ID...

→ store ΔChipID , Δcolumn , row , patternID

- 👍 Most frequent ~98% of patterns (220) have unique IDs ⇒ lossless
- 👎 To keep encoding alphabet short, group rare patterns with similar properties (e.g. b-box dimensions) under the same ID ⇒ lossy (although lost information is not used in standard reconstruction)

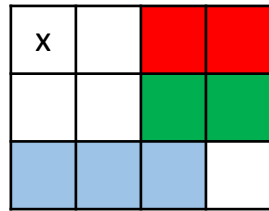




ITS (MFT) clusters compression



Find cluster, define b-box

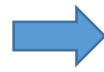


Convert to bitmap



calculate its hash

#



assign ID of this pattern from pre-calculated list sorted in frequency



For short patterns ($N_{col} \times N_{row} \leq 8$) extract ID directly from LUT

Clusters within the chip are sorted in column, chips sorted in ID...

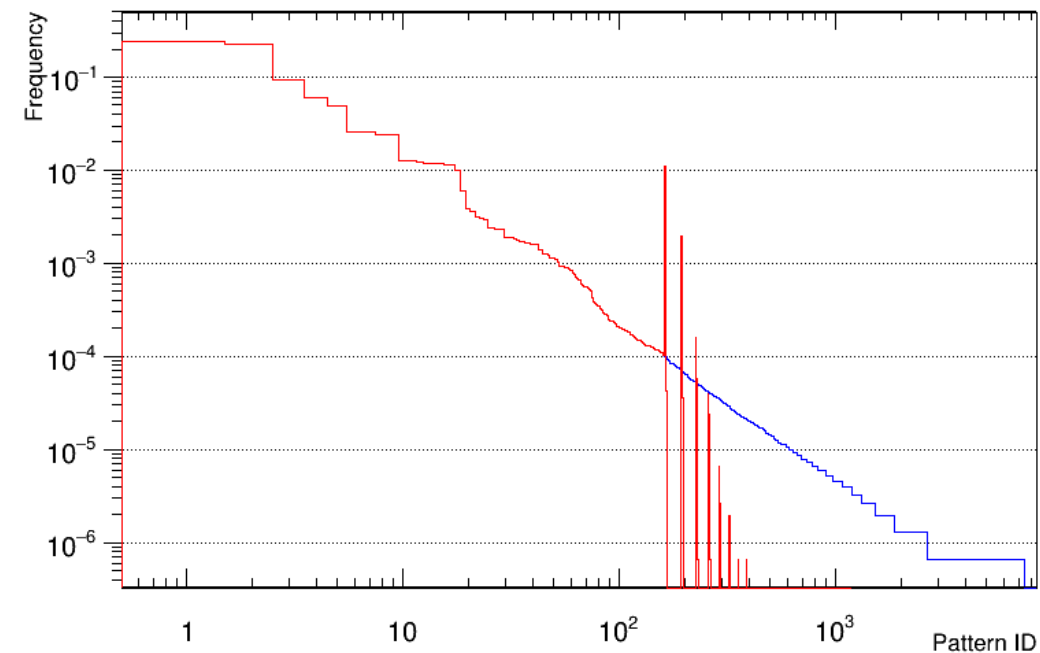
➔ store ΔChipID , Δcolumn , row , patternID

👍 Most frequent ~98% of patterns (220) have unique IDs \Rightarrow lossless

👎 To keep encoding alphabet short, group rare patterns with similar properties (e.g. b-box dimensions) under the same ID \Rightarrow lossy (although lost information is not used in standard reconstruction)

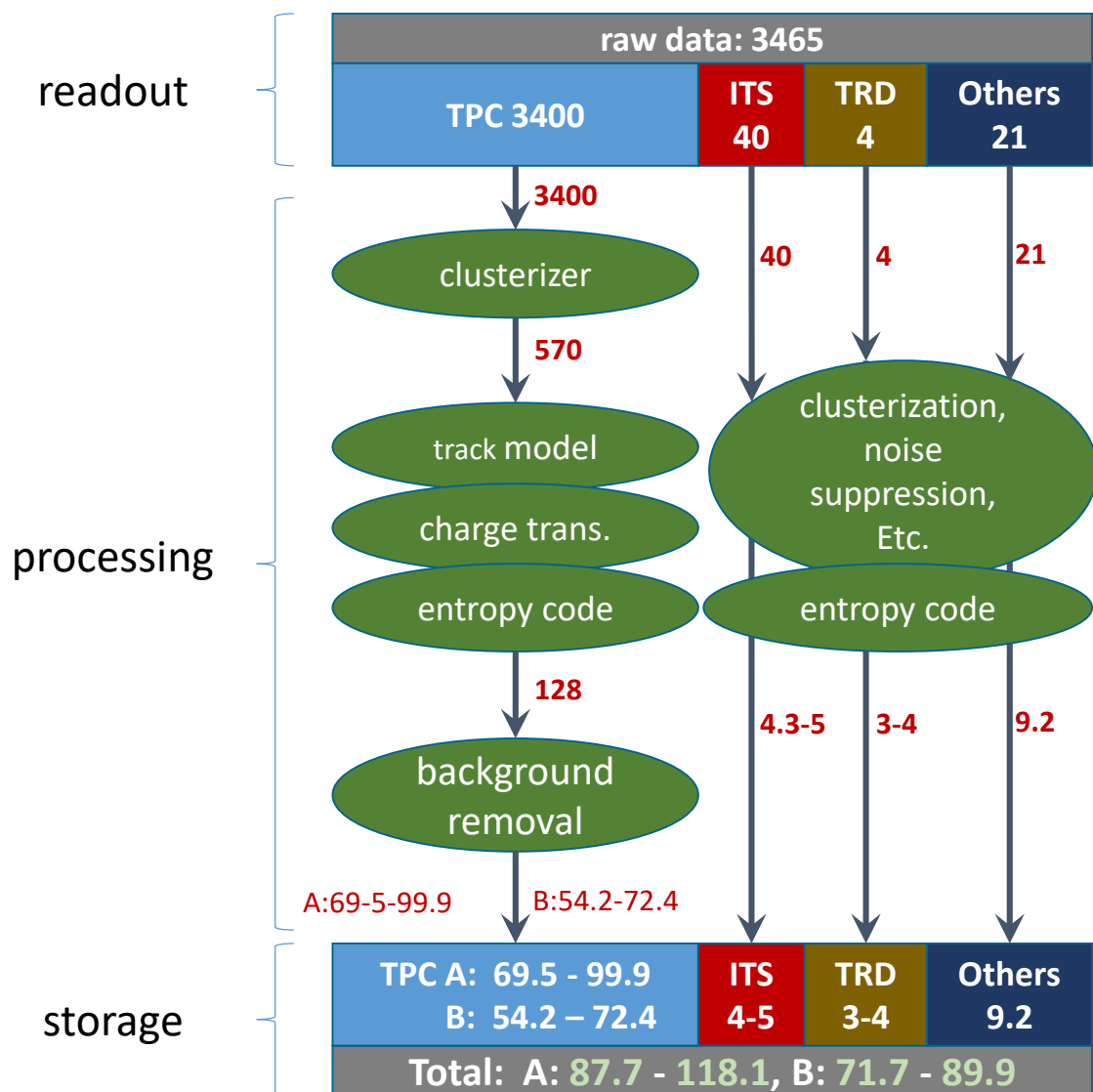
➔ Store separate stream with row/column increments wrt reference pixel for clusters with these rate topologies \Rightarrow ~5% overhead

➔ ~ 22-30 bits/cluster \Rightarrow ~30 GB/s \rightarrow ~4.5 GB/s





Data rates (in GB/s) at online stage



- For remaining detectors data compression consists of
- reformatting raw data (TRD, TOF)
 - Storing signal extracted from ADC samples (EMCal, PHOS: 30 10-bit samples → 2 floats per cell)
 - Storing clusters discarding pads data (MUON, feasibility still being evaluated)

Two alternative scenarios for TPC data lossy compression (both will be implemented)

- A. Reject only clusters or identified noise and background tracks (loopers):
Rejects: 12.5% - 39.1%
- B. Keep only clusters attached or in the proximity of identified signal tracks.
Rejects: 37.3% - 52.5%



ALICE Upgrade presentations at CHEP 2019



November 4:

[ALFA: A framework for building distributed applications](#), M. Al-Turany, T5, 11:30

[Jiskefet, a bookkeeping application for ALICE](#), M.Teitsma, T4, 11:45

[AliECS: a New Experiment Control System for the ALICE Experiment](#), T. Mrnjavac, T1, 14:00

[The ALICE data quality control system](#), P. Konopka, T1, 15:15

November 5:

[Assessment of the ALICE O² readout servers](#), F. Costa, T1, 11:00

[A VecGeom navigator plugin for Geant4](#), S. Wenzel, T2, 11:30

[Design of the data distribution network for the ALICE Online-Offline \(O²\) facility](#), G. Neskovic, T1, 12:00

[Data Analysis using ALICE Run3 Framework](#), G.Eulisse, T6, 11:45

[System simulations for the ALICE ITS detector upgrade](#), S. Nesbo, T2, 12:15

[GPU-based reconstruction and data compression at ALICE during LHC Run3](#), D.Rohr, TX, 14:15

[Running synchronous detector reconstruction in ALICE using declarative workflows](#), M. Richter, TX, 16:30

[Running ALICE Grid Jobs in Containers - A new approach to job execution for the next generation ALICE](#), M.Melnik, T7, 17:45

[Using multiple engines in the Virtual Monte Carlo package](#), B.Volkel, T2, 17:45

Posters:

[Fast and Efficient Entropy Compression of ALICE Data using ANS Coding](#), M.Lettrich, T1

[Space point calibration of the ALICE TPC with track residuals](#), O. Schmidt, T1

[The evolution of the ALICE O² monitoring system](#), A. Wegrzynek, T1



ALICE

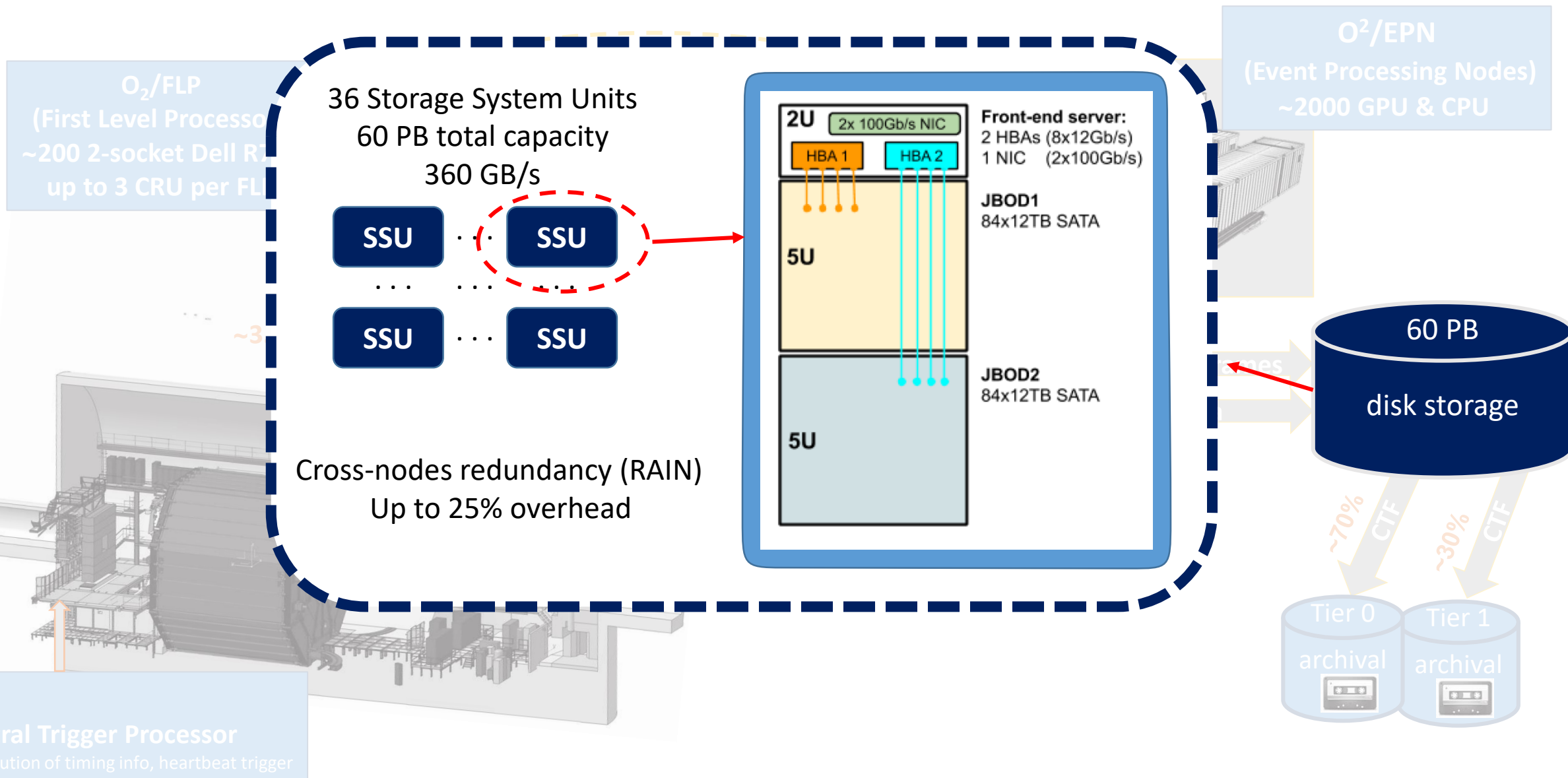


BACKUP





ALICE raw data flow in Run3





CRU and FLP minimum needs per detector



Detector	Throughput GB/s)	Number of read-out boards		Number of FLPs with		
		CRU	CRORCs	2 CRUs or CRORCs	3 CRUs or CRORCs	4 CRUs or CRORCs
ACO (*)	0.01		1	1	1	1
CPV	0.09	1		1	1	1
CTP	0.02	1		1	1	1
DCS		1		1	1	1
EMC (*)	4.00		4	2	2	1
FIT	0.12	2		1	1	1
HMP (*)	0.06		4	2	2	1
ITS	40.0	24		12	8	6
MCH	2.2	32		16	11	8
MFT	10.0	11		6	4	3
MID	0.03	2		1	1	1
PHS (*)	2.0		4	2	2	1
TOF	2.50	4		2	2	1
TPC	570.0	360		180	120	90
TRD	4.0	36		18	12	10
ZDC	0.06	1		1	1	1
Total	635	475	13	247	170	127

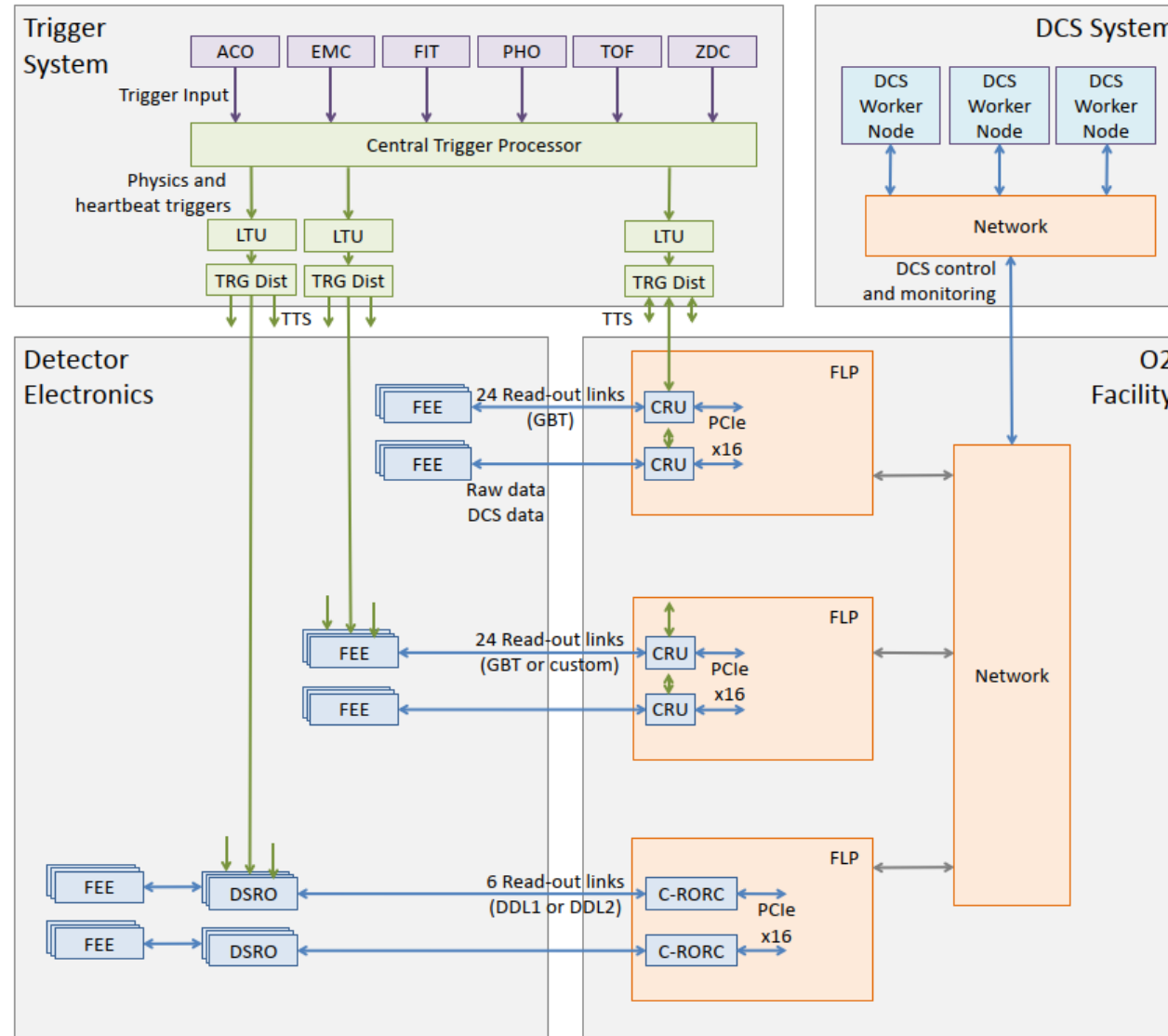


Figure 3.1: Detector read-out and interfaces of the O² system with the trigger, detector electronics and DCS.



QC

