



ALICE

ALICE computing: status and perspectives

**Domenico Elia, I.N.F.N. Bari
Massimo Masera University of Torino and I.N.F.N.**

Outline

- ALICE Computing
 - ➔ present situation
 - ➔ LHC Run 2
- The High Level Trigger: a laboratory for future computing solutions
- The ALICE Upgrade for the LHC Run 3:
 - ➔ the O2 system
 - ➔ a new computing model
 - ➔ a first estimate of the needed resources
 - ➔ INFN activities for Run 3
- Conclusions

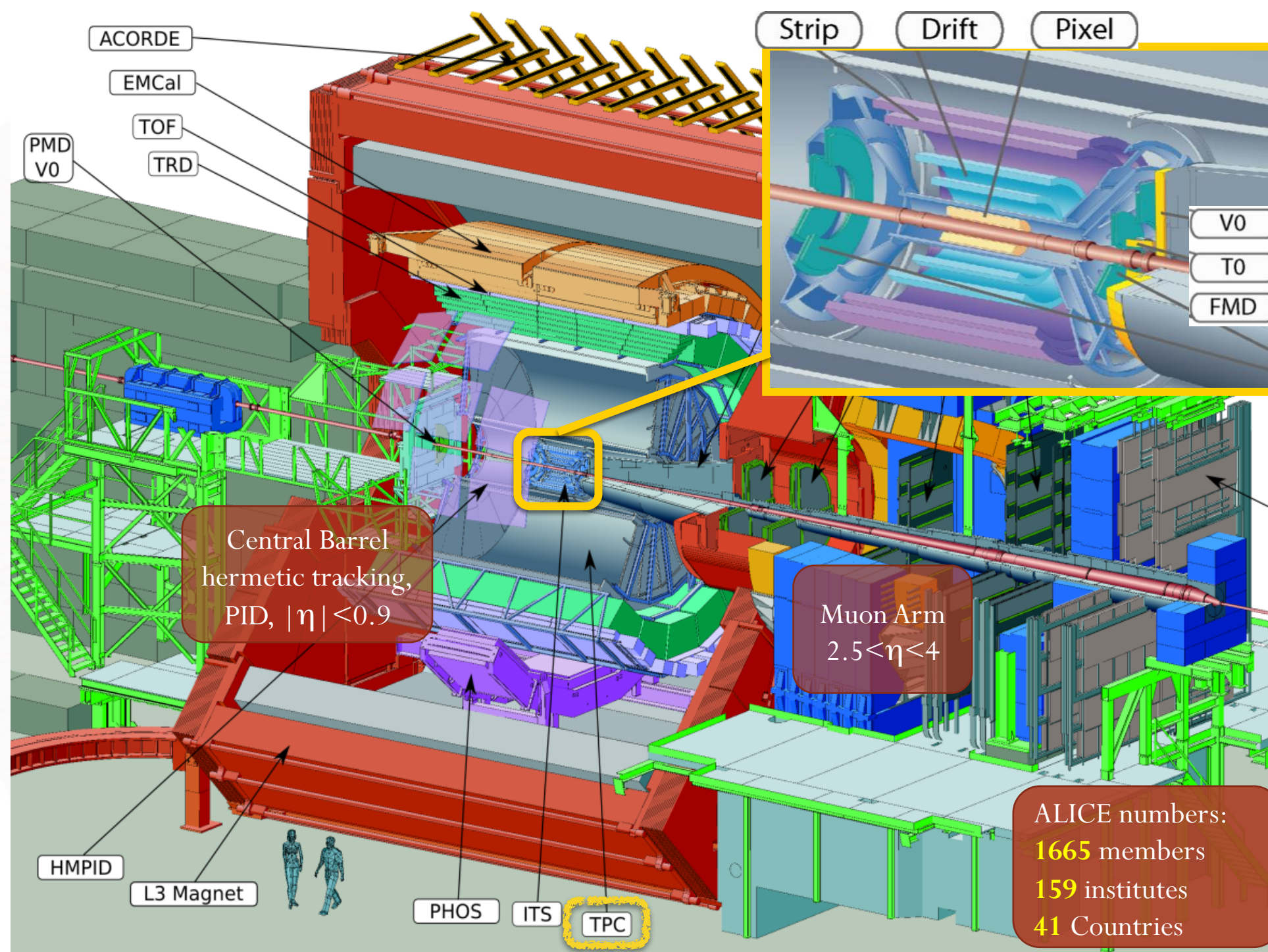
Pb-Pb @ $\sqrt{s} = 2.76$ ATeV
2011-11-12 06:51:12
Fill : 2290
Run : 167693
Event : 0x3d94315a

ALICE

- Designed to reconstruct and identify charged particles in a central rapidity window \rightarrow central barrel down to low transverse momentum ($p_T \sim 100$ MeV/c for pions)
- Main vertexing and tracking detectors: ITS and TPC
- Event recording bandwidth: **1.25 GB/s for Pb-Pb events**
- Data (raw and reconstructed) on permanent storage: few PB/year. Overall stored data:
 - » **tape: ~45 PB**
 - » **storage: ~55 PB**
- Reconstruction: almost completely offline

Acronyms:

- ITS - Inner Tracking System
- TPC - Time Projection Chamber



ALICE numbers:
1665 members
159 institutes
41 Countries

ALICE: present status and LHC Run 2



ALICE

HEICE

Pb-Pb @ $\sqrt{s} = 2.76$ ATeV

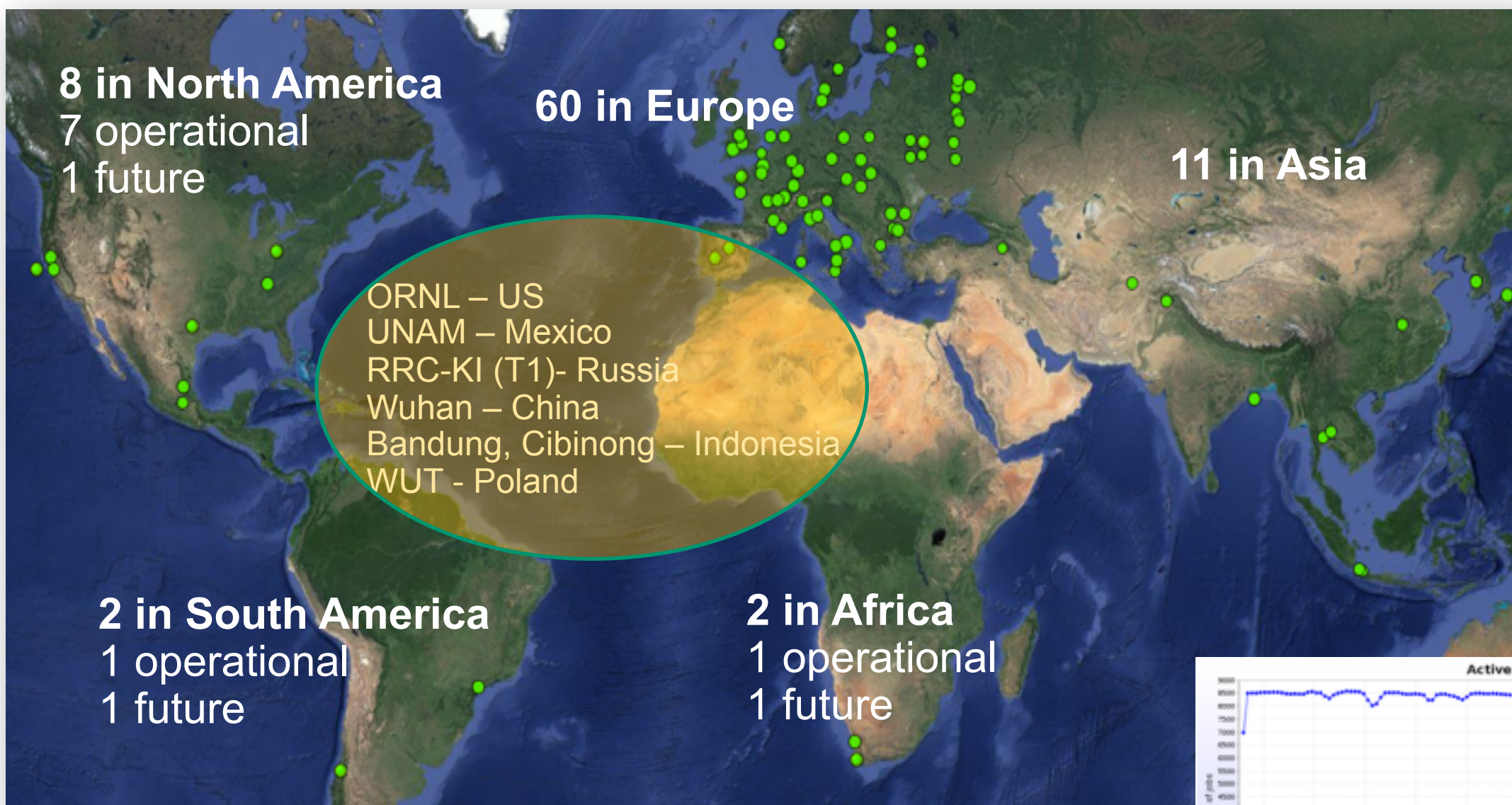
2011-11-12 06:51:12

Fill : 2290

Run : 167693

Event : 0x3d94315a

ALICE Grid



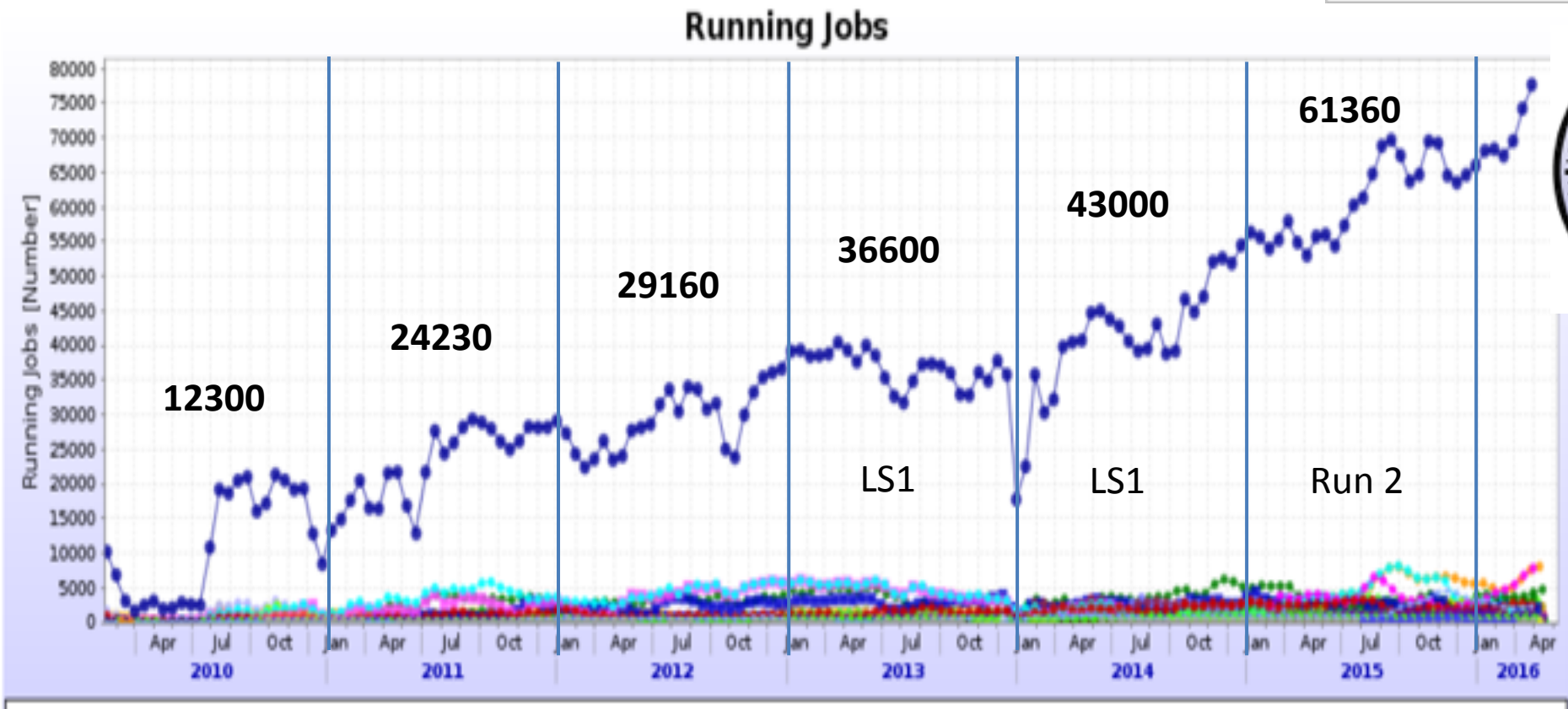
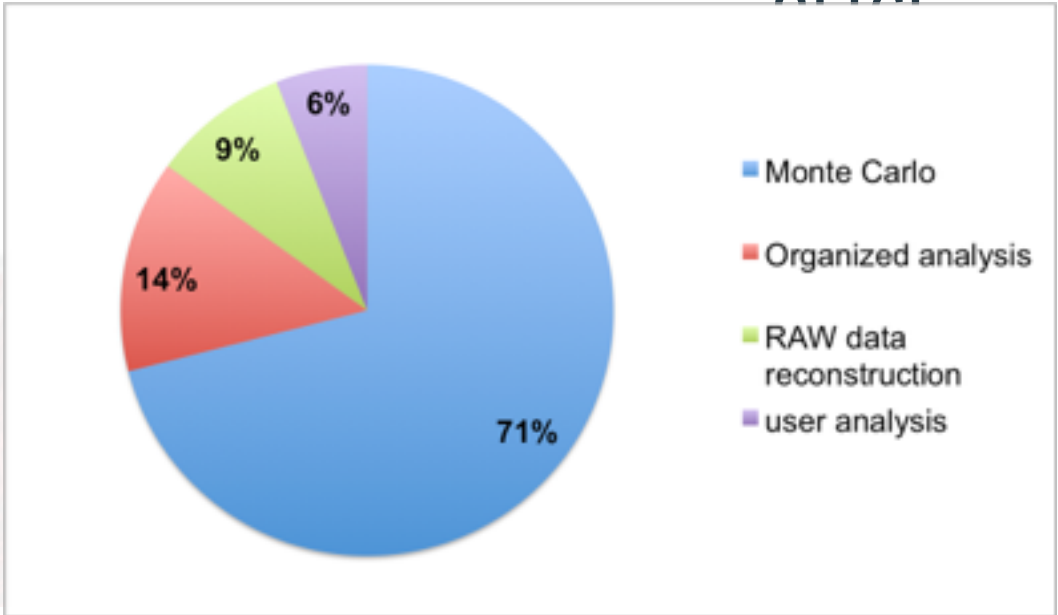
- ➔ Computing centres mostly in Europe
- ➔ There are new entries in 2015-2016
- ➔ The High Level Trigger Farm is now used, when not in run, for offline activities. It is a fully virtual site





CPU resources evolution

- ➔ Steady growth of the number of active jobs
- ➔ More than 280 million jobs have been successfully run so far
- ➔ System scaled from 500 up to 100000 concurrently running jobs
- ➔ Scheduled analysis is now prevailing on chaotic (user) analysis: **+60% wrt 2014 ⇒ better efficiency**



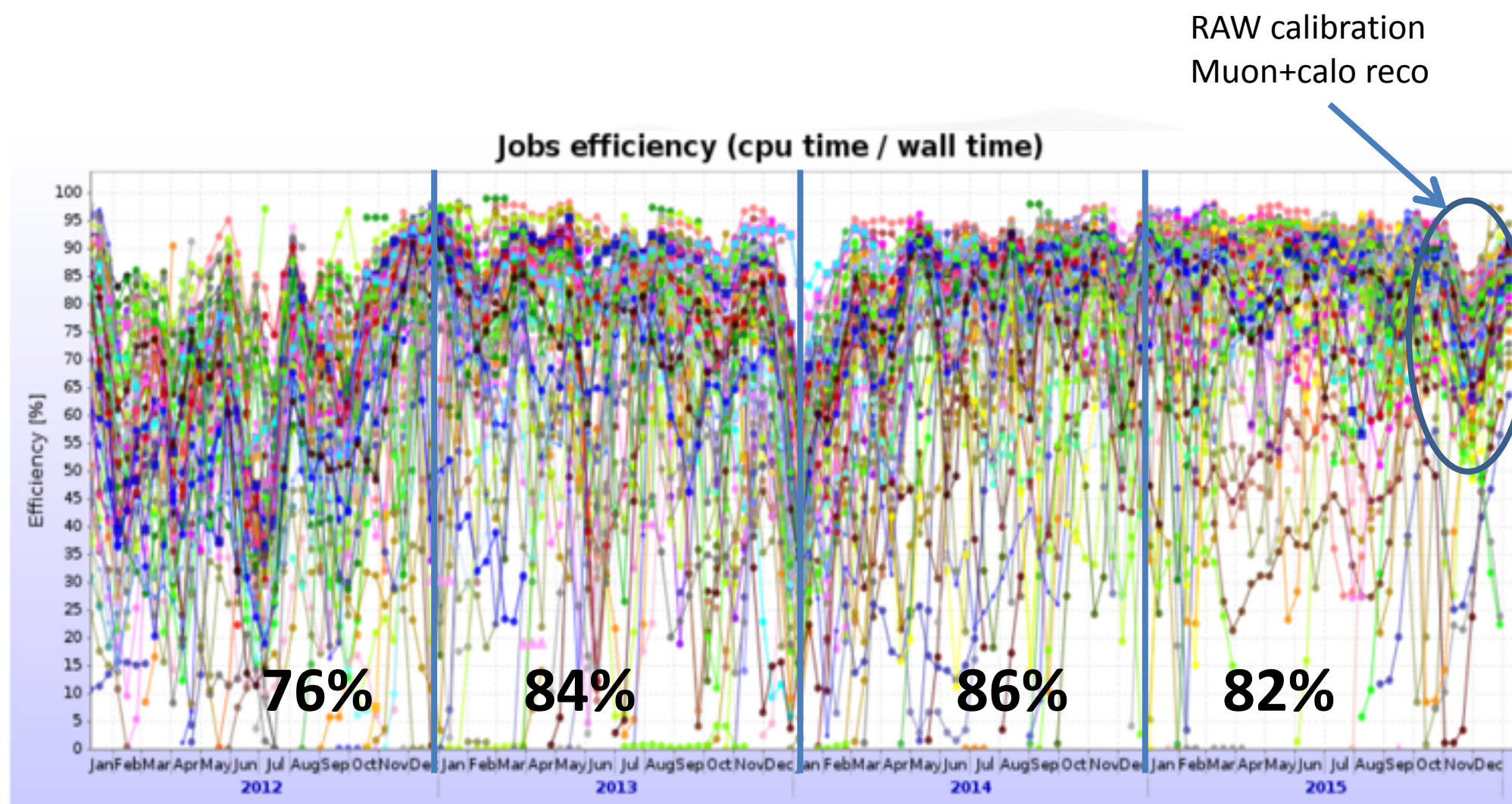
Year	Chaotic analysis (# jobs)	Scheduled analysis (# jobs)
2012	6900	3000
2015	3900	9100

Number of individual users: ~450 flat in the last 4 years

Year on year increase

↑ +97% ↑ +20% ↑ +26% ↑ +17% ↑ **+43%**

Job efficiency



Year on year change

↑
+8%

↑
+2%

↑
-4%

Steadily above 80% in the last three years

Pb-Pb @ $\sqrt{s} = 2.76$ ATeV
2011-11-12 06:51:12
Fill : 2290
Run : 167693
Event : 0x3d94315a

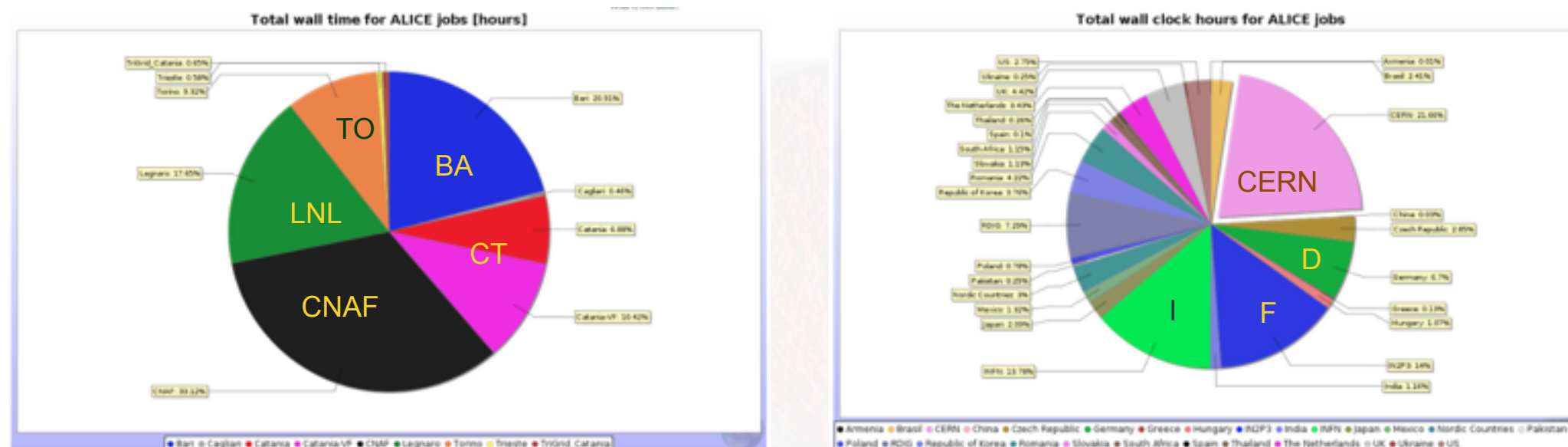
Job efficiency evolution

- Efficiency has been rather flat for the last 4 years
- Decrease in 2015 is largely due to specific RAW data reconstruction cycles
 - Increase of scheduled analysis w.r.t. individual analysis helps to compensate the lower RAW reconstruction efficiency
- Specific effort to increase the efficiency of offline calibration tasks
- **~85% efficiency** is a good benchmark value for ALICE

Pb-Pb @ \sqrt{s} = 2.76 ATeV
2011-11-12 06:51:12
Fill : 2290
Run : 167693
Event : 0x3d94315a

Resources for Run 2

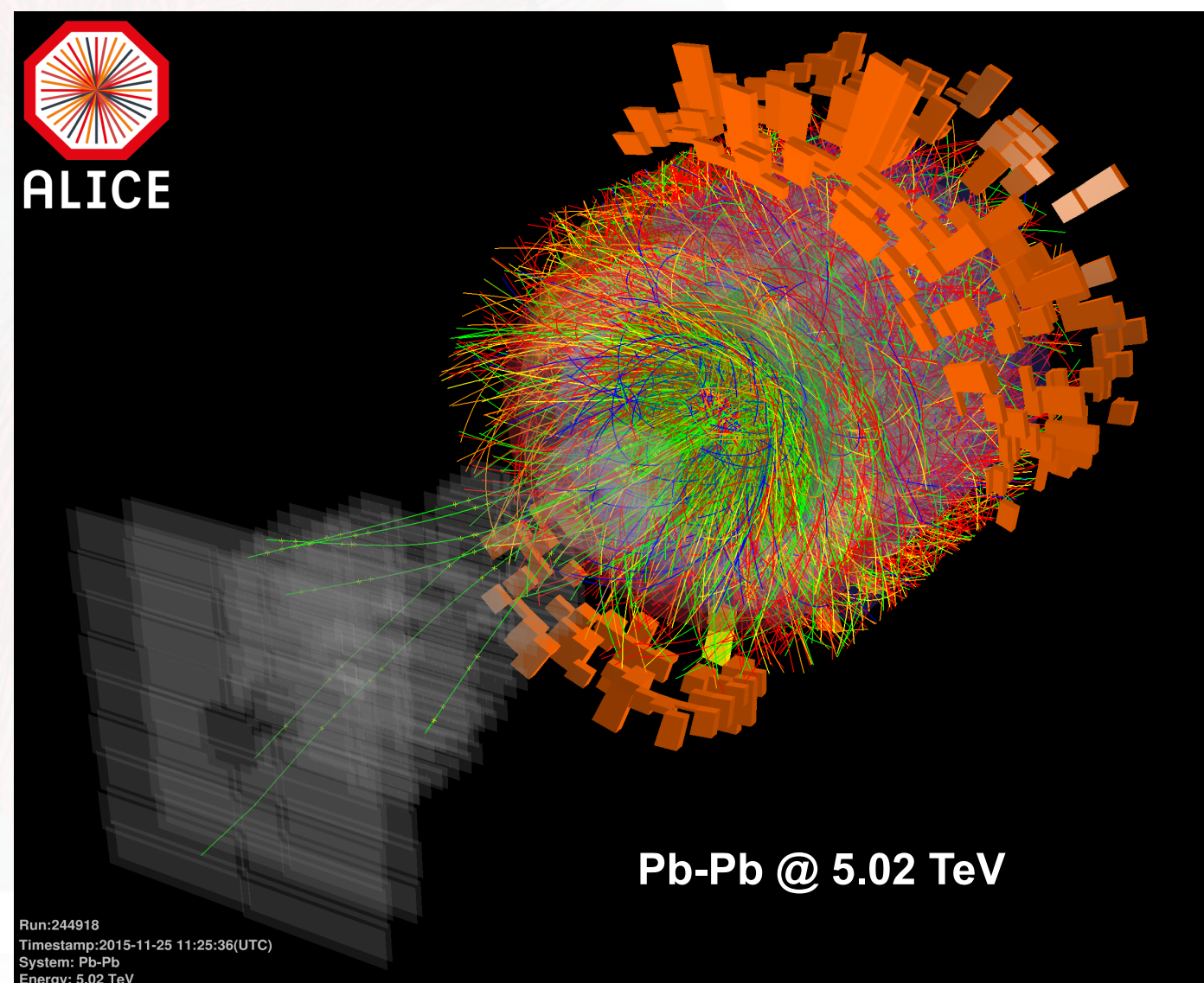
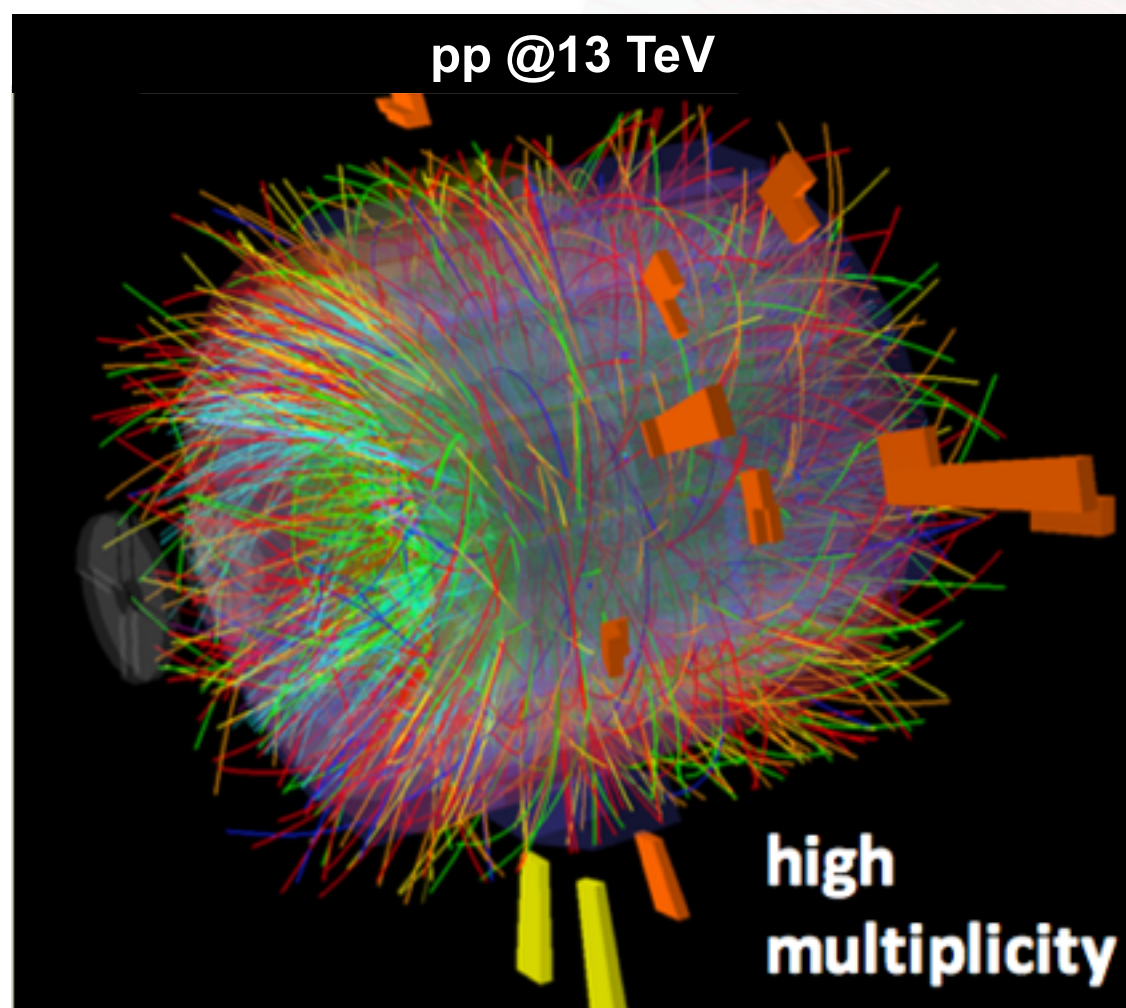
Sharing during the last year



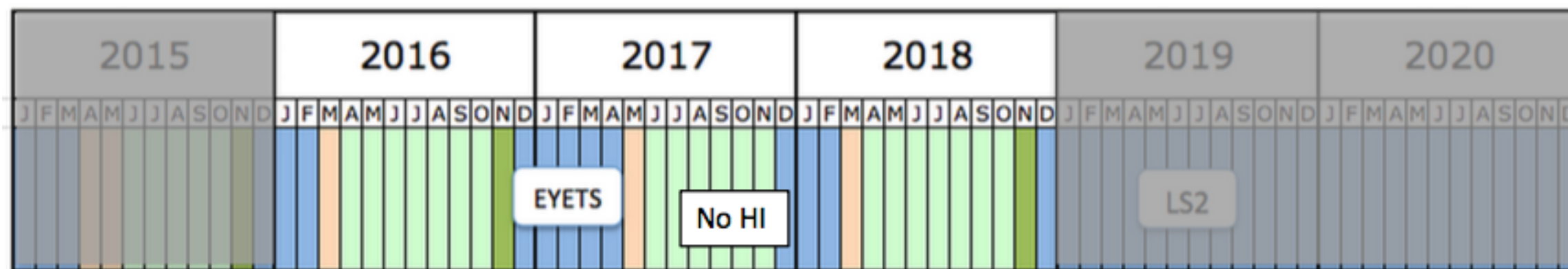
		2016	2017	2018
CPU (kHS06)	T0	224	255	318
	T1 + T2	403	496	604
	INFN (T1 + T2)	76.2	93.7	114.2
DISK (PB)	T0	16.8	21.4	27.7
	T1 + T2	47.5	56.3	70.7
	INFN (T1 + T2)	8.9	10.6	13.4
TAPE (PB)	T0	26.3	34.4	45.9
	T1	20.3	28.4	39.9
	INFN (T1)	6.5	9.1	12.8

Run 2 - The first year

Year	System	E [TeV]	Lumi [cm ⁻² s ⁻¹]	Rate [kHz]	Time
2015	pp	13	5×10^{30}	300	14w
	Pb-Pb	5.02	1×10^{27}	8	3w
	pp-ref	5.02	1×10^{30}	50	4d

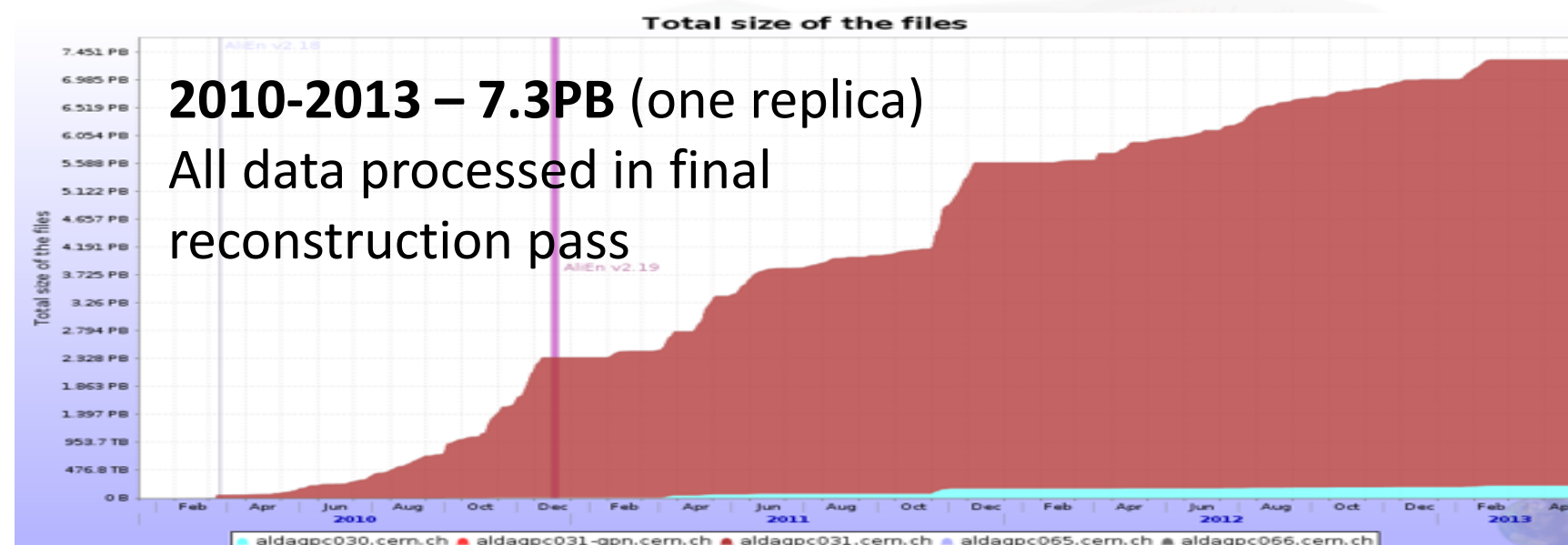


RUN 2 Overview

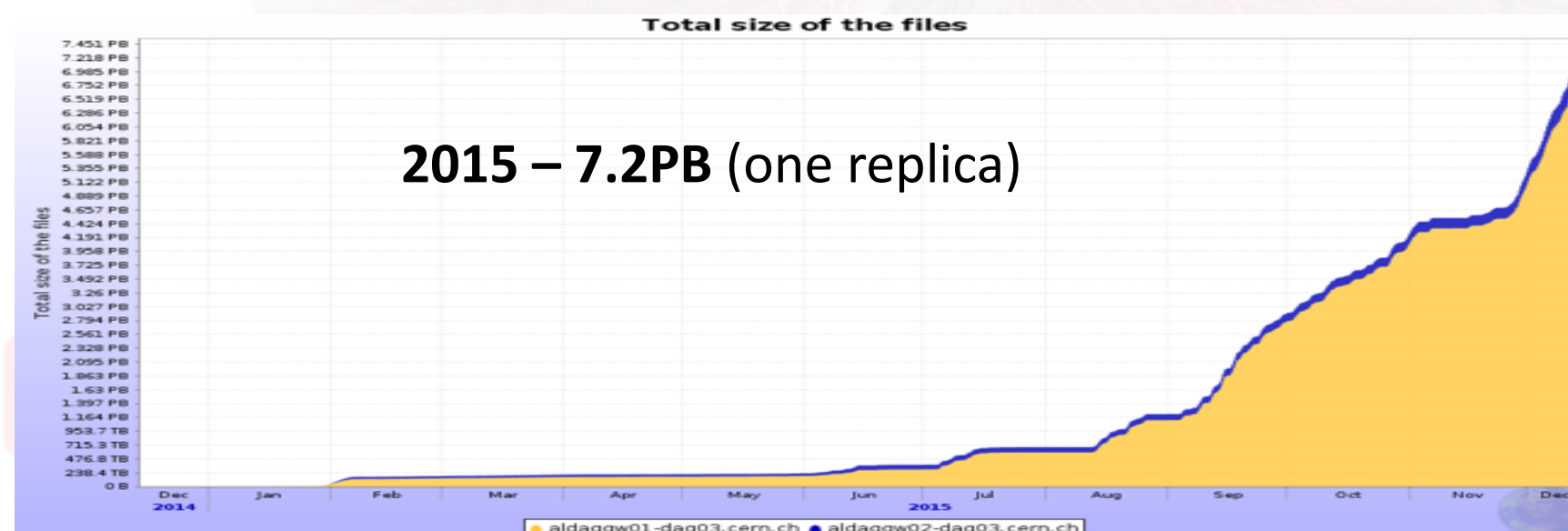


Year	System	E [TeV]	Lumi [cm ⁻² s ⁻¹]	Rate [kHz]	Time
2016	p-p	13	5x10 ³⁰	300	28w
	p-Pb	5.02(8.16)	1x10 ²⁸ + 1x10 ²⁹	20(mB)/200	4w
2017	p-p	13	5x10 ³⁰	300	24w
	pp-ref(?)	5.02	1x10 ³⁰	50	7d
2018	p-p	13	5x10 ³⁰	300	28w
	Pb-Pb	5.02	1x10 ²⁷	8	4w
	pp-ref(?)	5.02	1x10 ³⁰	50	7d

2015 RAW data collection



Run 1



Run 2

Current activities (1)

- Monte Carlo: about 150 individual MC cycles in the last year
 - Physics papers in preparations
 - *First physics* analysis of 2015 data
- RAW data processing:
 - Code improved to reduce the memory consumption: now 2 GB/job
 - 2015 data reconstructed partially
 - Distortions in the TPC detector occur in runs with high interaction rate
 - Specific corrections have been developed and are currently being validated

Pb-Pb @ $\sqrt{s} = 2.76$ ATeV
2011-11-12 06:51:12
File : 2290
Run : 167693
Event : 0x3d94315a

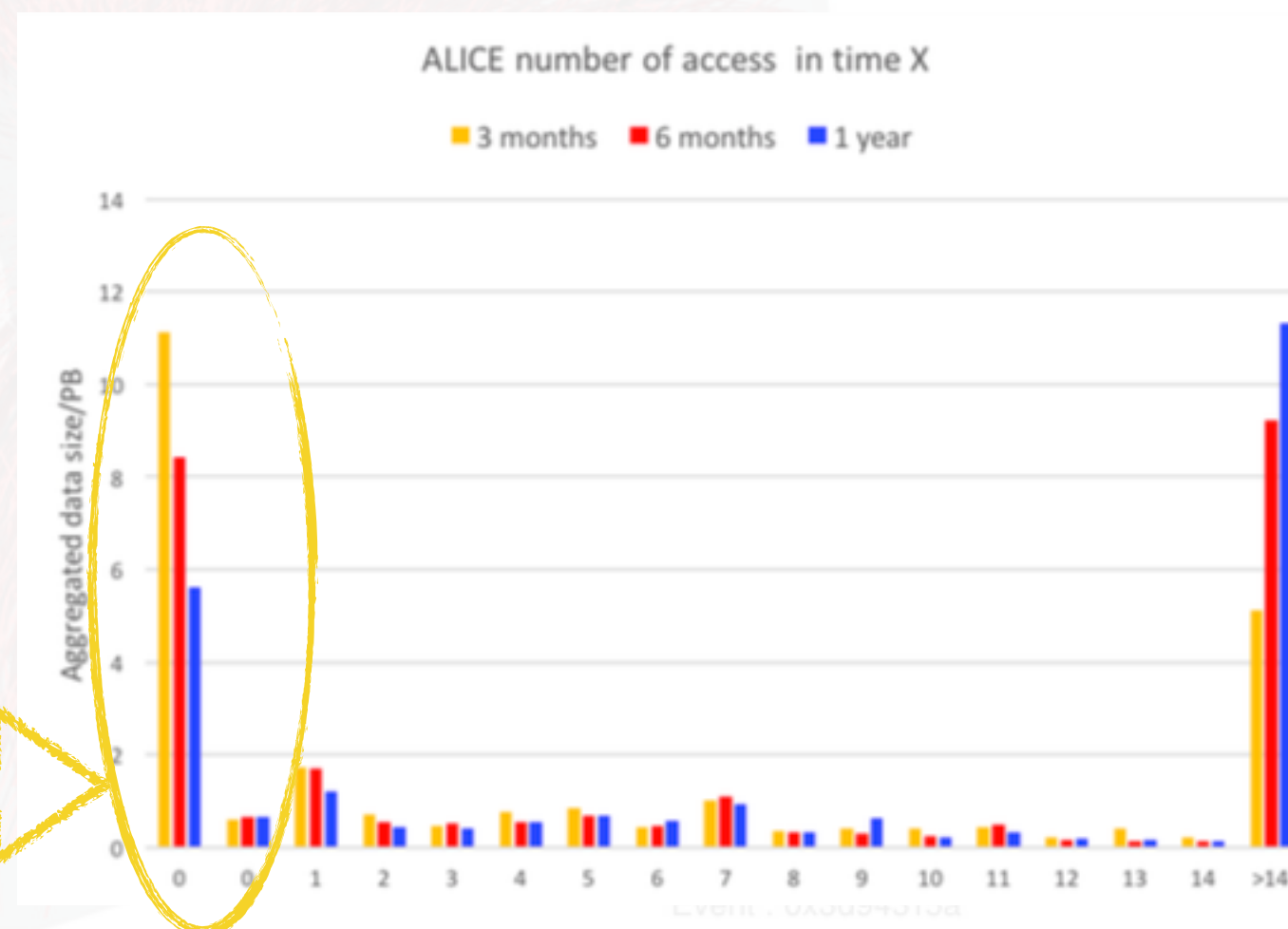
Current activities (2)

- ➔ Replication policy has been changed to cope with the available storage
 - ✓ ESD are replicated only once
- ➔ Global disk space needed for 2015 data processing:
 - ✓ 5 ÷ 6 PB (RAW+MC)
 - ✓ barely feasible with the expected resources

Popularity and cleanup

- ➔ Really old MC productions have been removed
- ➔ Second ESD replica for low access productions have been removed
- ➔ A list of dataset to be “sanitized” is kept and available

- ➔ Volumes of data vs number of accesses in X=3 - 6 - 12 months
- ➔ The first bin is for data created before the period X began and **not** accessed during that period



The High Level Trigger: a starting point for the future LHC Run 3



ALICE

HEICE

Pb-Pb @ $\sqrt{s} = 2.76$ ATeV

2011-11-12 06:51:12

Fill : 2290

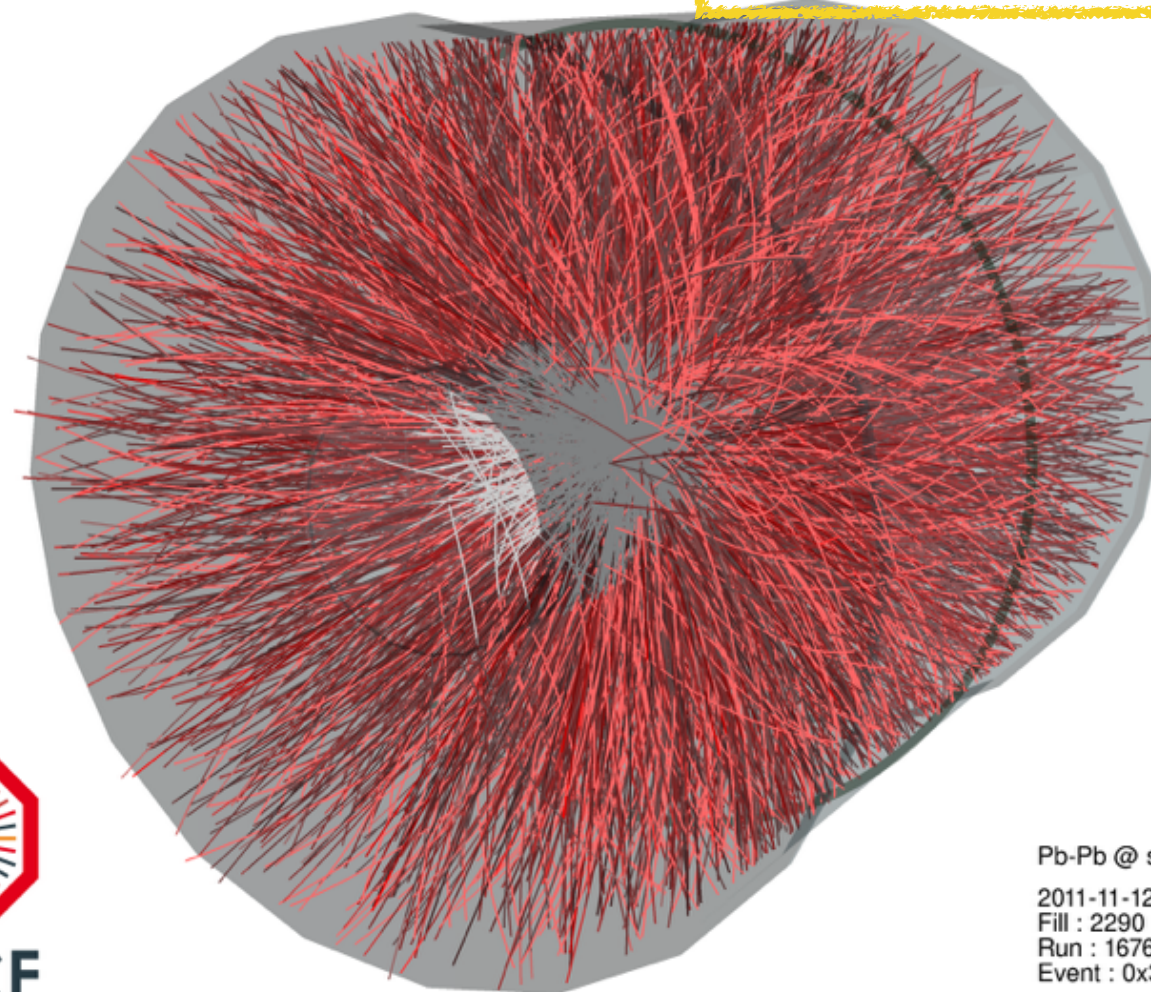
Run : 167693

Event : 0x3d94315a

ALICE - Data reconstruction challenges

Charged particle multiplicity: $dN/dy \sim 1400$
Number of track in TPC up to 20000
Number of TPC clusters $\sim 10^6$

- The High Level Trigger (HLT) system is capable to carry out the event reconstruction in real time.
- Designed for online event selection.
- It is currently also used to perform the local reconstruction (clustering) of the TPC \rightarrow TPC clusters are stored instead of raw data to reduce the event size
- The track reconstruction in a high multiplicity environment is a challenging issue

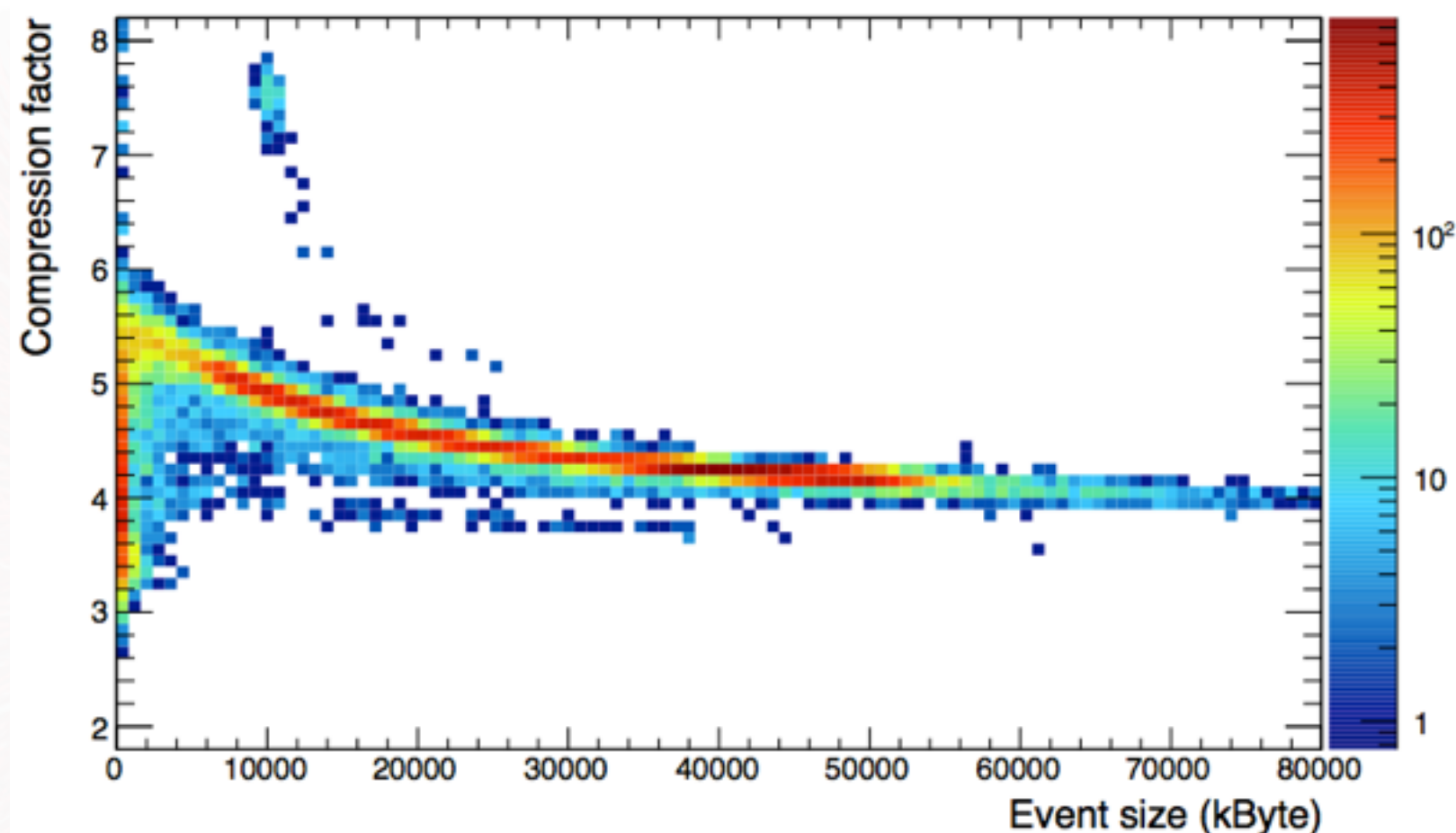


Pb-Pb @ $\sqrt{s} = 2.76$ ATeV
2011-11-12 06:51:12
Fill : 2290
Run : 167693
Event : 0x3d94315a

Reconstructed charged particle trajectories (tracks) in the ITS (white) and TPC detectors for a Pb-Pb event

ALICE - Data reconstruction challenges

- The High Level Trigger (HLT) system is capable to carry out the event reconstruction in real time.
- Designed for online event selection.
- It is currently also used to perform the local reconstruction (clustering) of the TPC -> **TPC clusters are stored instead of raw data to reduce the event size**
- The track reconstruction in a high multiplicity environment is a challenging issue



TPC cluster compression factor as a function of the event size (Pb-Pb data - 2011)

Pb-Pb @ $\sqrt{s} = 2.76$ ATeV

2011-11-12 06:51:12

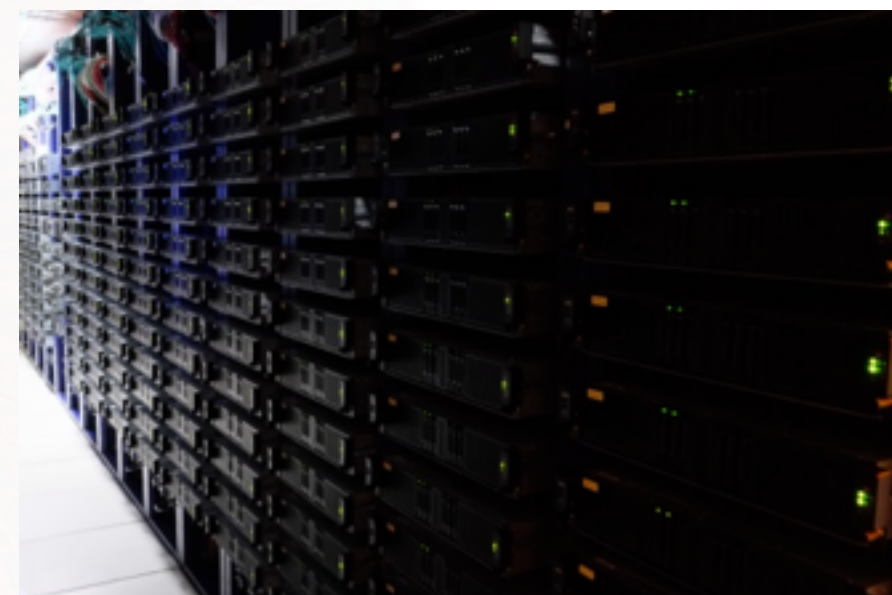
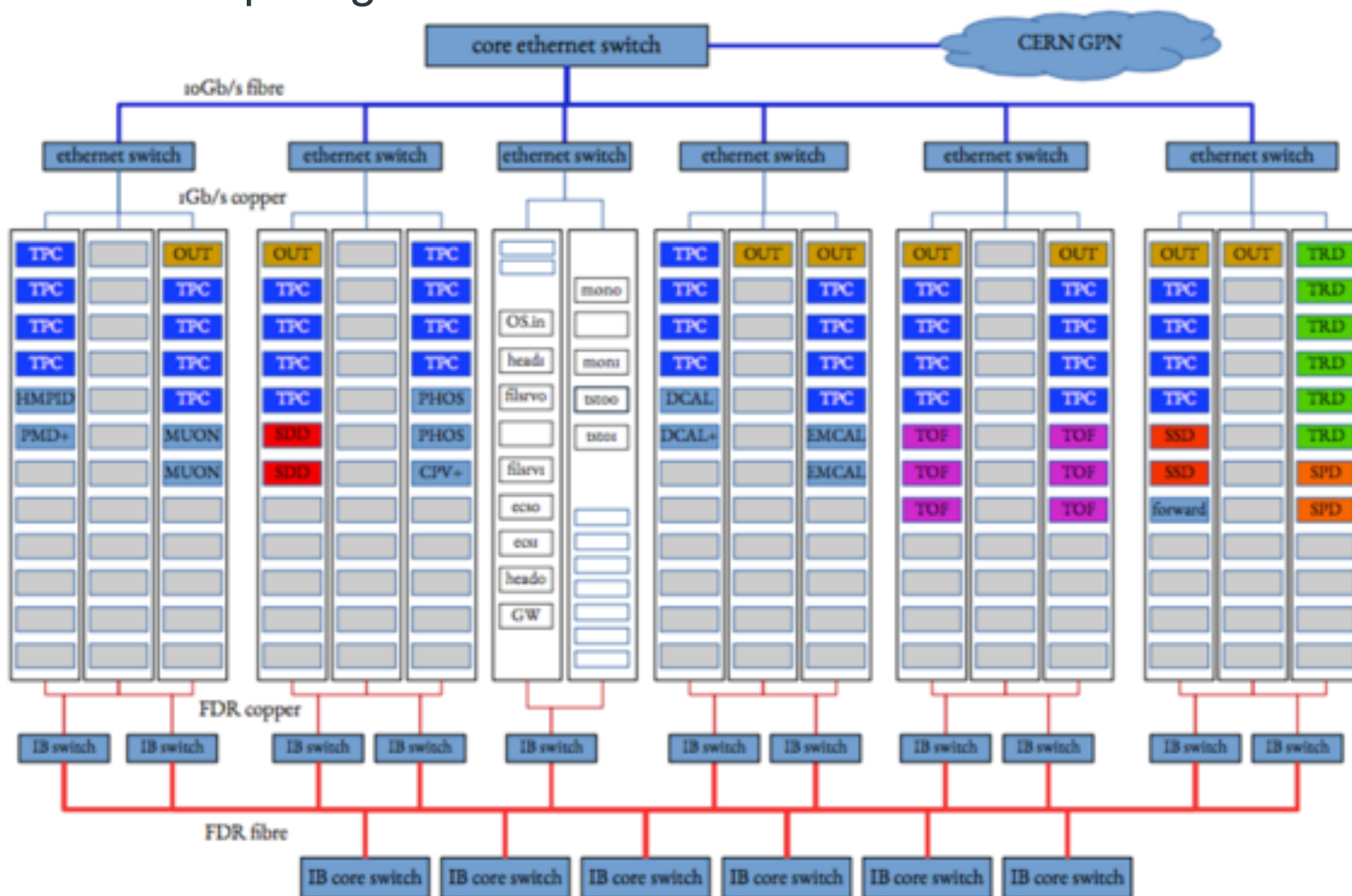
Fill : 2290

Run : 167693

Event : 0x3d94315a

Online computing in ALICE: the HLT

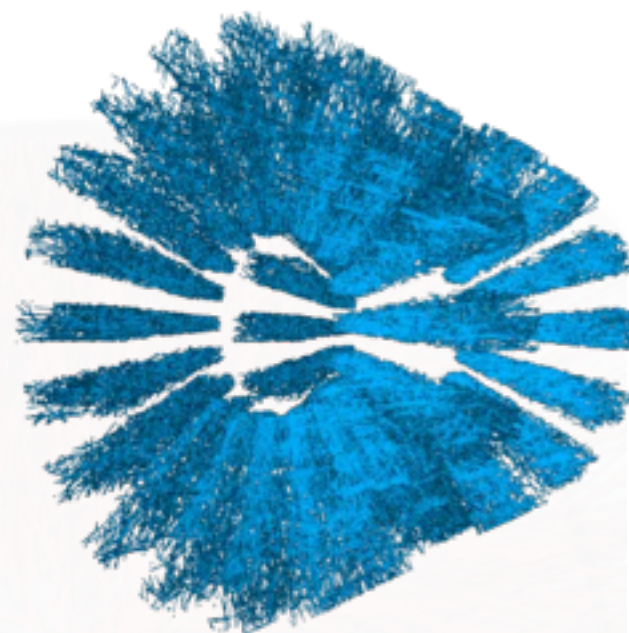
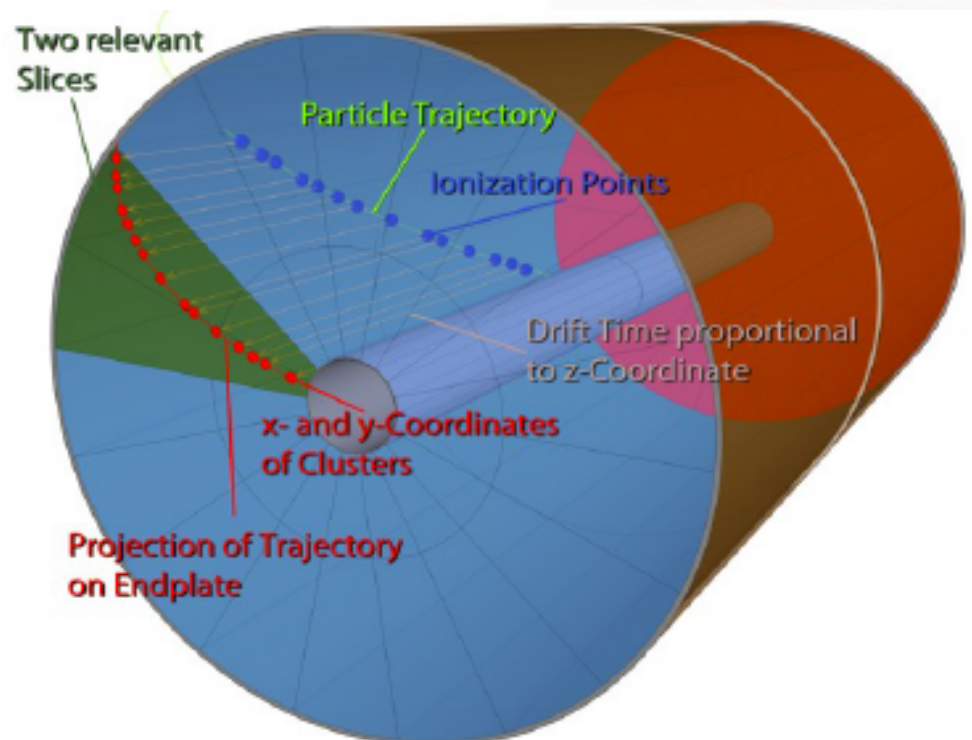
- The HLT computing farm has to process the events in real time.
- Event rate 300 Hz Pb-Pb (2 kHz pp)
- Data rate ~30 GB/s
- The cluster finder algorithm for the TPC is coded on FPGAs connected to sub-elements of the detector, i.e. before the event building (in blue)
- Up to 20000 tracks in the TPC/event
- Up to 159 clusters/track
- Tracking: combining clusters to reconstruct particle trajectories
 - » high combinatorics
 - » computing intensive



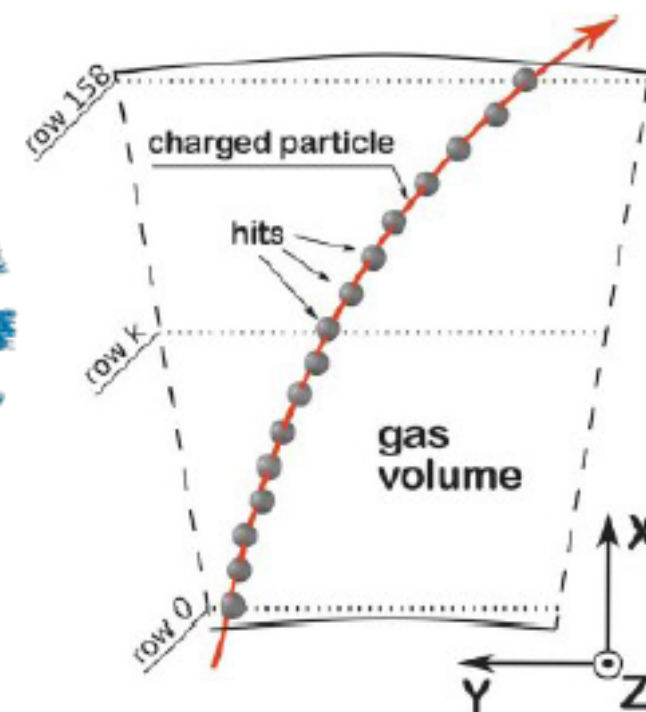
Layout of the HLT farm. Worker nodes are in gray. In colour, the nodes equipped with Readout receiver cards, hosting a Xilinx Virtex FPGA

HLT: TPC tracking

- TPC tracking with the HLT is the most relevant example of HPC in ALICE
- It is the basis for the reconstruction code in the forthcoming LHC Run 3, with an upgraded ALICE apparatus at a much higher interaction rate
- The TPC volume is split in 36 sectors: tracking is done in each sector individually.



Each sector, 159 rows



- The radial and azimuthal coordinates of the clusters are measured by charge collection in 159 rows.
- Inner radius: 85 cm
- Outer radius: 250 cm
- The coordinate along the beam axis is measured via the drift time

Pb-Pb @ $\sqrt{s} = 2.76$ ATeV

2011-11-12 06:51:12

Fill : 2290

Run : 167693

Event : 0x3d94315a

HLT: TPC tracking

Neighbour finder

- For each hit at row k , the best pair of neighbouring hits from row $k+1$ and $k-1$ is found (best=straight line)

Evolution

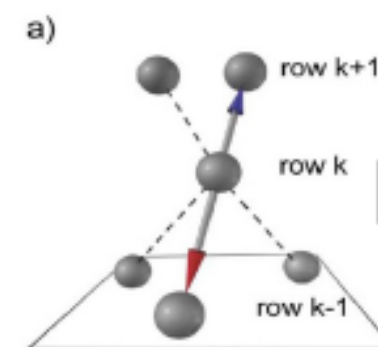
- Reciprocal links are determined and saved

Tracklet reconstruction

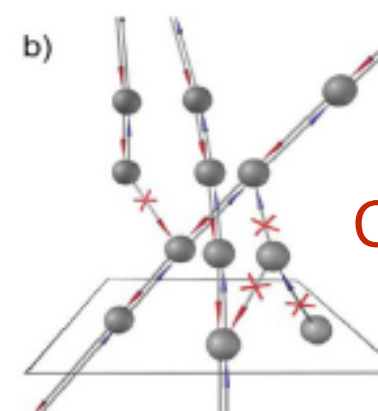
- Tracklets are created following hit-to-hit links; Kalman filter to fit geometrical trajectories

Tracklet selection

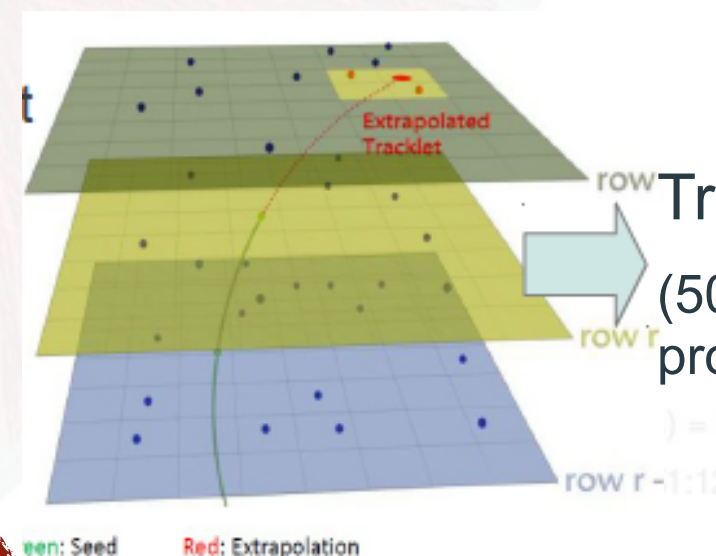
- In case of tracks with overlapping parts, the longest is kept



Every hit in parallel



Cellular automaton

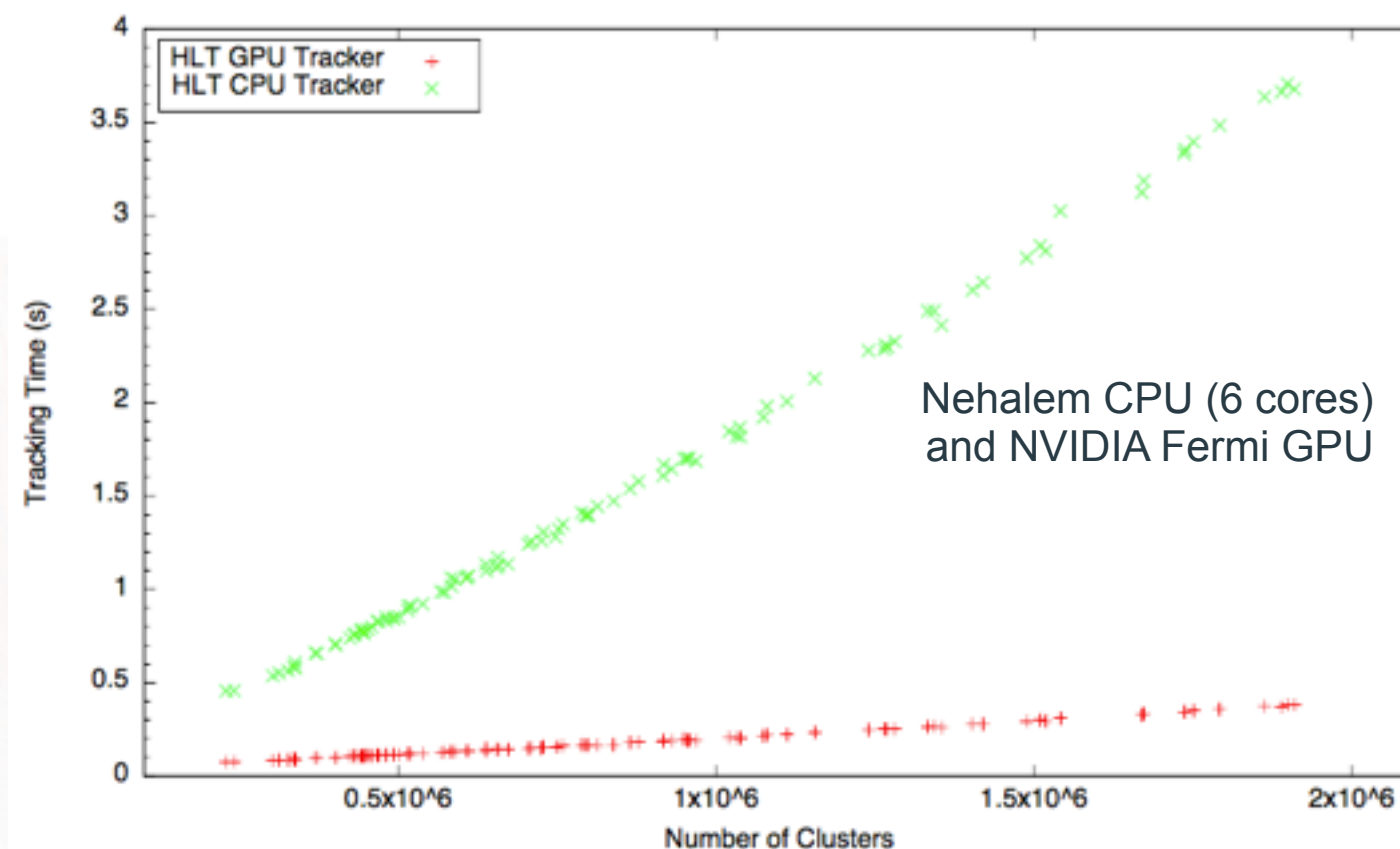


Track following
(50% of the total processing time)

Event: 0x3d94315a

HLT: TPC tracking

- GPU tracker is about 10 times faster w.r.t. the CPU version
- Track merging and fitting is done on CPU (data transfer time cancels the GPU speed gain)
- The whole tracking is done with 1 GPU +3 CPU cores
- Performance of 1GPU+3 CPU cores ~ 27 CPU cores



Task	How	Locality	Description	Time	Where
Seeding	Cellular Automaton	Local	Find track candidates (3-10 clusters)	~30%	GPU or CPU
Track following	Kalman filter	Sector	Fit parameters to candidate, find full track segment in one sector via track following with simplified Kalman filter (e.g. constant B-field, y and x uncorrelated)	~60%	
Track merging	Combinatorics	Global	merge track segments within a sector and between sectors	~2%	CPU (GPU version exists - not used)
Track fit	Kalman filter	Global	Full track fit with full Kalman filter (polynomial approximation of B-field)	~8%	

Code management

- The HLT farm is **heterogeneous**: CPU and GPU accelerators are available
- The GPU code may be vendor bound: e.g. Cuda in case of Nvidia cards
- A CPU version of the code must be maintained to be able to run the same reconstruction code on other facilities
 - » on WLCG infrastructure for MC data, for instance.
- For the HLT TPC tracking, **CPU and GPU codes share common source files.**
- Specialized wrappers for CPU, Cuda **and OpenCL** are provided. They include the common files.
- The fraction of **common source code is above 90%**
- The experience gained with the HLT is the basis for the new **Online+Offline (O²)** infrastructure that will be used for LHC Run 3 (starting from 2020)

LHC Run 3



ALICE

HEICE

Pb-Pb @ $\sqrt{s} = 2.76$ ATeV

2011-11-12 06:51:12

Fill : 2290

Run : 167693

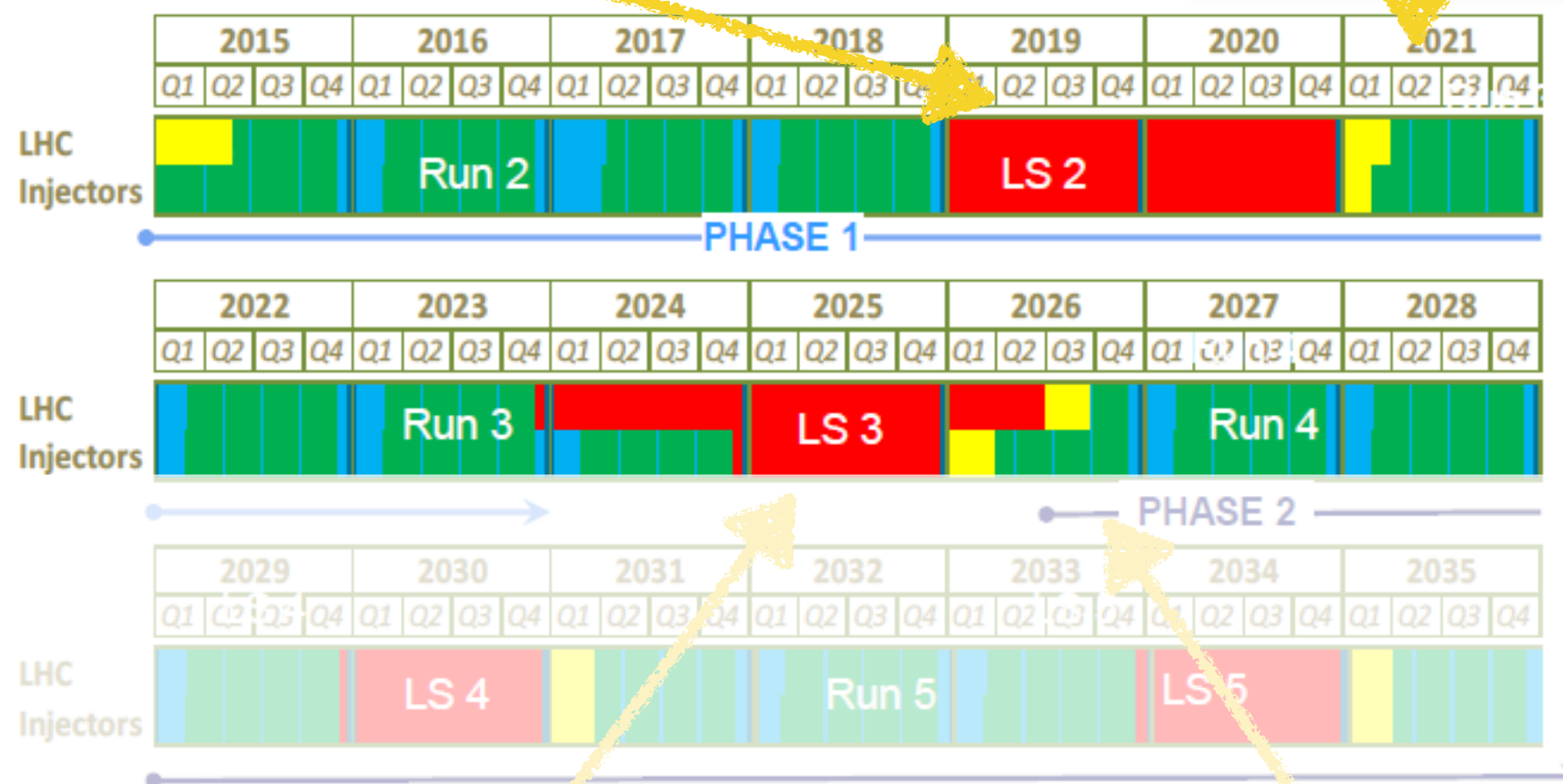
Event : 0x3d94315a

LHC Run 3 and ALICE Upgrade

PHASE I Upgrade

ALICE, LHCb major upgrade
ATLAS, CMS ,minor' upgrade

Heavy Ion Luminosity
from 10^{27} to 7×10^{27}



PHASE II Upgrade

ATLAS, CMS major upgrade

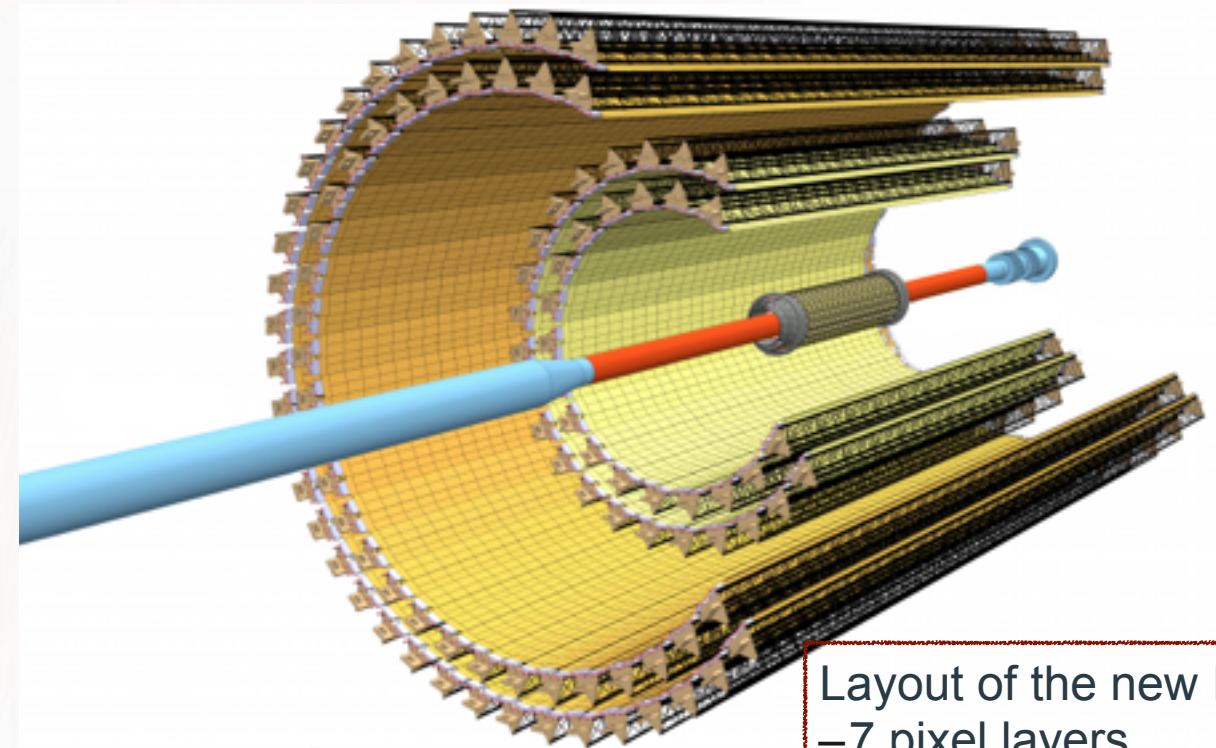
HL-LHC, pp luminosity
from 10^{34} (peak) to 5×10^{34} (levelled)

ALICE Upgrade

- After the second LHC Long Shutdown (2019-20), new conditions are expected for the subsequent Run3:
 - » Expected Pb-Pb peak interaction rate: 50 kHz (now it is 8 kHz)
- Presently ALICE readout rate is limited to ~1 kHz
- Goal for Run3:
 - » no reliable triggering strategies for several physics channels —> increase the readout rate to **50 kHz**
 - » improve pointing resolution both in the barrel (**new ITS**) and in the Muon Arm (new Muon Forward Tracker)

The ALICE upgrade requires major improvements for the TPC and other detectors in order to increase the readout rate

Capability of reducing online the data volume delivered by the detectors, since the expected integrated luminosity is $> 10 \text{ nb}^{-1}$ for Pb-Pb (x100 w.r.t. Run 1)



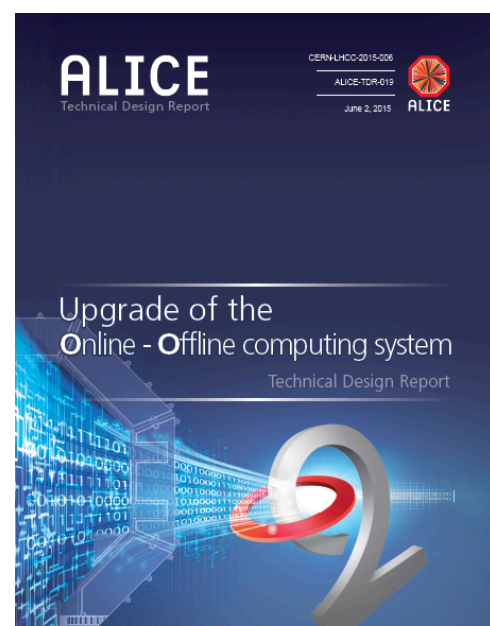
Layout of the new ITS:
– 7 pixel layers
– 10 m² of silicon
– 12.5 Gpixel

ALICE Upgrade: the O² system

- The expected data rate for Pb-Pb collisions at 50 kHz is ~1.1 TB/s
- The TPC alone accounts for 1 TB/s
- The O² project aims to integrate in a single infrastructure the present DAQ, HLT and Offline (for the reconstruction part) systems

Detector	Average event size (MB)	Data rate for Pb-Pb @ 50 kHz (GB/s)
TPC	20.7	1012
ITS	0.8	40
TRD	0.5	20
MFT	0.2	10
Others	0.3	12.2
Total	22,5	1094,2

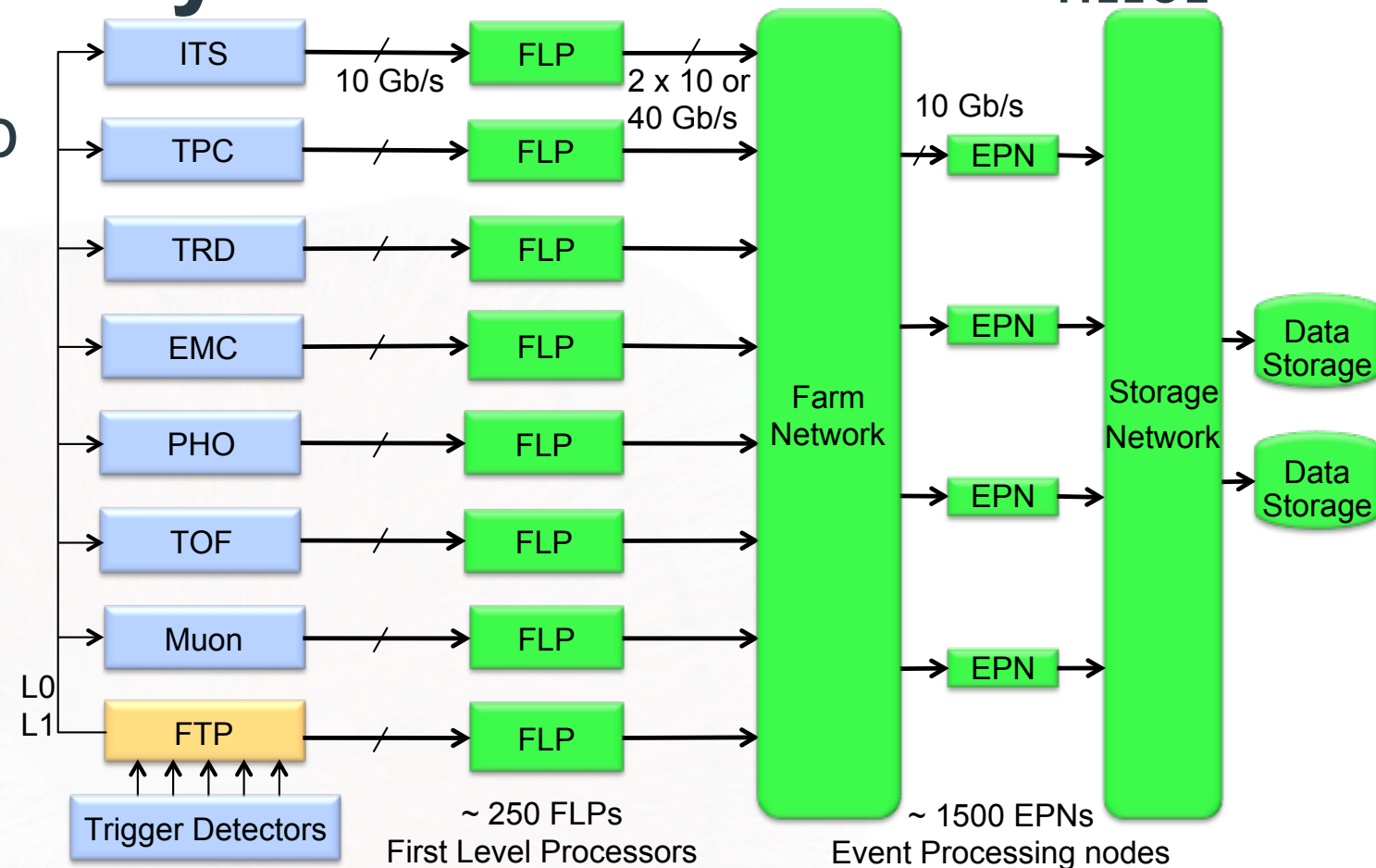
- The data volume coming from the detectors must be substantially reduced before sending the data to the mass storage.
- Online processing is the only option
- The computing strategy must rely on a heterogeneous architecture to match the interaction rate:
 - » ~250 FLP worker nodes (First Level Processors) equipped with FPGA
 - » ~1500 EPN worker nodes (Event Processing Nodes) equipped with GPU
 - » yearly amount of data (2020, 2021): 54 PB



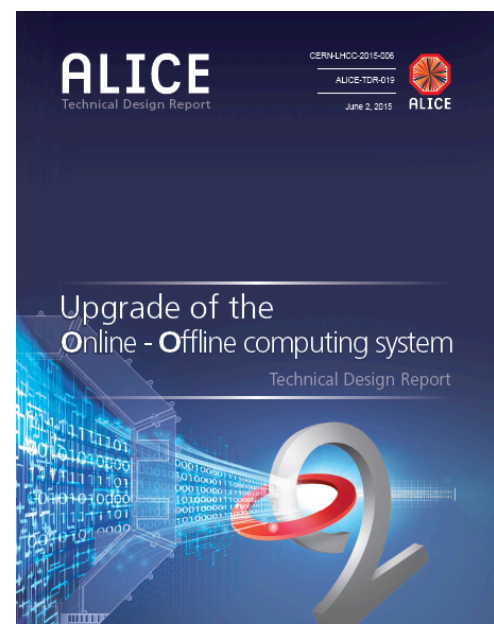
ALICE Upgrade: the O² system



- The expected data rate for Pb-Pb collisions at 50 kHz is ~1.1 TB/s
- The TPC alone accounts for 1 TB/s
- The O² project aims to integrate in a single infrastructure the present DAQ, HLT and Offline (for the reconstruction part) systems



- The data volume coming from the detectors must be substantially reduced before sending the data to the mass storage.
- Online processing is the only option
- The computing strategy must rely on a heterogeneous architecture to match the interaction rate:
 - » ~250 FLP worker nodes (First Level Processors) equipped with FPGA
 - » ~1500 EPN worker nodes (Event Processing Nodes) equipped with GPU
 - » yearly amount of data (2020, 2021): 54 PB



Online data volume reduction

- The impressive reduction factor that can be obtained for the TPC is based on:
 - » zero suppression
 - » clustering and compression
 - » removal of clusters non associated to interesting particle tracks (e.g. very low momentum electrons)
 - » data format optimization
- Largely based on the present HLT results

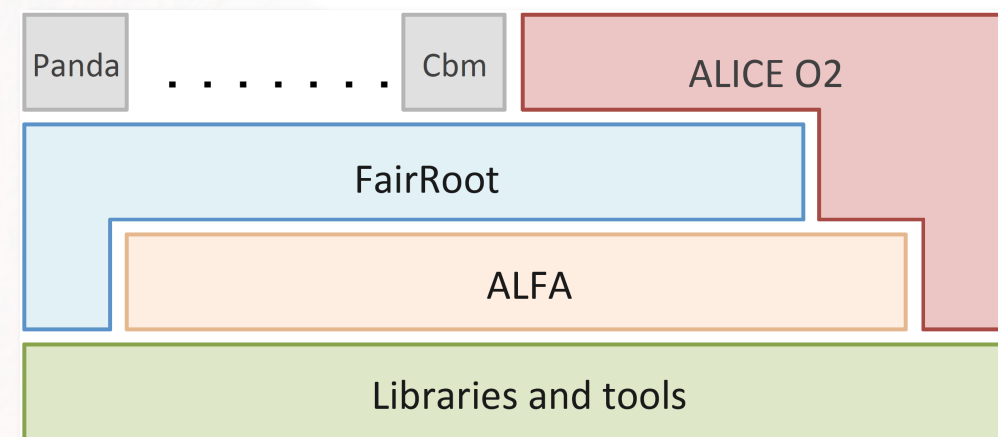
Still uncertainties for the ITS:

- » The contribution from noisy clusters is unknown: here a pessimistic estimate of a **probability of 10^{-5} per pixel** has been made
- » If full synchronous reconstruction will be feasible a higher reduction factor will be achieved (noise removal)

Detector	Data rate for Pb-Pb @ 50 kHz (GB/s)	Compressed data rate (GB/s)	Data reduction
TPC	1012	50	20.2
ITS	40	26 (8)	1.5 (5)
TRD	20	3	6.7
MFT	10	5	2
Total	1082	84 (66)	12.9 (16.4)

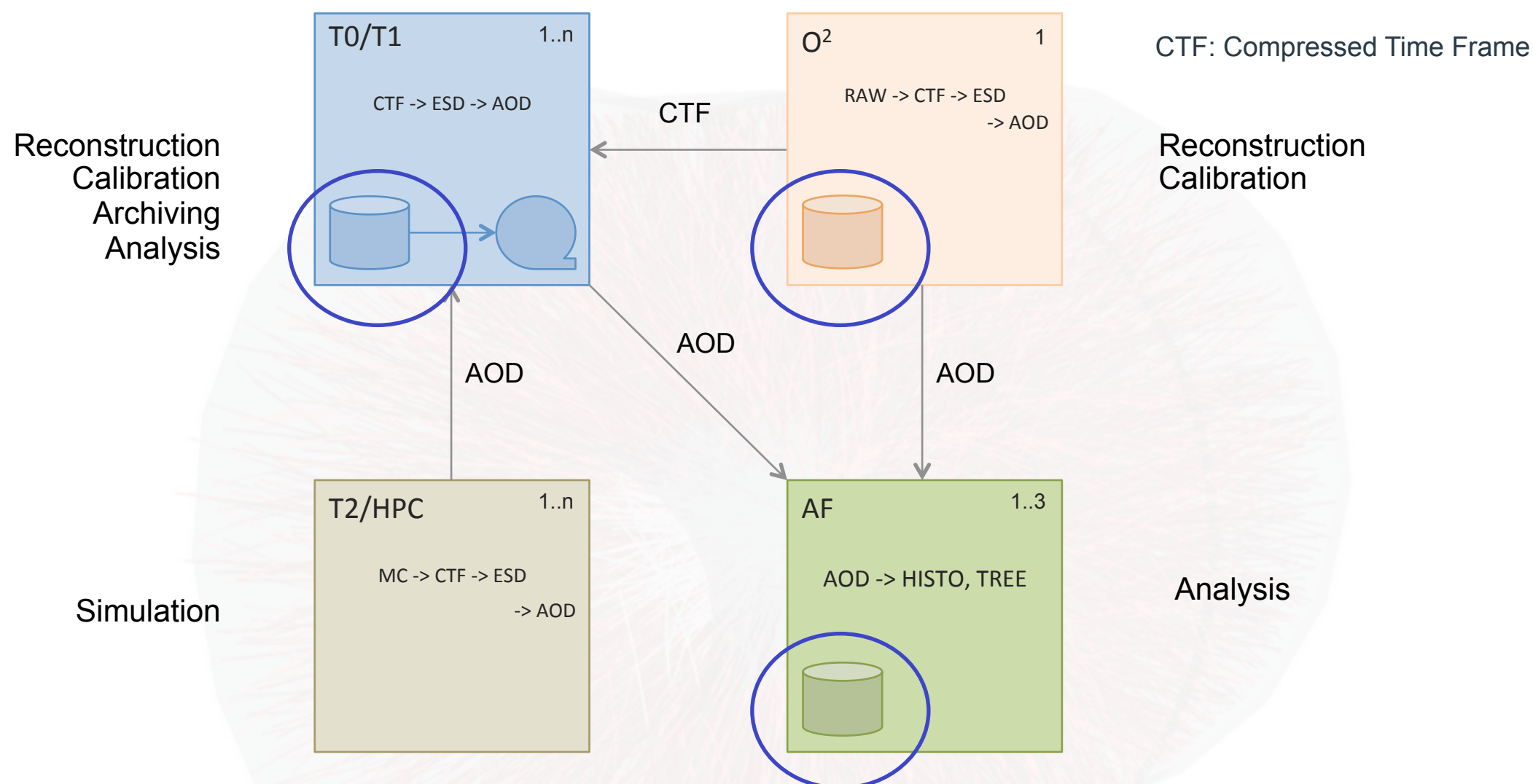
Software architecture

- Simulations for the expected performance of the ALICE upgrade have been carried out within the present framework, AliRoot
- The new O^2 framework will run on both the O^2 facility and on the external sites
- Message-based multi-processing
 - ease of development
 - scalability
 - possibility to extend to different hardware
 - Multi-threading within processes possible
- ALFA (ALICE + FAIR)
 - developed in common by experiments at FAIR and ALICE
 - based on message transport packages
 - data transport
 - Dynamic Deployment System



Pb-Pb @ $\sqrt{s} = 2.76$ ATeV
2011-11-12 06:51:12
File : 2290
Run : 167693
Event : 0x3d94315a

Roles of computing centers in Run 3

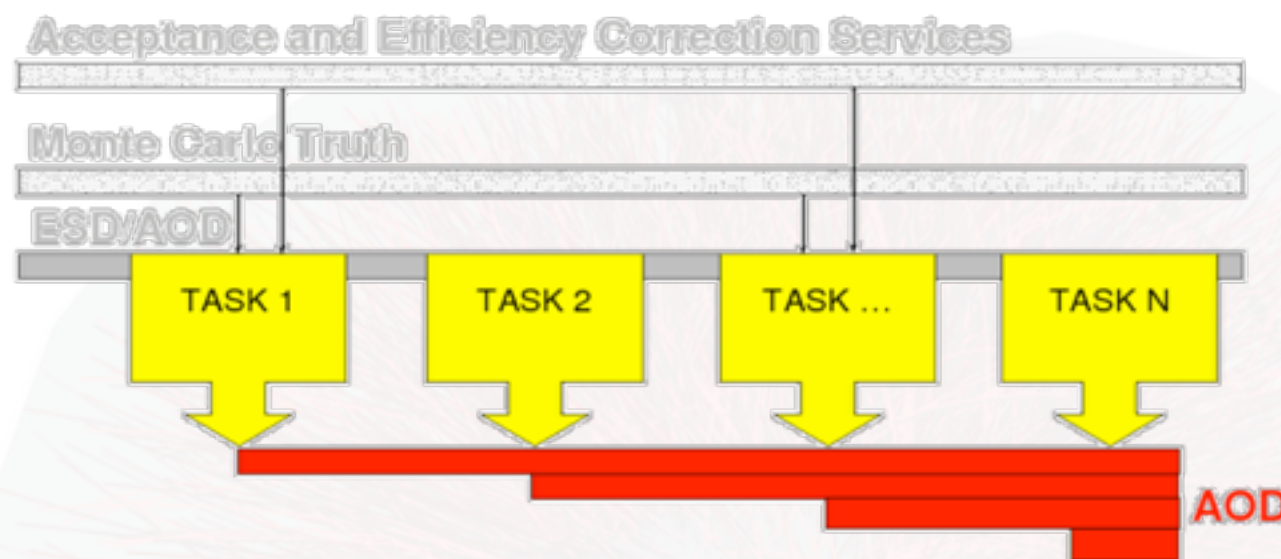


Grid Tiers will be **mostly specialized for given role**

- O² facility (2/3 of reconstruction and calibration), T1s (1/3 of reconstruction and calibration, archiving to tape), T2s (simulation)
- All AODs will be collected on the specialized Analysis Facilities (AF) capable of processing ~5 PB of data within ½ day timescale

The goal is to minimize data movement and optimize processing efficiency

Analysis facilities



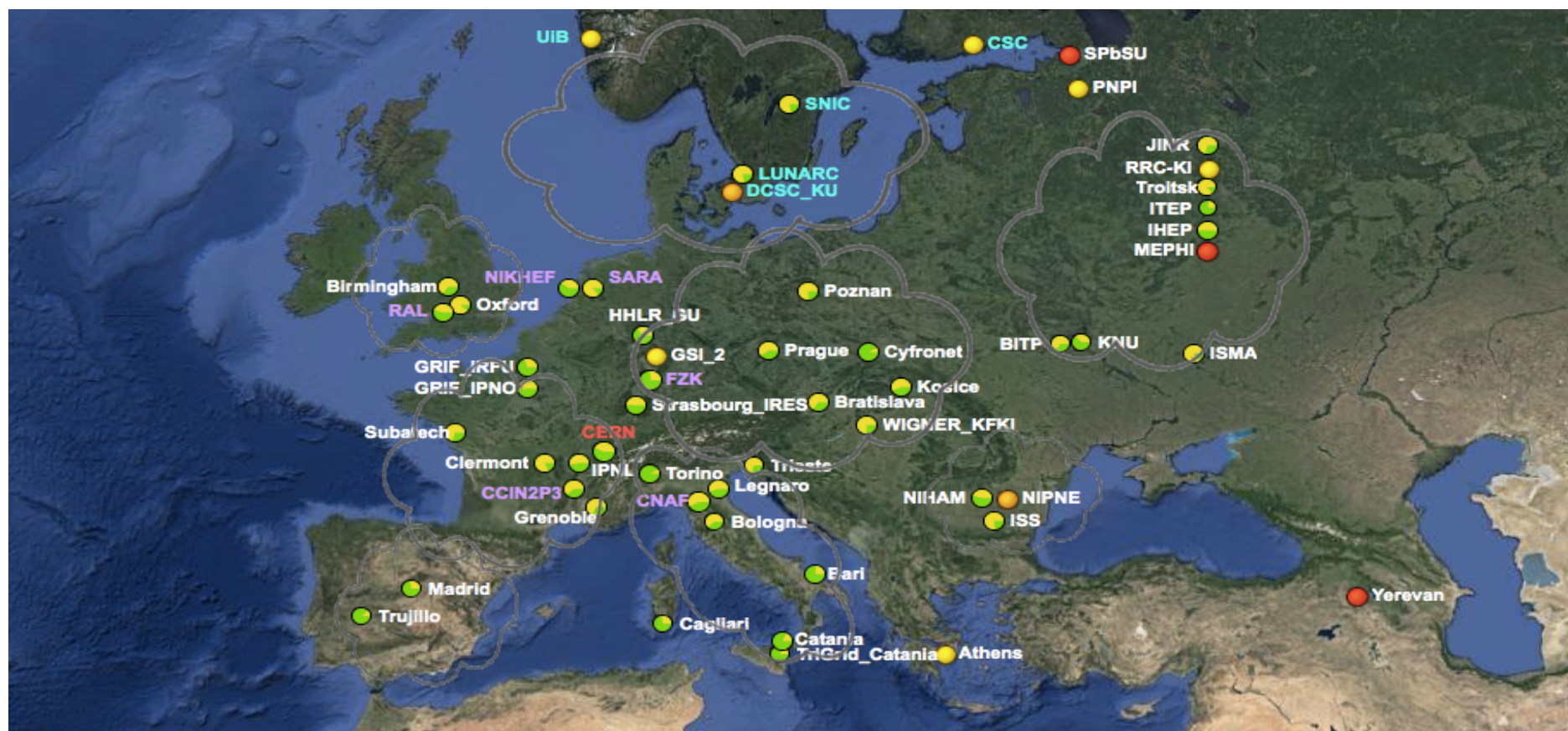
- Analysis Facility
 - Collect AODs on a **few dedicated sites** (AFs) that are capable of locally processing quickly large data volume
 - The AF needs to be able to digest more than 4PB of AODs in a 12 hours period
 - Typically (a fraction of) HPC facility (**20-30'000 cores**) and **5-10 PB** of disk on very performant file system
 - Analysis trains need on average 5 MB/s per job slot to be reasonably efficient.
 - We require the cluster file system able to serve 20,000 job slots at an aggregate throughput of 100 GB/s.

2011-11-12 06:51:12

Run : 167693

Event : 0x3d94315a

Reducing complexity



- Virtually joining together the sites based on proximity (latency) and network capacity into Regional Data Clouds
- Each cloud/region has to provide reliable data management and sufficient processing capability
 - Dealing with handful of clouds/regions instead of the individual sites

No Replication policy

- Due to a substantial increase in data volume, in Run3, there will be only **one instance** of each raw data file (CTF - Compressed Time Frame) stored on disk with a backup on tape
 - In case of data loss, we will restore lost files from the tape
- The size of O² disk buffer should be sufficient to accommodate CTF data from the entire period.
 - As soon as it is available, the CTF data will be archived to the Tier 0 tape buffer or moved to the Tier 1

Pb-Pb @ \sqrt{s} = 2.76 ATeV
2011-11-12 06:51:12
File : 2290
Run : 167693
Event : 0x3d94315a

Deletion policy



- With the exception of raw data (CTF) and derived analysis data (AOD), all the other intermediate data created at various processing stages are transient (removed after a given processing step) or temporary (with limited lifetime)
 - CTF and AODs are archived to tape
- Given the limited size of the disk buffers in O² and Tier 1s, all CTF data collected in the previous year, will have to be removed before new data taking period starts.
- All data not finally processed during this period will remain parked on tapes until the next opportunity for re-processing arises: LS3



Needed resources

- 2019-2020: essentially replacements of obsolete HW with a possible marginal increase of the resources
- Estimates for computing needs based on:
 - no replication + deletion policies
 - an online compression factor ~ 16
 - **$\sim 20\%$ yearly growth during Run 3**
 - **x 2** of the resources at the end of Run 3 w.r.t. Run 2
 - Pessimistic estimates based on an online compression factor of 12%
 - **$\sim 27\%$ yearly growth during Run 3**
 - **x 2.5** of the resources at the end of Run 3 w.r.t. Run 2
- **Caveat:** AF are out of these evaluations:
 - $2 \div 3$ centers with:
 - target size: overall 20000-30000 cores and 5-10 PB of disk storage
 - progressive deployment - no need full size AF at Run 3 start
 - impact: $10 \div 15\%$ of the total WLCG resources (T0+T1+T2)
 - AF to be provided by FA as in-kind contribution to the experiment

Current activities in Italy for the upgrade

- INFN is not an official member of the O² project
- However detector groups have to contribute to the core offline software related to their detectors
- Currently our community is active in the ITS-Upgrade software development
- **Vertexing** and standalone **track reconstruction** for the ITS are currently coded in Italy
- Also contributing to the porting of the code to the new O² infrastructure

Pb-Pb @ \sqrt{s} = 2.76 ATeV

2011-11-12 06:51:12

File : 2290

Run : 167693

Event : 0x3d94315a

ITS standalone tracking

- TPC tracks can be prolonged inwards to the Inner Tracking System.
- However the ITS can be used as a **standalone detector** and tracks found in the ITS can be prolonged to the TPC.
- Since in Run 3 we will need to calibrate the TPC online, the ITS track seeds will be needed for TPC calibration
- An ITS Tracker based on a **Cellular Automaton** has been coded and tested on CPU, within the present ALICE offline framework, AliRoot.
- The next step is to port this code to the new O² framework and to a heterogeneous **CPU-GPU computing environment**
- The goal is to have a demonstrator of TPC+ITS tracking in 2016.

Pb-Pb @ \sqrt{s} = 2.76 ATeV

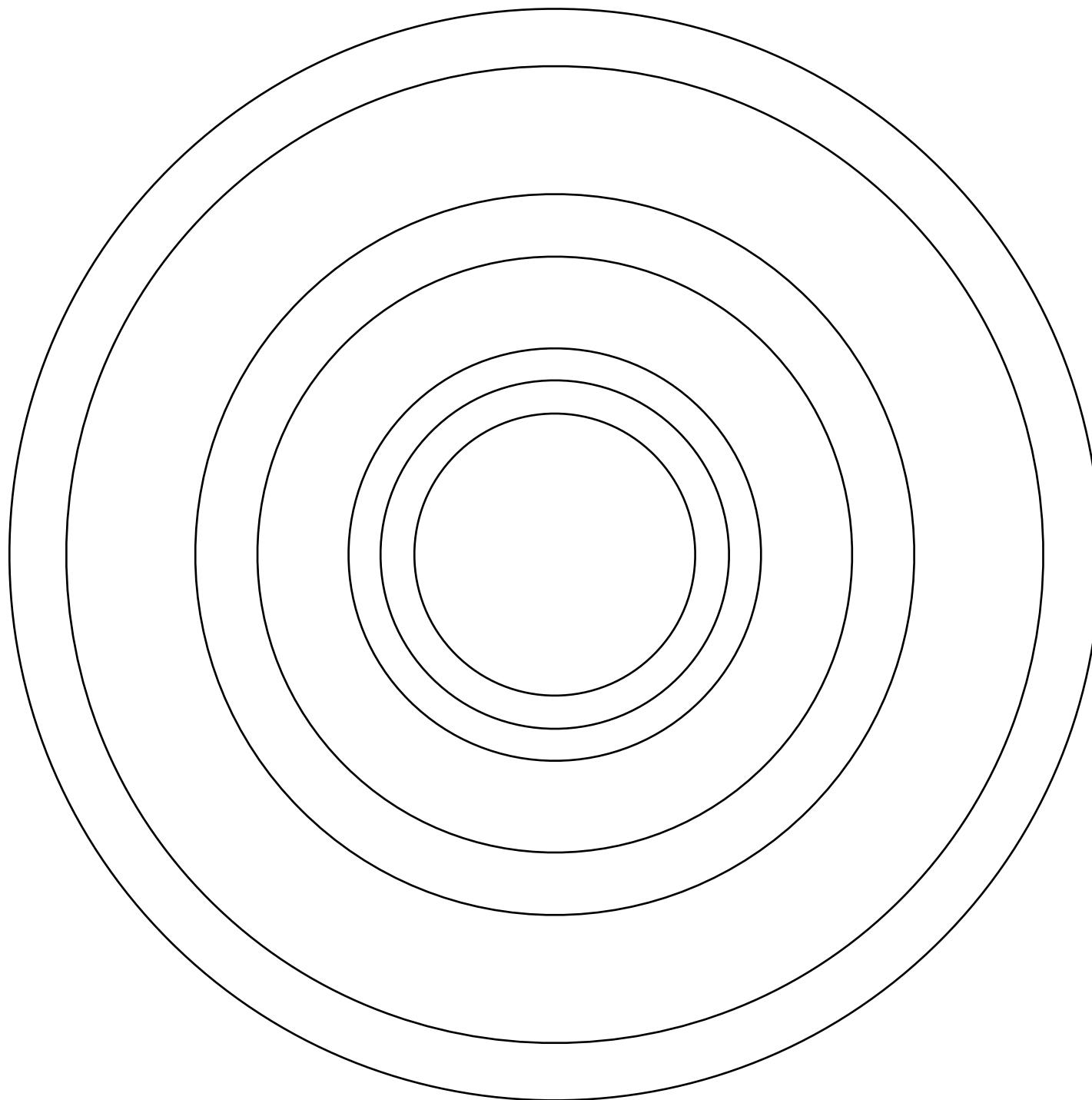
2011-11-12 06:51:12

File : 2290

Run : 167693

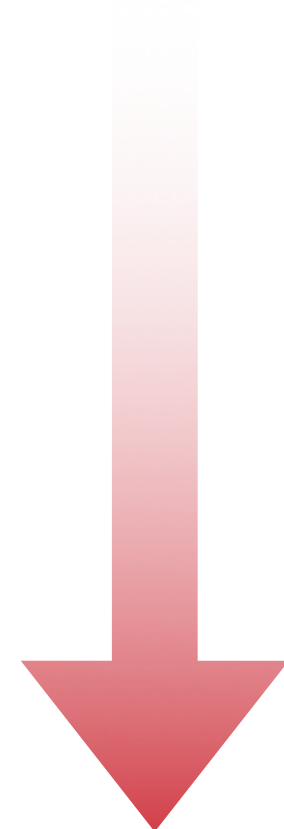
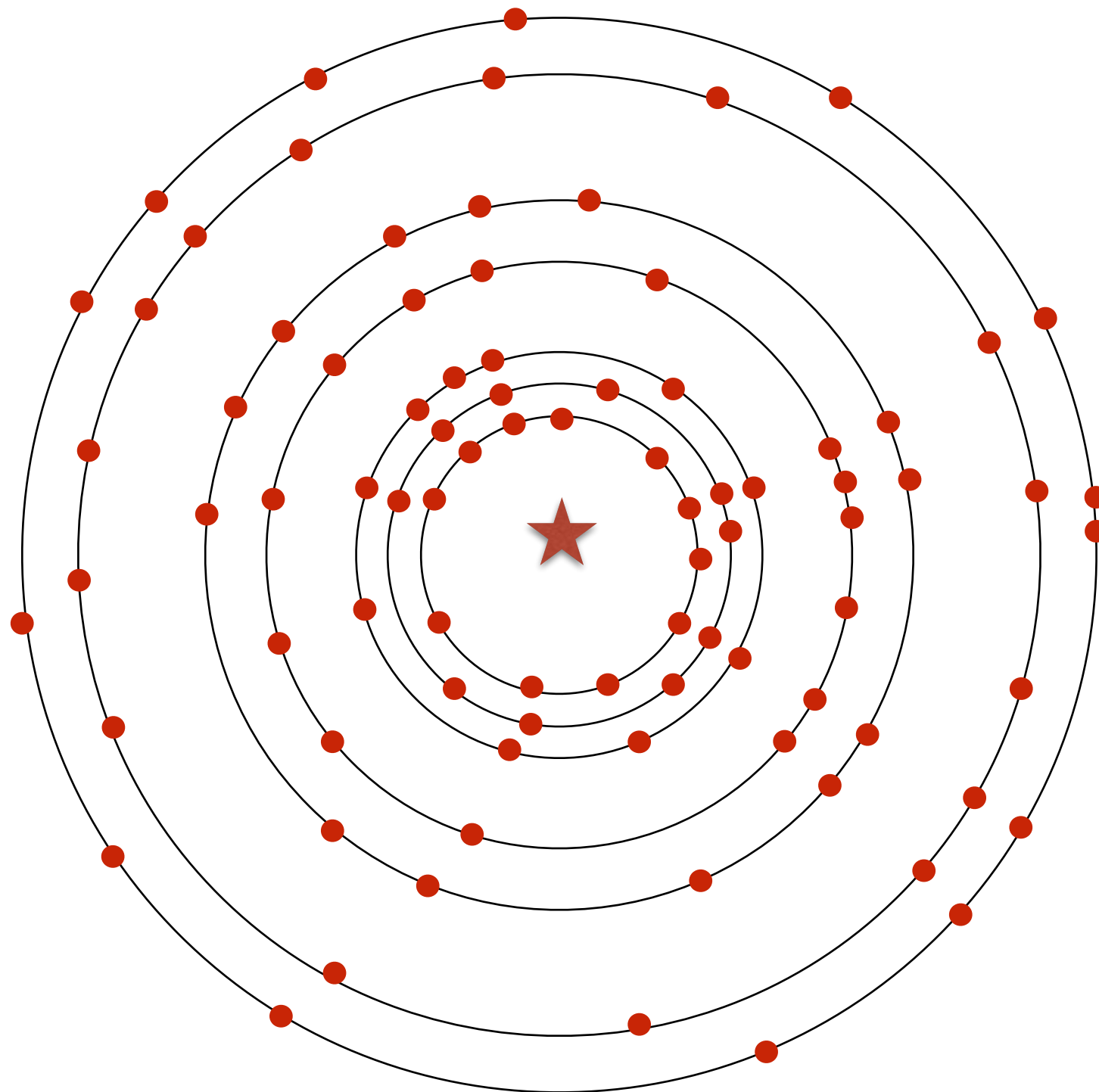
Event : 0x3d94315a

Tracking with ITS Upgrade

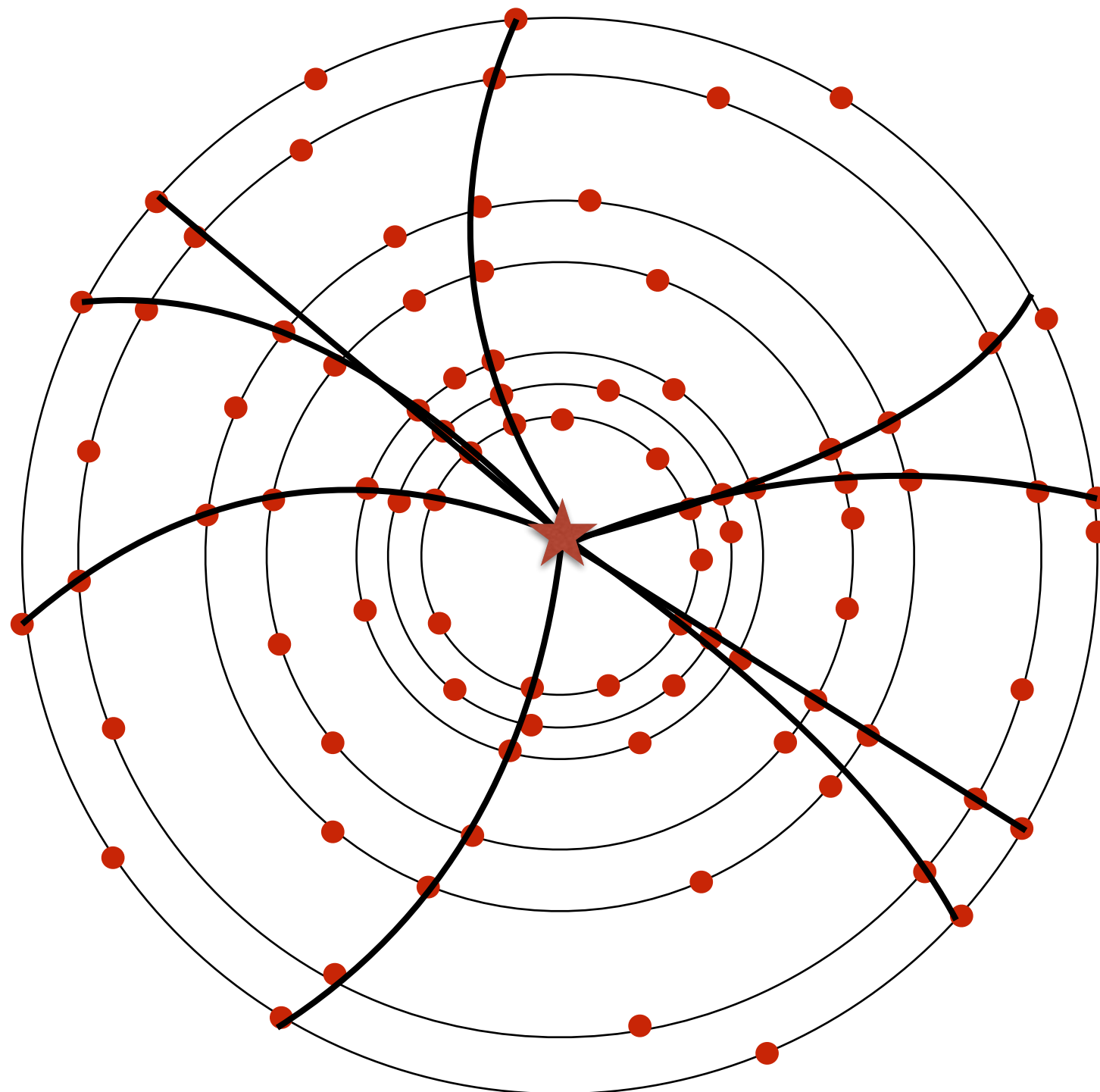


DISCLAIMER:
For easiness in explanations,
in the following I will show only
the transverse section of ITSU.

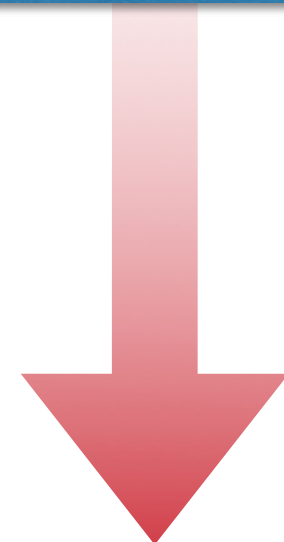
Tracking with ITS Upgrade



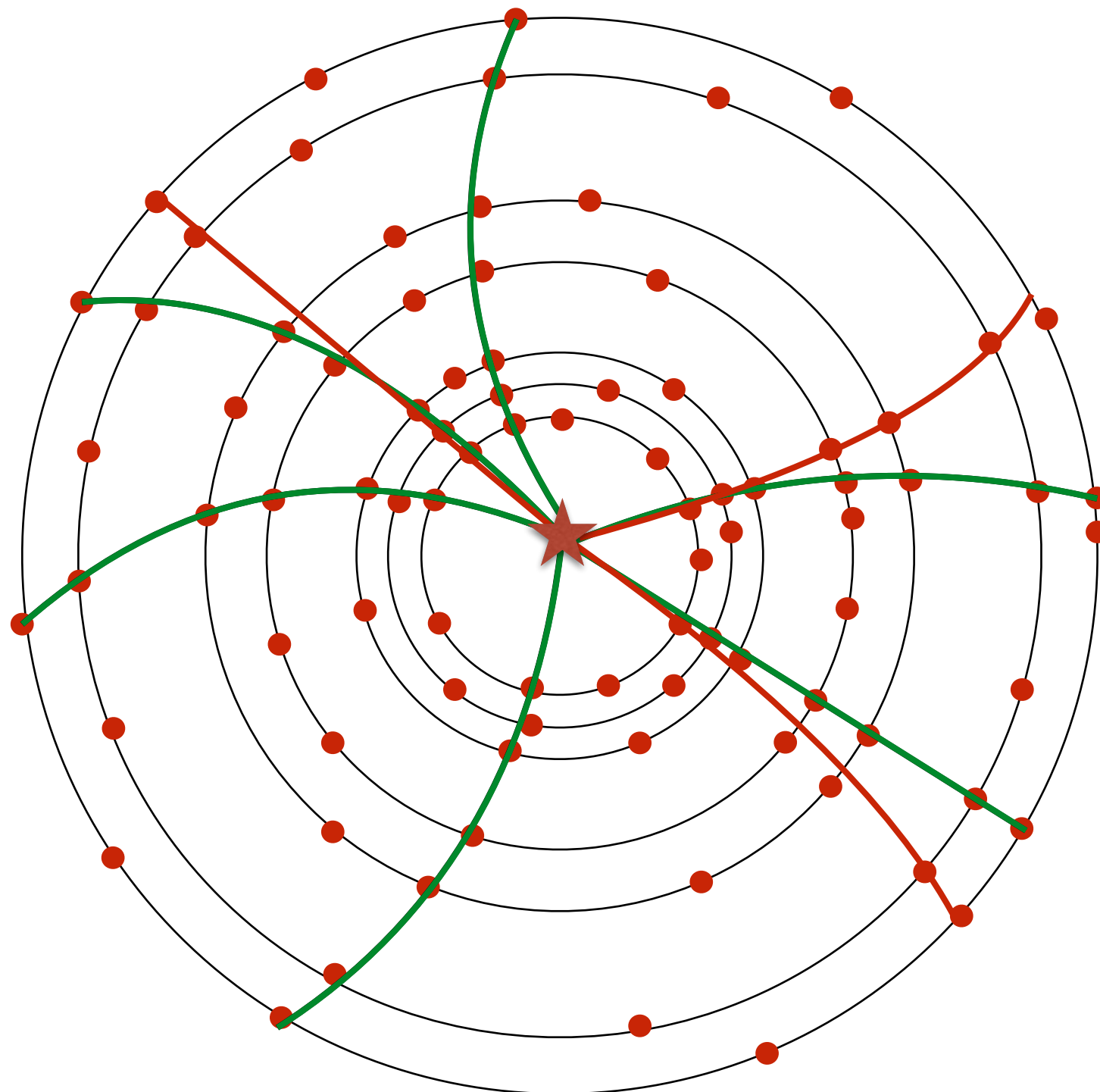
Tracking with ITS Upgrade



Using a pattern recognition method, find track candidates

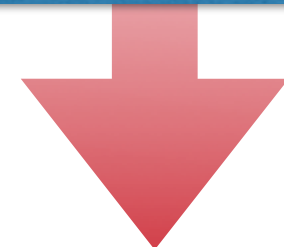


Tracking with ITS Upgrade

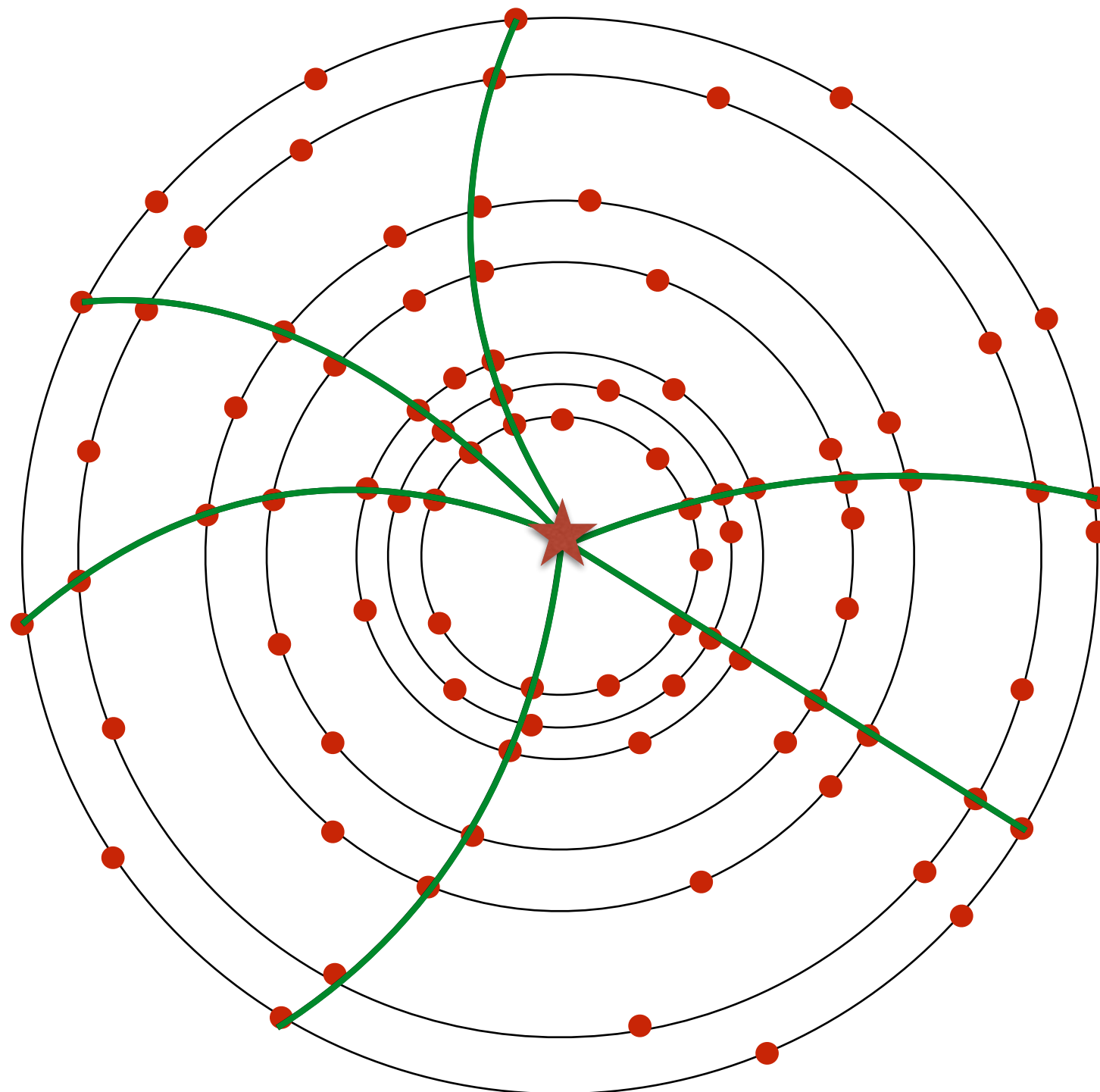


Using a pattern recognition method, find track candidates

Fitting of the candidates using Kalman Filter in three passes (inward, outward, inward)



Tracking with ITS Upgrade



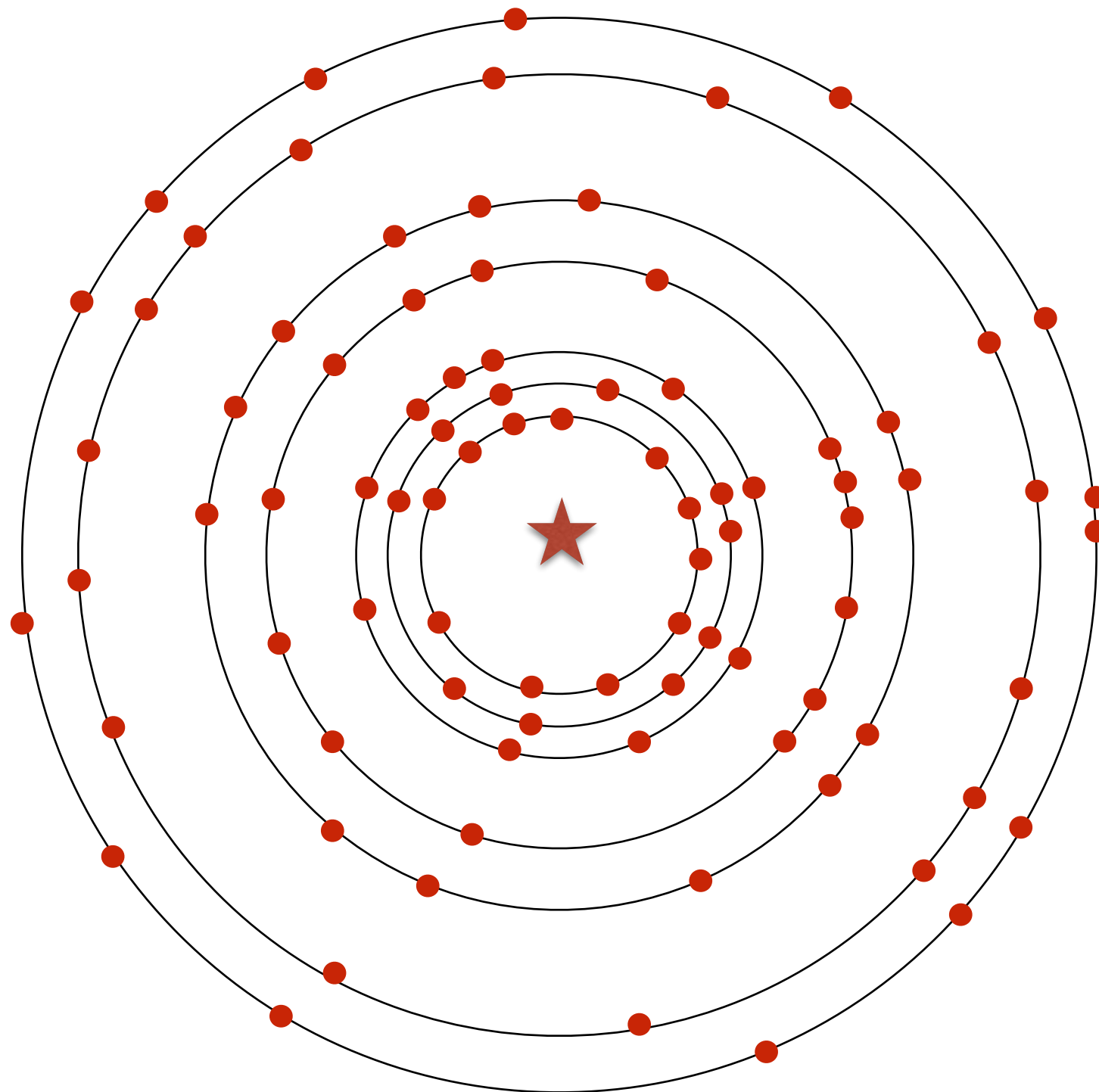
Using a pattern recognition method, find track candidates

Fitting of the candidates using Kalman Filter in three passes (inward, outward, inward)

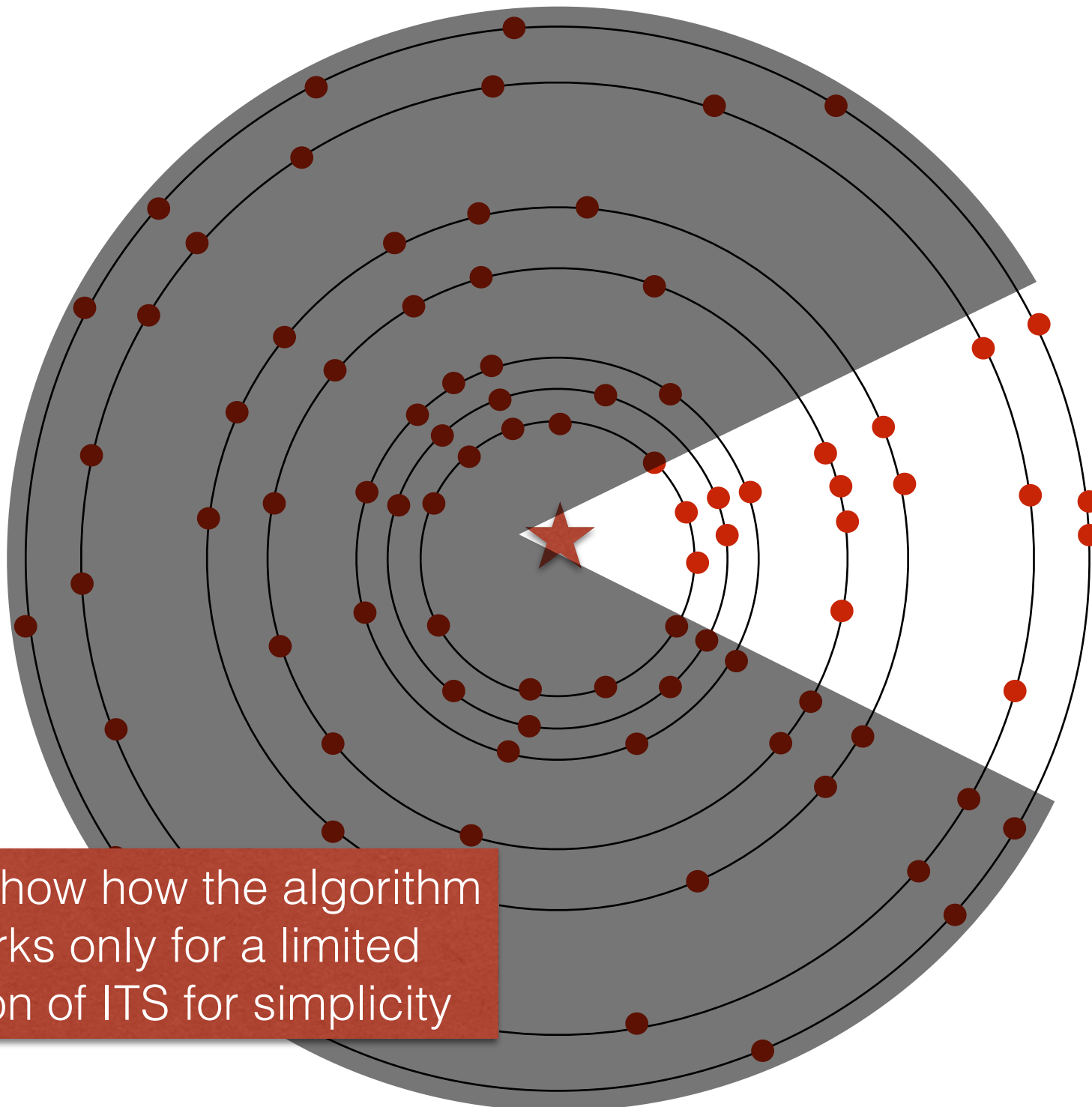
Candidates with the best χ^2 values are stored as reconstructed tracks

Currently two different approaches to the pattern recognition step are implemented for ITS Upgrade

The Cellular Automaton approach



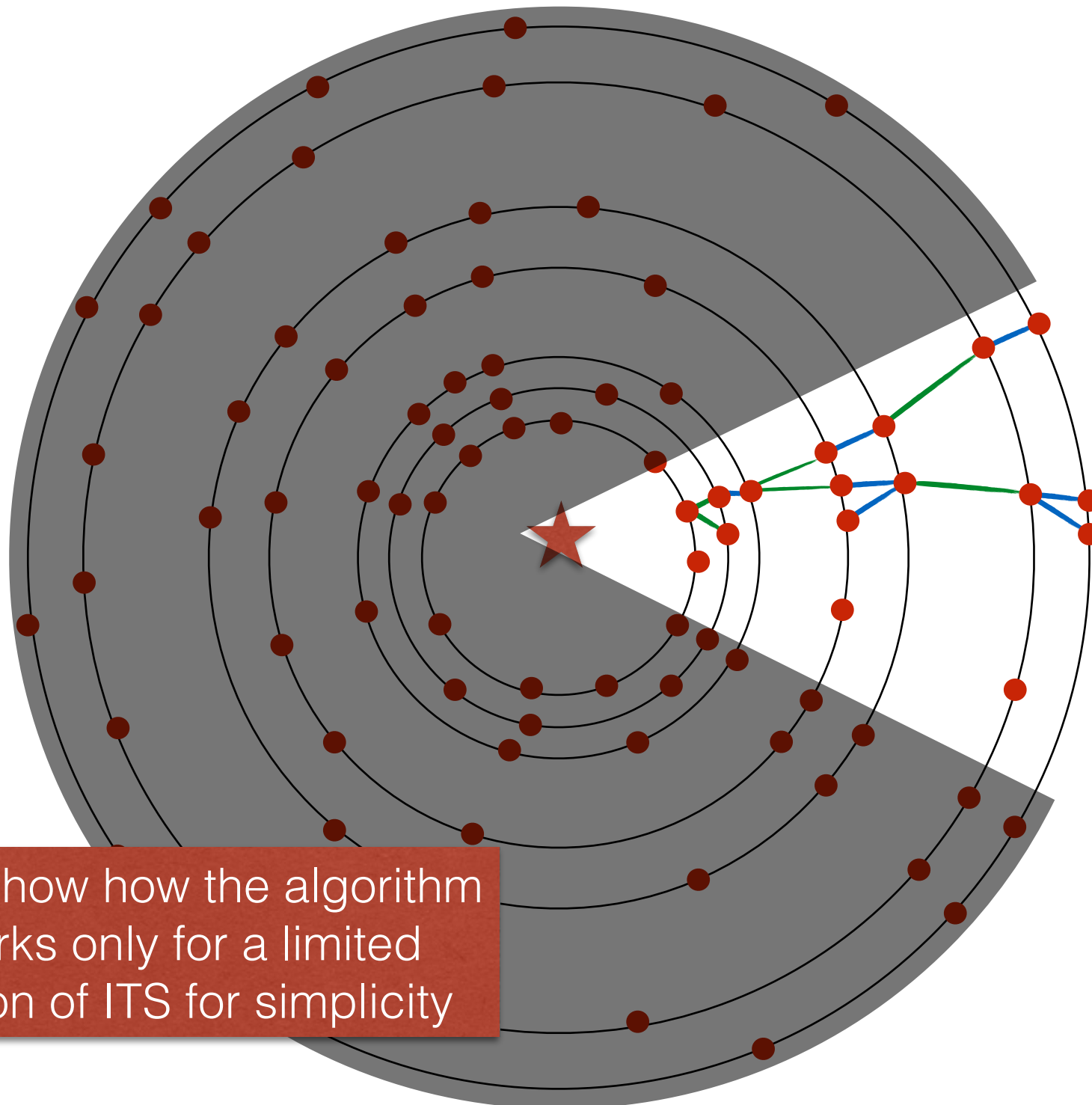
The Cellular Automaton approach



I will show how the algorithm works only for a limited region of ITS for simplicity

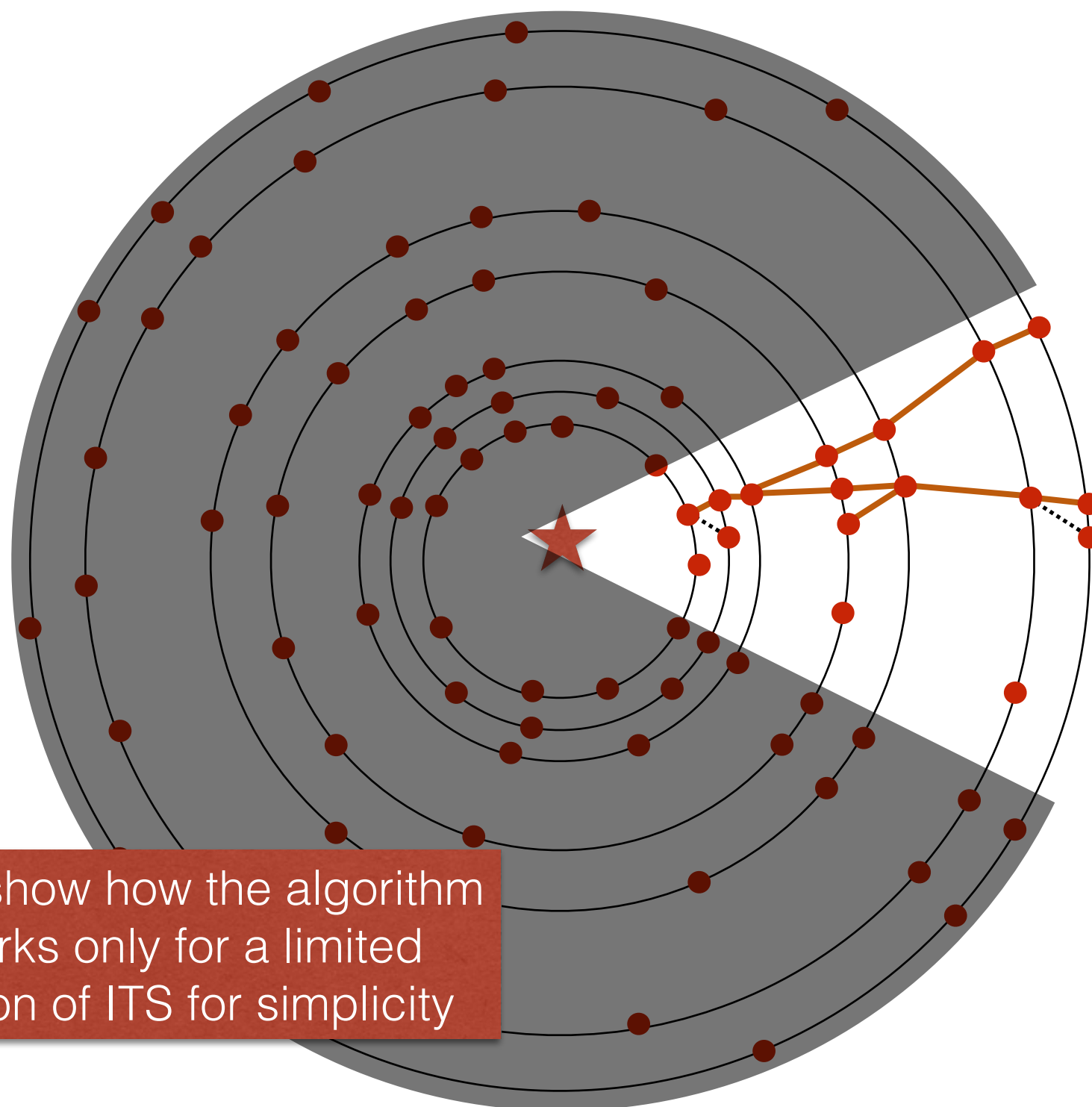
The Cellular Automaton approach

For each cluster on each layer a 2D window is opened. Then the clusters are joined with those on the next layer within the window



I will show how the algorithm works only for a limited region of ITS for simplicity

The Cellular Automaton approach

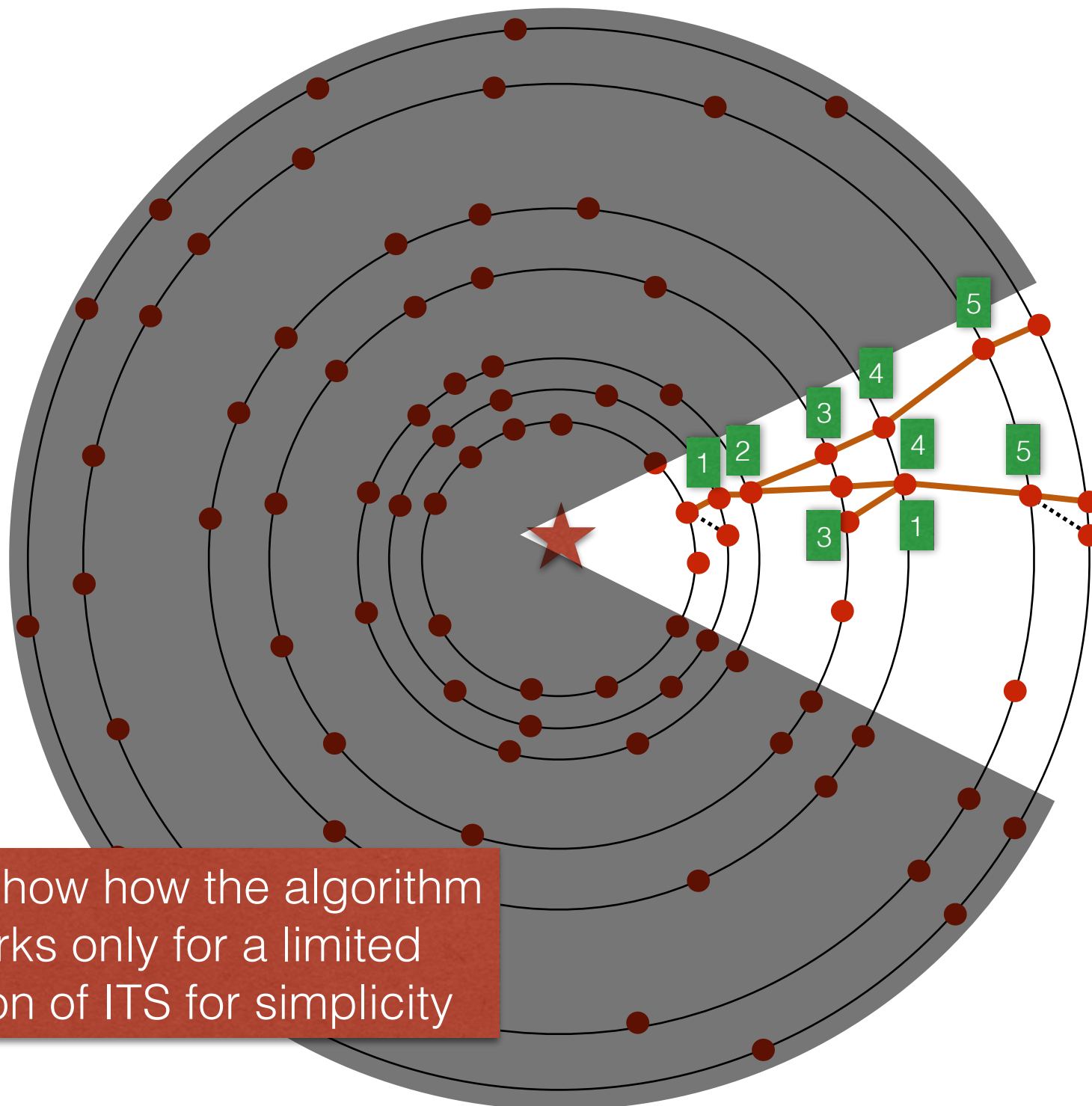


For each cluster on each layer a 2D window is opened. Then the clusters are joined with those on the next layer within the window

Subsequent doublets are combined in cells (3 points seed) and track params are computed

I will show how the algorithm works only for a limited region of ITS for simplicity

The Cellular Automaton approach



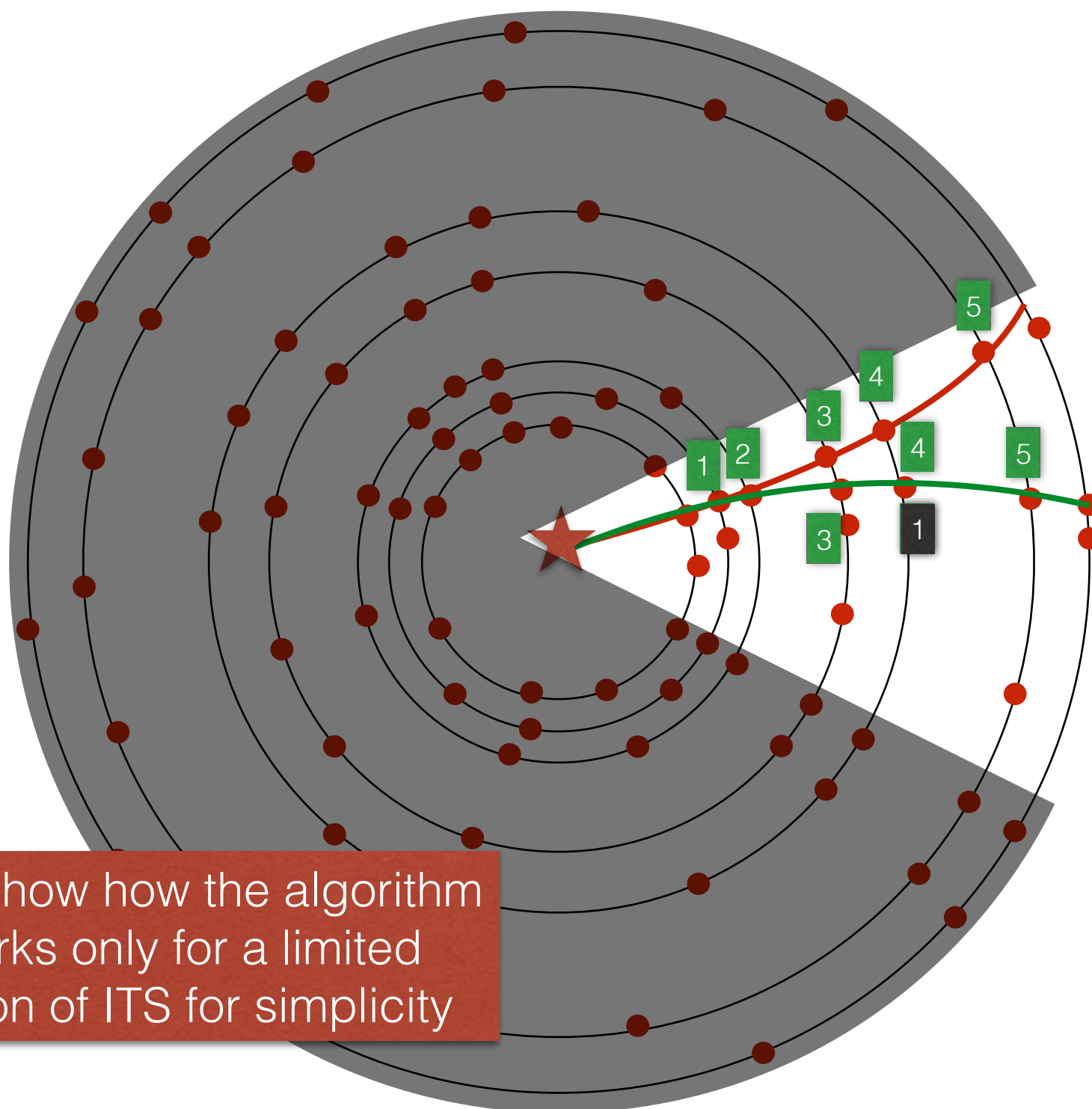
For each cluster on each layer a 2D window is opened. Then the clusters are joined with those on the next layer within the window

Subsequent doublets are combined in cells (3 points seed) and track params are computed

Each cell has an index representing the number of connected inner cells + 1

I will show how the algorithm works only for a limited region of ITS for simplicity

The Cellular Automaton approach



I will show how the algorithm works only for a limited region of ITS for simplicity

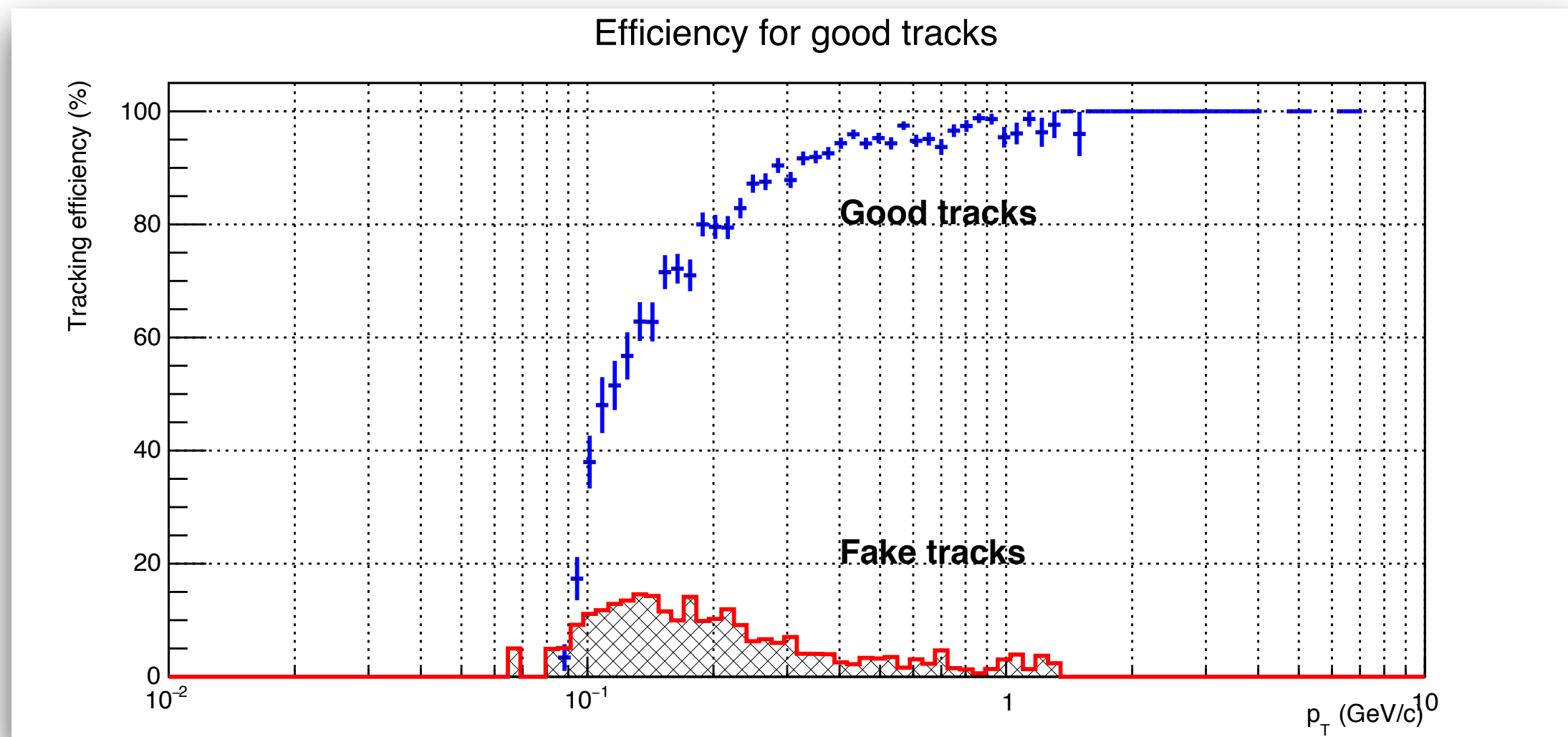
For each cluster on each layer a 2D window is opened. Then the clusters are joined with those on the next layer within the window

Subsequent doublets are combined in cells (3 points seed) and track params are computed

Each cell has an index representing the number of connected inner cells + 1

Longest, continuous sequences of indices represent candidates

Tracking performance Pb-Pb



	$p_T > 0.6$ GeV/c	$p_T > 2$ GeV/c	Full
Central Pb-Pb	~0.3 s	~0.2 s	~0.7 s
Central Pb-Pb with noise	~1.3 s	~1.0 s	~5 s
p-p with noise	~0.6 s	~0.7 s	~2.3 s

Conclusions

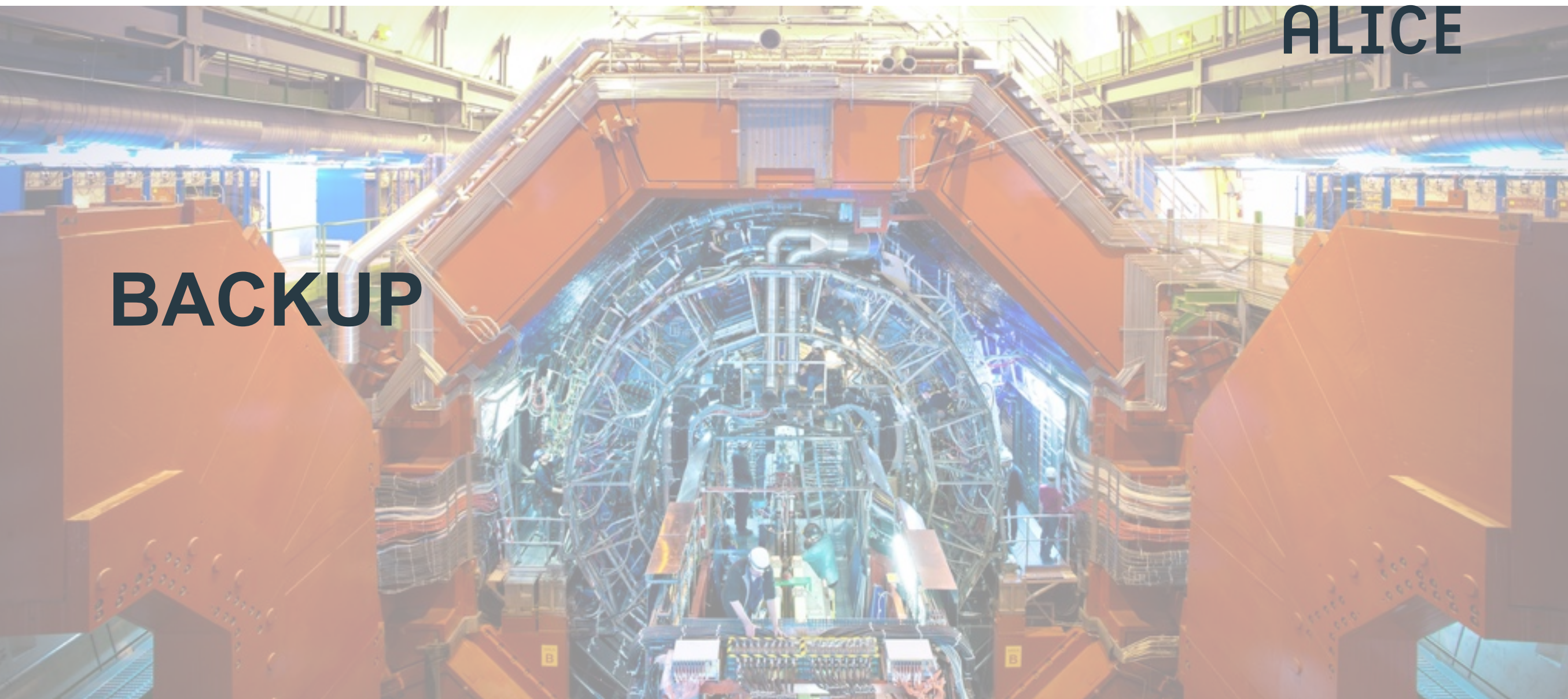
- Computing in ALICE is still largely based on the intrinsic parallelism offered by the independence of different events.
- The necessity of performing a synchronous data reconstruction online, during data taking triggered a computing approach that exploits the potentials of a heterogeneous architecture.
- This solution is very successful in the present High Level Trigger system.
- In the next LHC Run 3 this approach will be extended to the whole ALICE data reconstruction.
- Run 3 will be possible with an expected max resource growth of a factor ~ 2.5 w.r.t. Run2
 - Dedicated analysis facilities are expected to be provided in addition to this figure

Pb-Pb @ $\sqrt{s} = 2.76$ ATeV
2011-11-12 06:51:12
Fill : 2290
Run : 167693
Event : 0x3d94315a



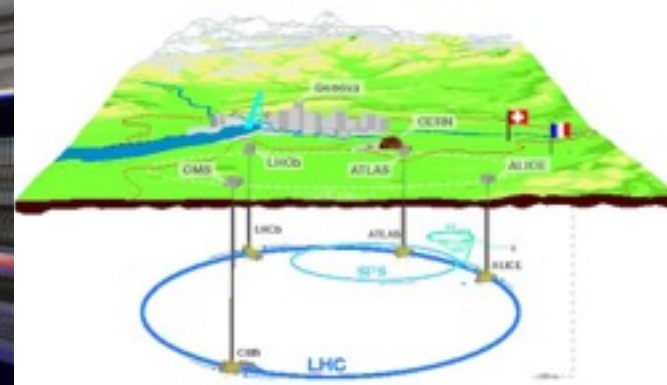
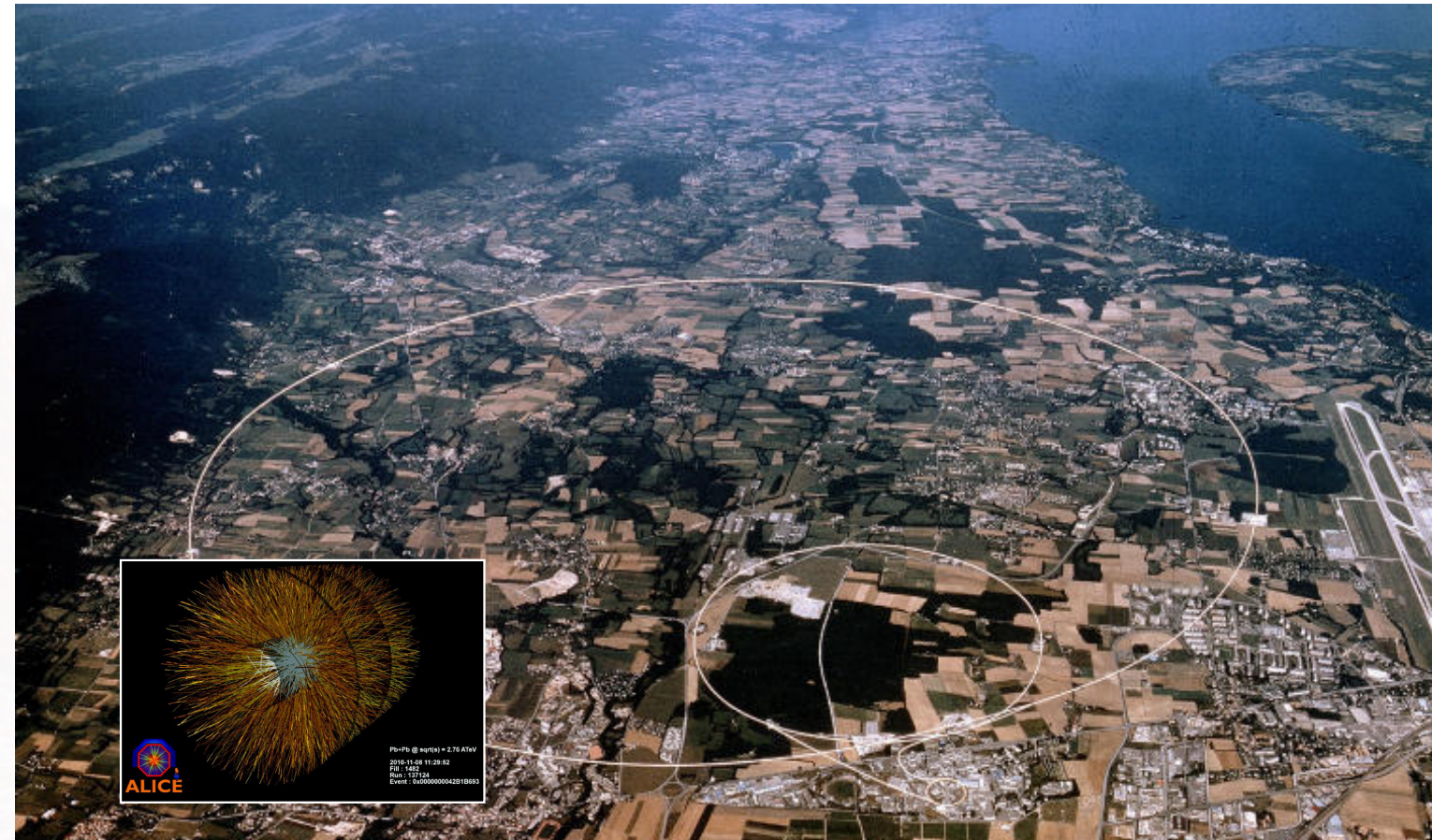
ALICE

BACKUP



LHC

- The Large Hadron Collider (LHC) is the largest and most powerful proton and ion collider in the world.
- The present centre-of-mass energy is:
 - » 13 TeV for pp collisions
 - » 5.02 TeV per nucleon pair for Pb-Pb collisions
- 4 major experiments: ALICE, ATLAS, CMS and LHCb
- ALICE (A Large Ion Collider Experiment) is designed primarily to study nucleus-nucleus collisions.



Run : 167693
Event : 0x3d94315a

Offline computing in ALICE

- When a Pb-Pb (pp or p-Pb) collision meets a predefined trigger condition, the corresponding data are collected: these data form an **event**.
- Ideally, an event corresponds to a single Pb-Pb (pp...) collision, even though **pile-up of several collisions may occur**, depending on the interaction rate.
- The data collected in a continuous data taking period (~ few hours) within one LHC fill, constitute a **run**.
- The events in a run must be processed independently; they share only the same data taking conditions.
- Different runs are **processed independently** on different computing cores (if no hyper-threading):
 - » at a given moment a computing farm with N cores processes N events in parallel;
 - » **embarrassingly parallel computing**.
- Presently, no further parallelization within a single event is attempted.



ALICE

Pb-Pb @ $\sqrt{s} = 2.76$ ATeV

2011-11-12 06:51:12

Fill : 2290

Run : 167693

Event : 0x3d94315a

Offline computing in ALICE

- When a Pb-Pb (pp or p-Pb) collision meets a predefined trigger condition, the corresponding data are collected: these data form an **event**.
- Ideally, an event corresponds to a single Pb-Pb (pp...) collision, even though **pile-up of several collisions may occur**, depending on the interaction rate.
- The data collected in a continuous data taking period (~ few hours) within one LHC fill, constitute a **run**.
- The events in a run must be processed independently; they share only the same data taking conditions.
- Different runs are **processed independently** on different computing cores (if no hyper-threading):
 - » at a given moment a computing farm with N cores processes N events in parallel;
 - » **embarrassingly parallel computing**.
- Presently, no further parallelization within a single event is attempted.

RUN N

RUN N+1

RUN N+2

Pb-Pb @ $\sqrt{s} = 2.76$ ATeV

2011-11-12 06:51:12

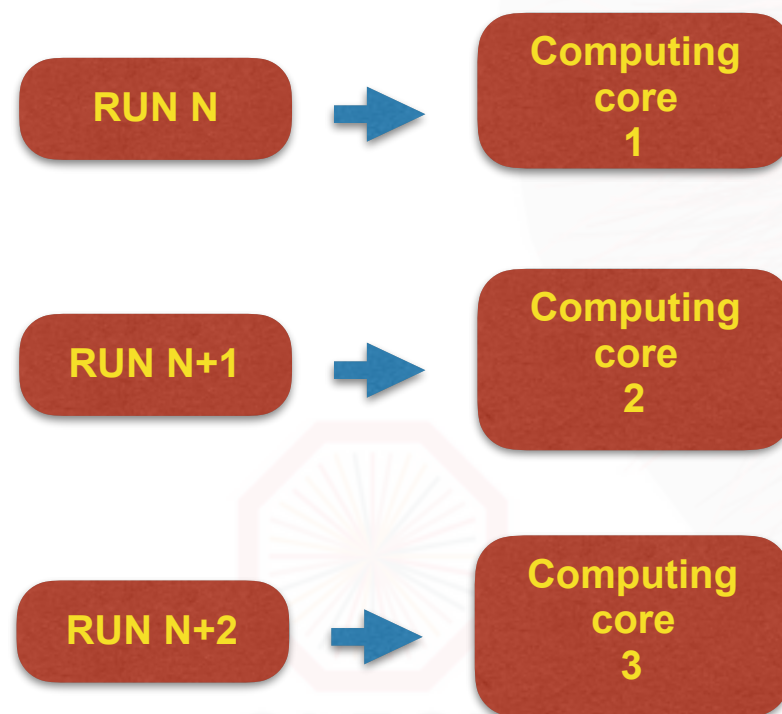
Fill : 2290

Run : 167693

Event : 0x3d94315a

Offline computing in ALICE

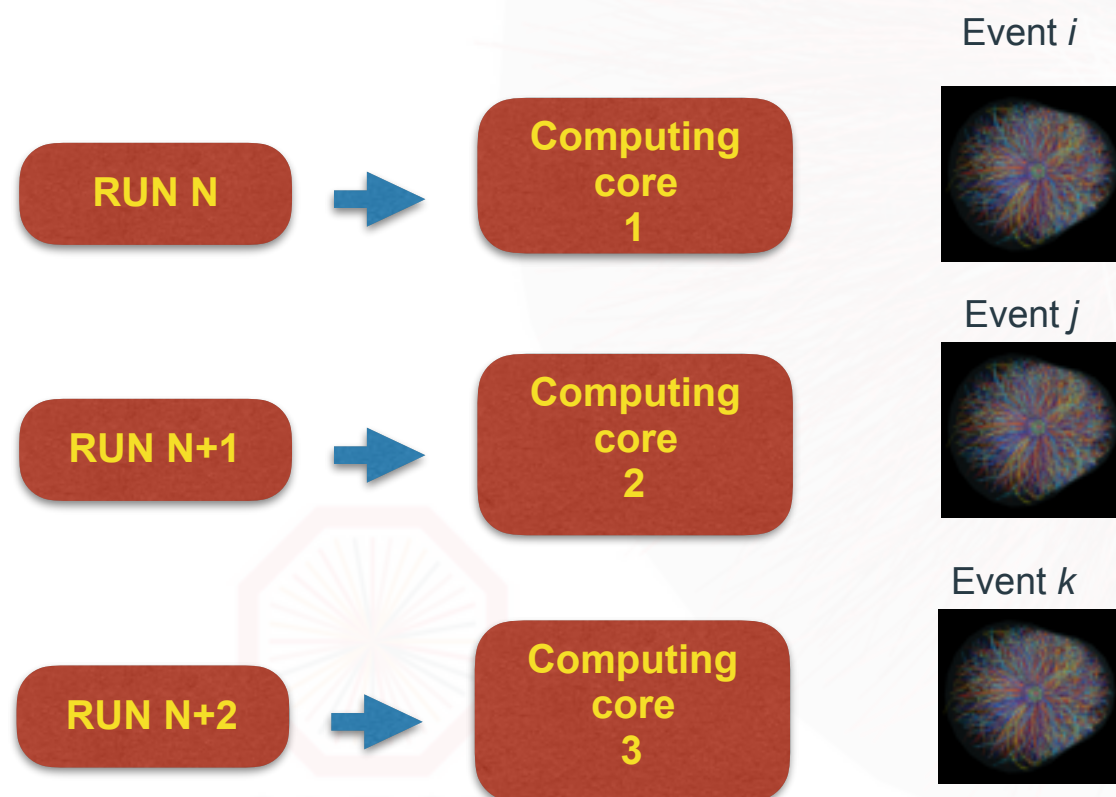
- When a Pb-Pb (pp or p-Pb) collision meets a predefined trigger condition, the corresponding data are collected: these data form an **event**.
- Ideally, an event corresponds to a single Pb-Pb (pp...) collision, even though **pile-up of several collisions may occur**, depending on the interaction rate.
- The data collected in a continuous data taking period (~ few hours) within one LHC fill, constitute a **run**.
- The events in a run must be processed independently; they share only the same data taking conditions.
- Different runs are **processed independently** on different computing cores (if no hyper-threading):
 - » at a given moment a computing farm with N cores processes N events in parallel;
 - » **embarrassingly parallel computing**.
- Presently, no further parallelization within a single event is attempted.



Pb-Pb @ $\sqrt{s} = 2.76$ ATeV
2011-11-12 06:51:12
Fill : 2290
Run : 167693
Event : 0x3d94315a

Offline computing in ALICE

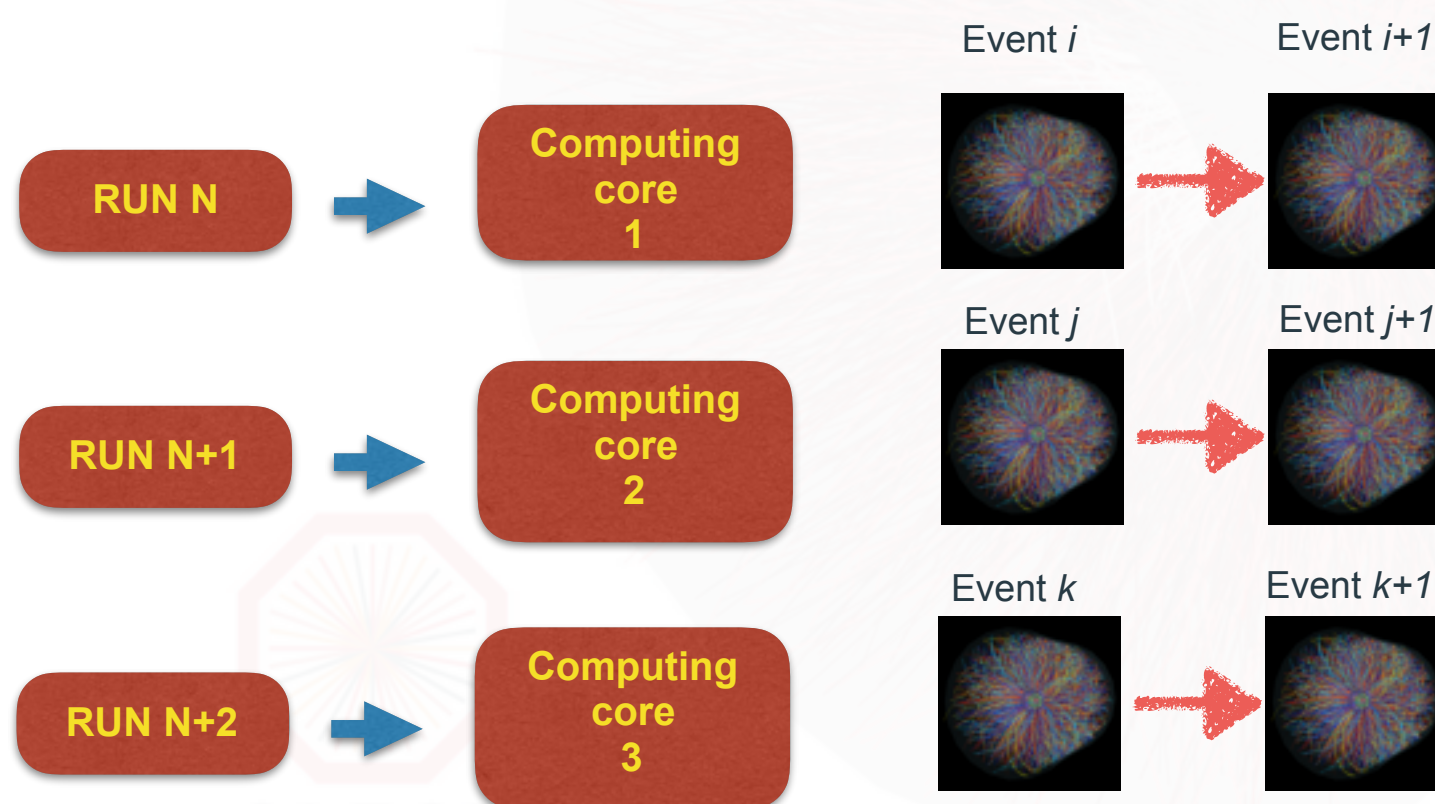
- When a Pb-Pb (pp or p-Pb) collision meets a predefined trigger condition, the corresponding data are collected: these data form an **event**.
- Ideally, an event corresponds to a single Pb-Pb (pp...) collision, even though **pile-up of several collisions may occur**, depending on the interaction rate.
- The data collected in a continuous data taking period (~ few hours) within one LHC fill, constitute a **run**.
- The events in a run must be processed independently; they share only the same data taking conditions.
- Different runs are **processed independently** on different computing cores (if no hyper-threading):
 - » at a given moment a computing farm with N cores processes N events in parallel;
 - » **embarrassingly parallel computing**.
- Presently, no further parallelization within a single event is attempted.



Pb-Pb @ $\sqrt{s} = 2.76$ ATeV
 2011-11-12 06:51:12
 Fill : 2290
 Run : 167693
 Event : 0x3d94315a

Offline computing in ALICE

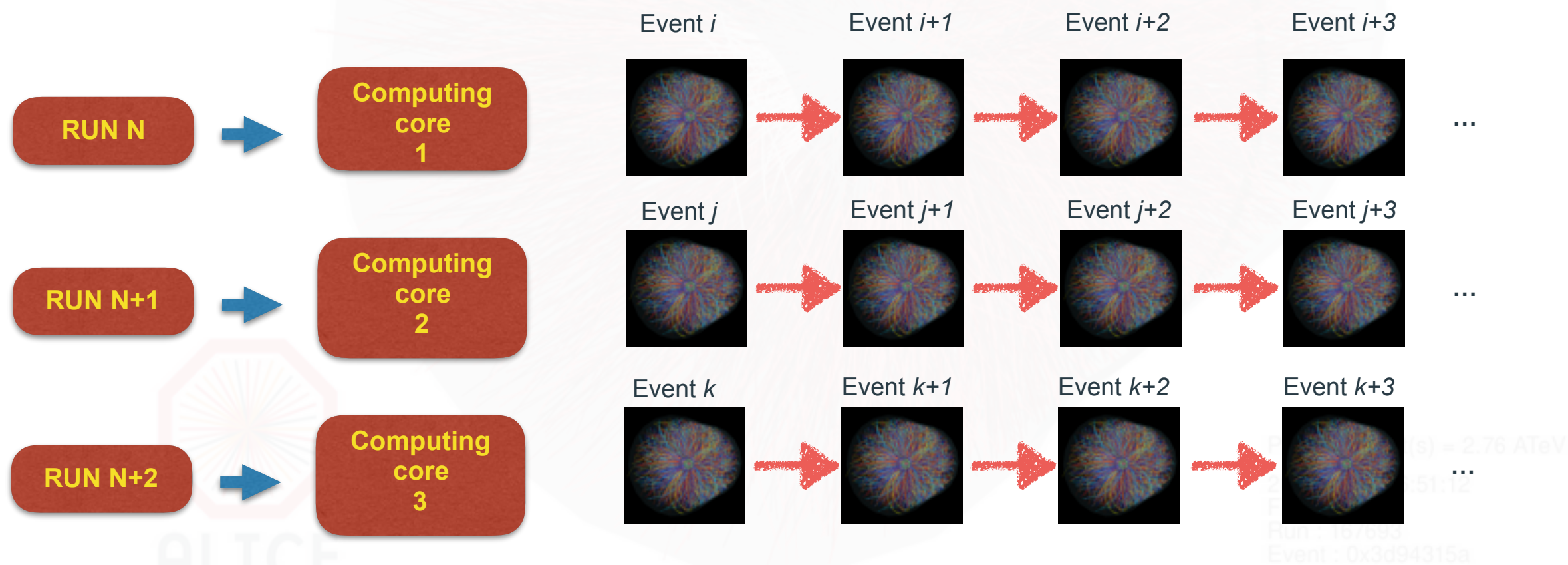
- When a Pb-Pb (pp or p-Pb) collision meets a predefined trigger condition, the corresponding data are collected: these data form an **event**.
- Ideally, an event corresponds to a single Pb-Pb (pp...) collision, even though **pile-up of several collisions may occur**, depending on the interaction rate.
- The data collected in a continuous data taking period (~ few hours) within one LHC fill, constitute a **run**.
- The events in a run must be processed independently; they share only the same data taking conditions.
- Different runs are **processed independently** on different computing cores (if no hyper-threading):
 - » at a given moment a computing farm with N cores processes N events in parallel;
 - » **embarrassingly parallel computing**.
- Presently, no further parallelization within a single event is attempted.



Pb-Pb @ $\sqrt{s} = 2.76$ ATeV
2011-11-12 06:51:12
Fill : 2290
Run : 167693
Event : 0x3d94315a

Offline computing in ALICE

- When a Pb-Pb (pp or p-Pb) collision meets a predefined trigger condition, the corresponding data are collected: these data form an **event**.
- Ideally, an event corresponds to a single Pb-Pb (pp...) collision, even though **pile-up of several collisions may occur**, depending on the interaction rate.
- The data collected in a continuous data taking period (~ few hours) within one LHC fill, constitute a **run**.
- The events in a run must be processed independently; they share only the same data taking conditions.
- Different runs are **processed independently** on different computing cores (if no hyper-threading):
 - » at a given moment a computing farm with N cores processes N events in parallel;
 - » **embarrassingly parallel computing**.
- Presently, no further parallelization within a single event is attempted.

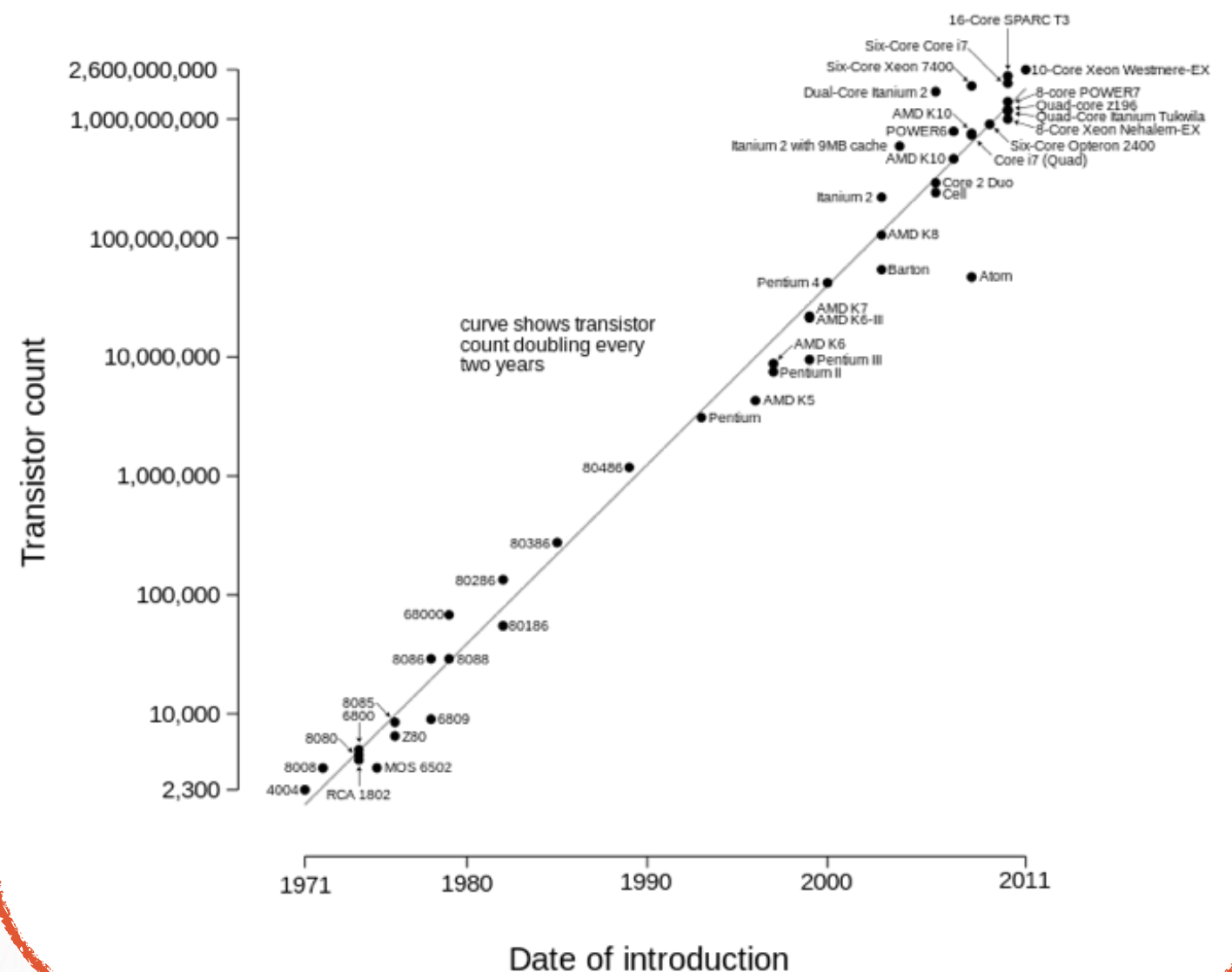


Moore's law

- According to Moore's law the number of transistors per chip doubles in 24 months
- Apart the fact that we are reaching a saturation due to physical limits, this growth does not necessarily imply an increase of performance for our software

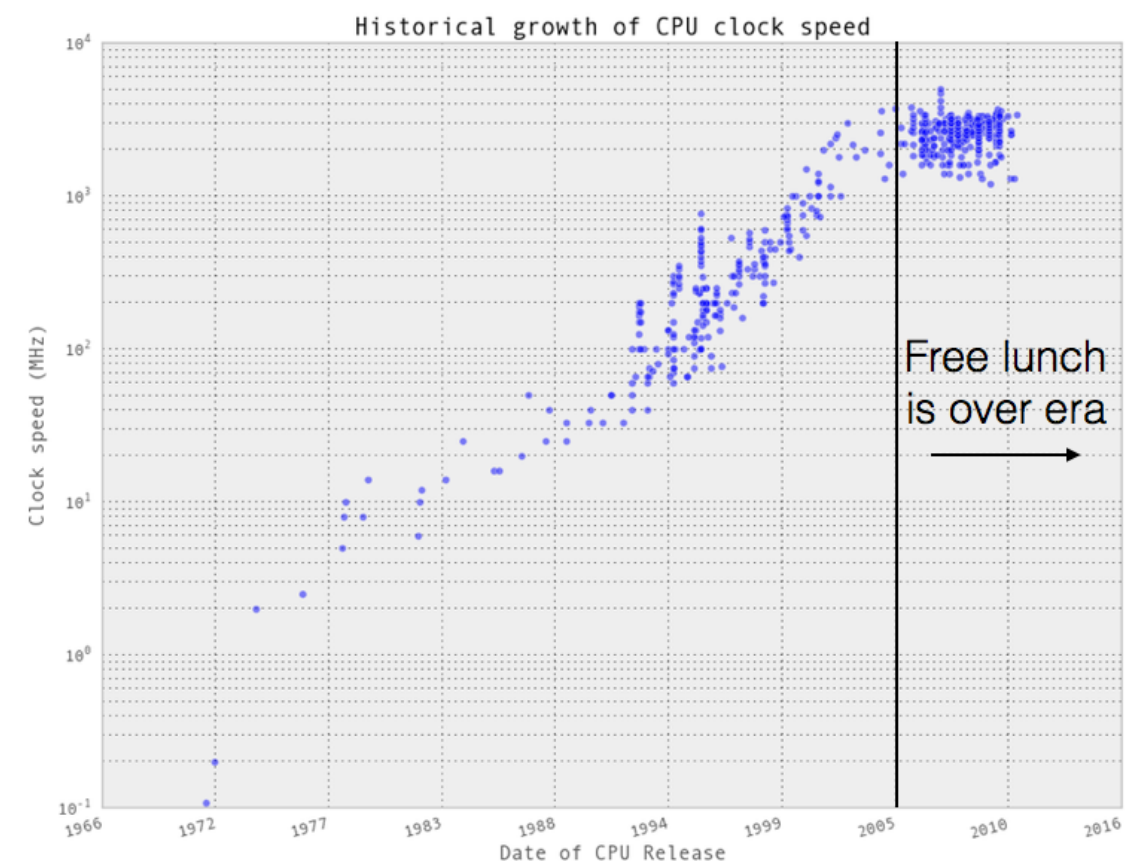


Microprocessor Transistor Counts 1971-2011 & Moore's Law



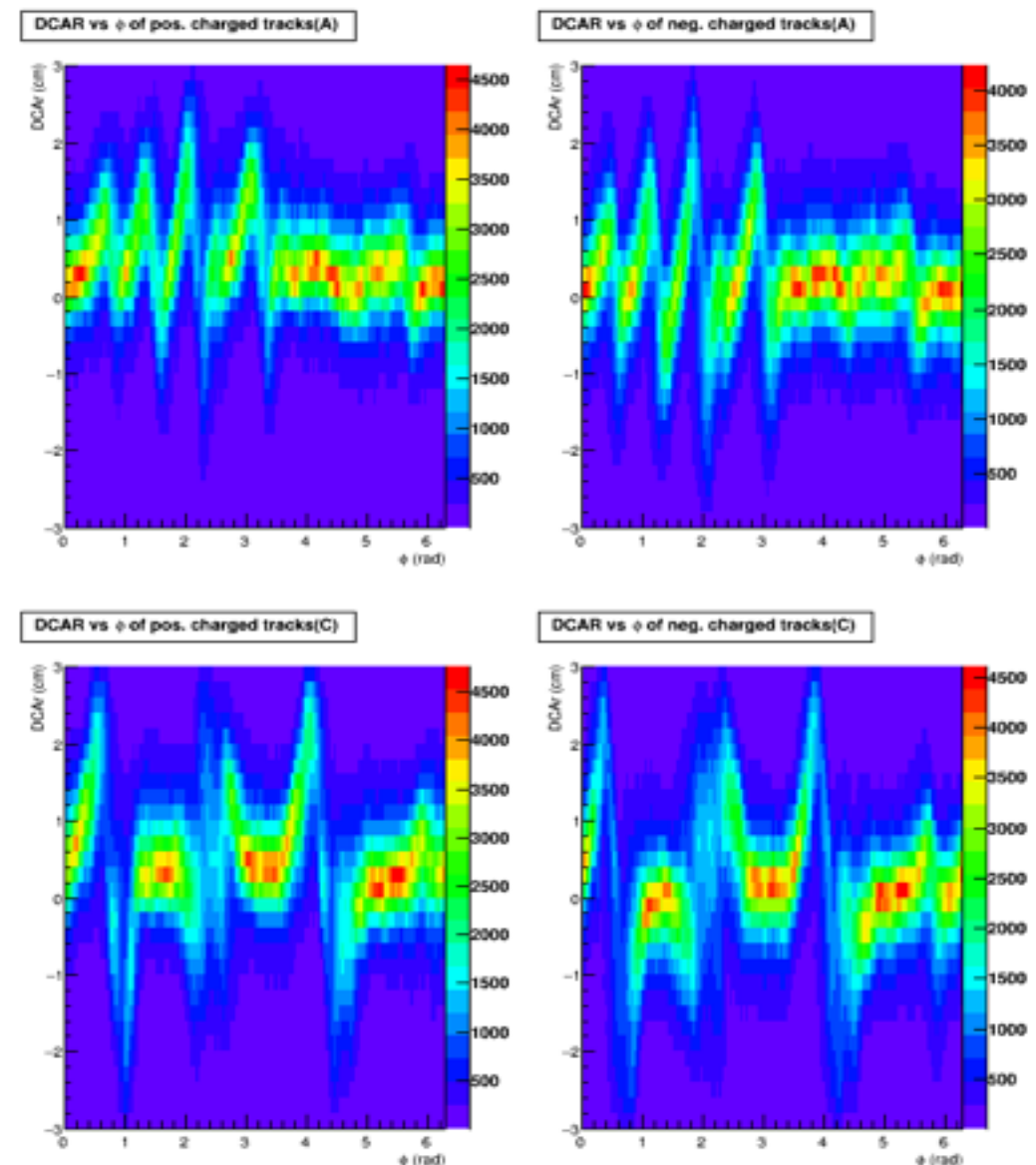
Why HPC in High energy Physics?

- Our software profited of an exponential growth of the CPU clock frequency.
- This growth ended about 10 years ago.
- CPU performance is still growing since the number of cores/CPU is growing.
- Due to the intrinsic parallelism of our data, we exploited this core growth by increasing the number of jobs/CPU.
- Considerations related to real time needs (online processing) and to budget evaluations (in terms of number of worker nodes and power/cooling costs) are pushing our community towards HPC solutions
- We are just at the beginning!



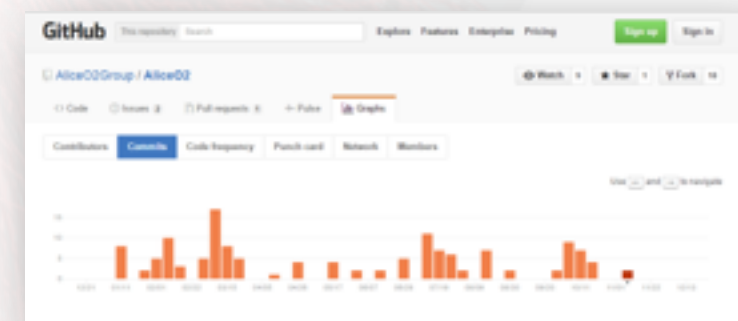
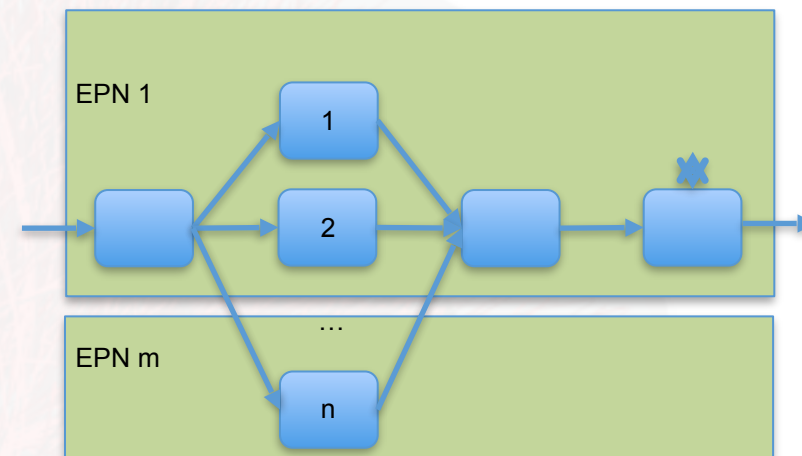
Status of 2015 data processing (2)

- Substantial IR-induced distortions in the TPC
- Affect both p-p and Pb-Pb data
- Sophisticated correction algorithms development in the past 6 months
- Data reconstructed partially (first physics, Lower IR runs)
- Bulk of reconstruction still pending



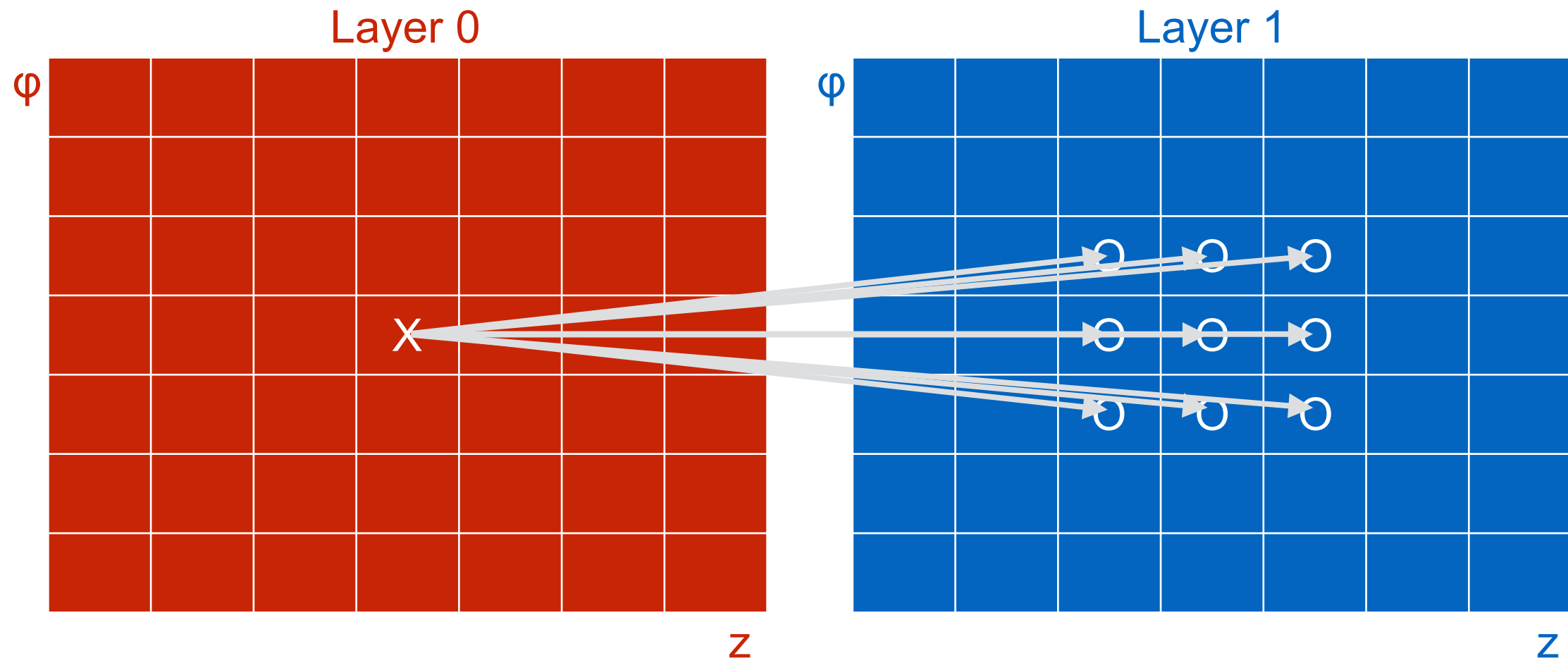
ALFA

- Large monolithic programs are divided in tasks
- Each task can
 - Run on multiple or different hardware (CPU or hw accelerator)
 - Be written in any of the supported language
 - Be multi-threaded if need be
- More documentation and examples
- New API for data serialization
 - Support of multipart messages
- Ongoing development
 - parameter manager
 - integration of key-value DB for parameter management
 - new libfabric based transport for NanoMsg protocol
- Dynamic Deployment System (DDS) prototype
 - Prototype released on November 20th
 - Available on <https://github.com/FairRootGroup/DDS/blob/1.0/ReleaseNotes.md>



Pb-Pb @ $\sqrt{s} = 2.76$ ATeV
2011-11-12 06:51:12
Fill : 2290
Run : 167693
Event : 0x3d94315a

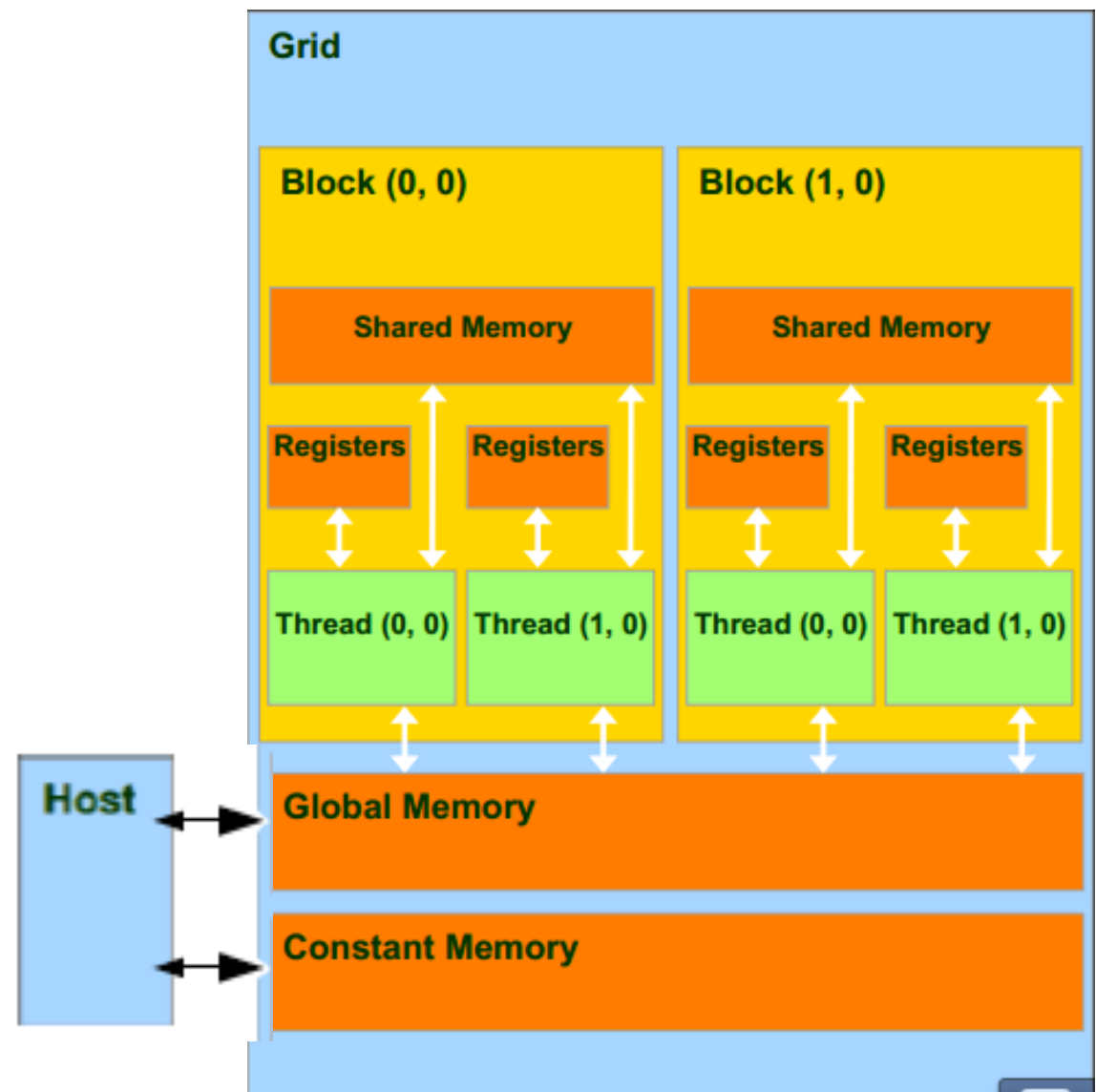
How to speed up: grids to browse data



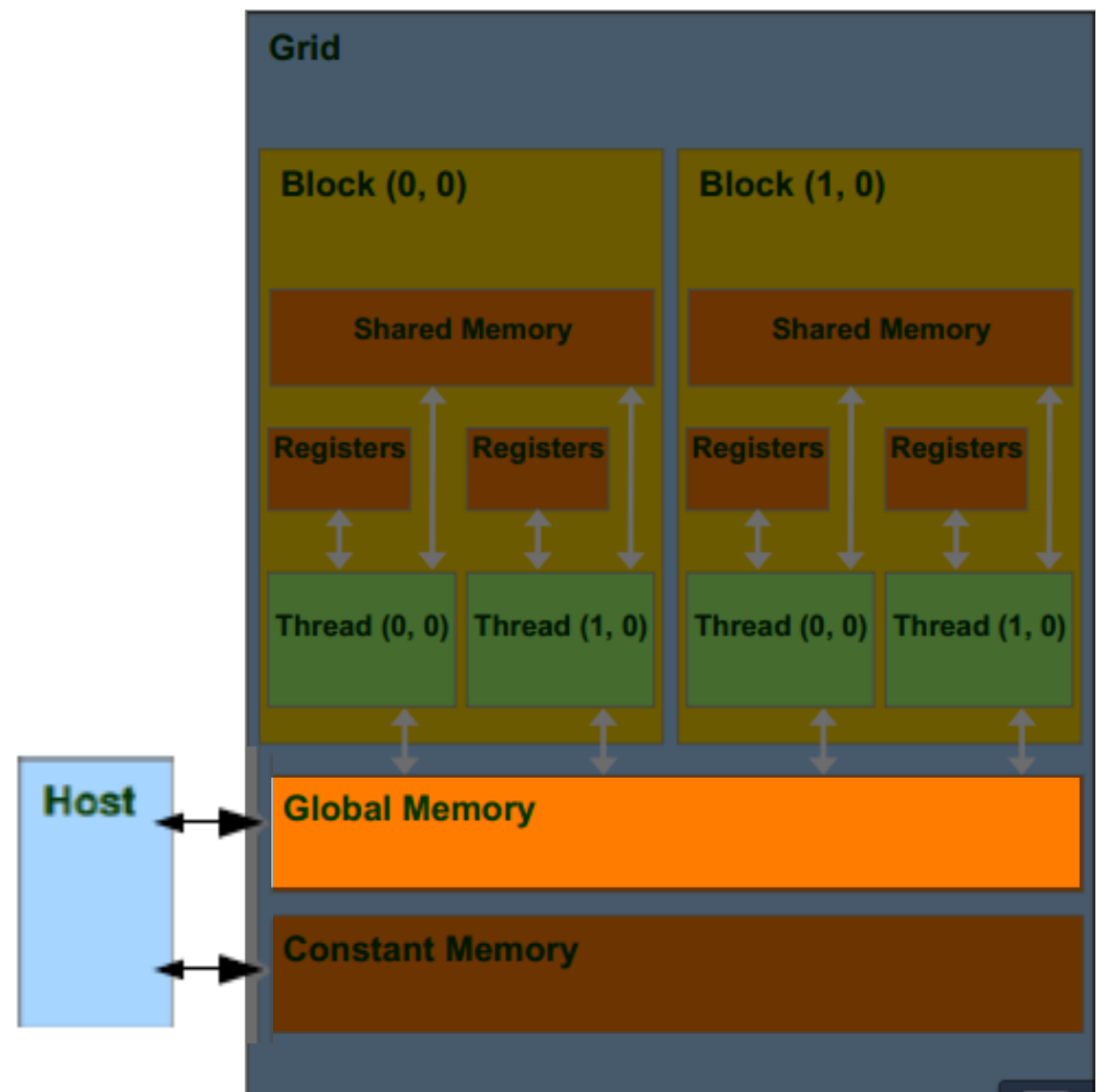
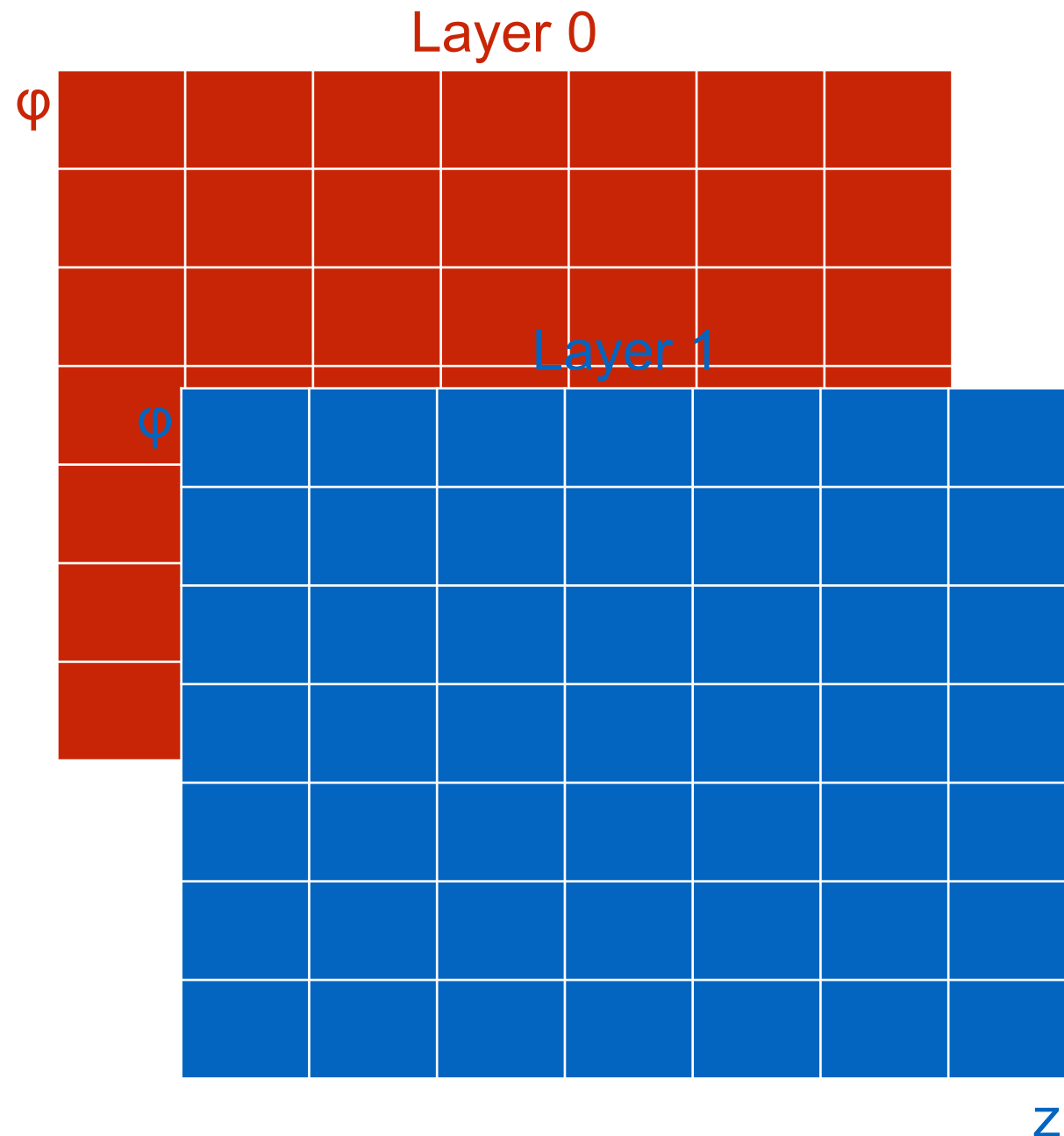
Lookup table approach. It permits:

- **Fine grained parallelism:** for each arrow connecting layer0 to layer1 a compute element is used
- **Coarse grained parallelism:** for each cell/row of layer 0 a compute element is used.

Porting on GPU: how to make doublets

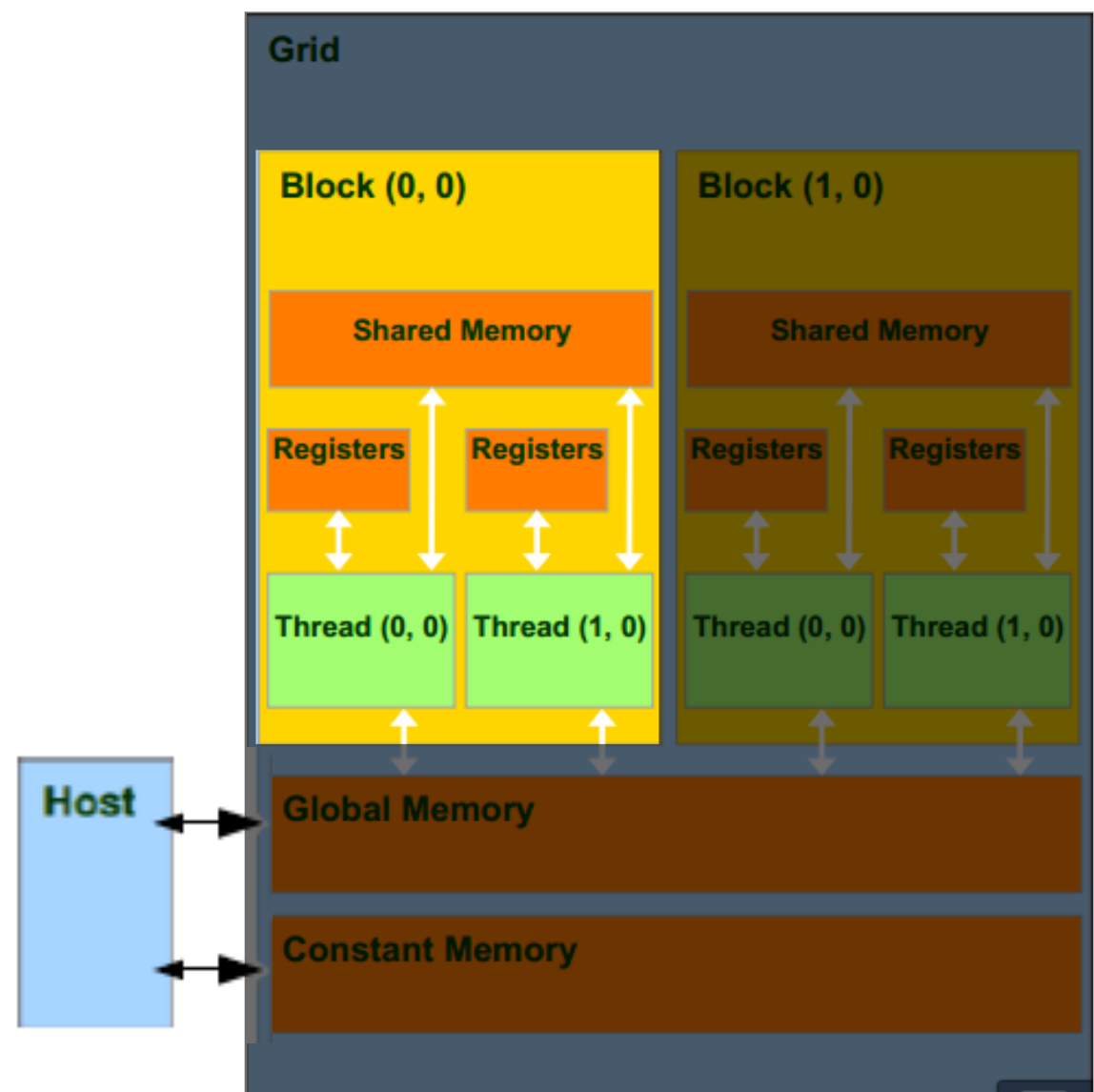
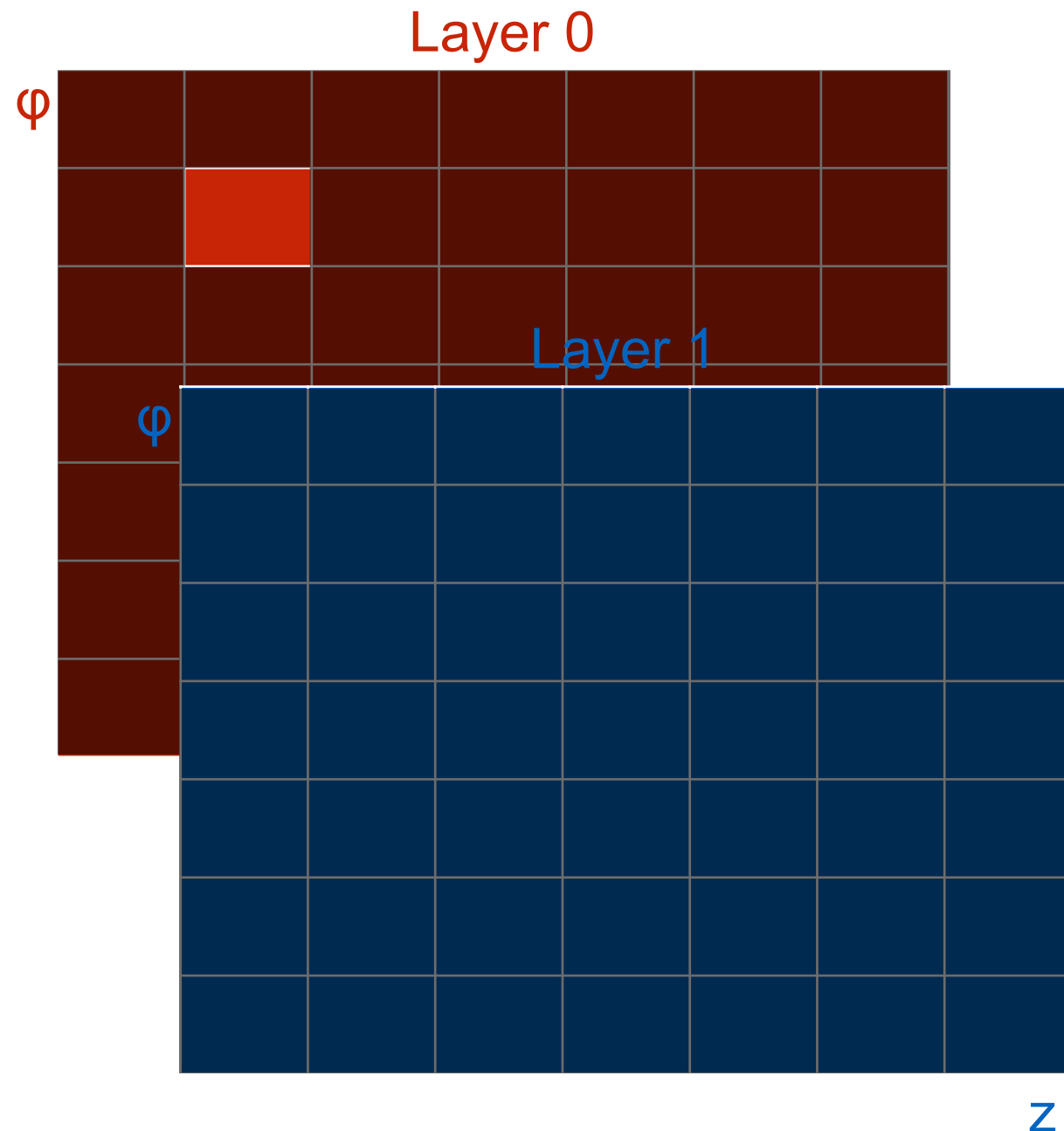


Porting on GPU: how to make doublets



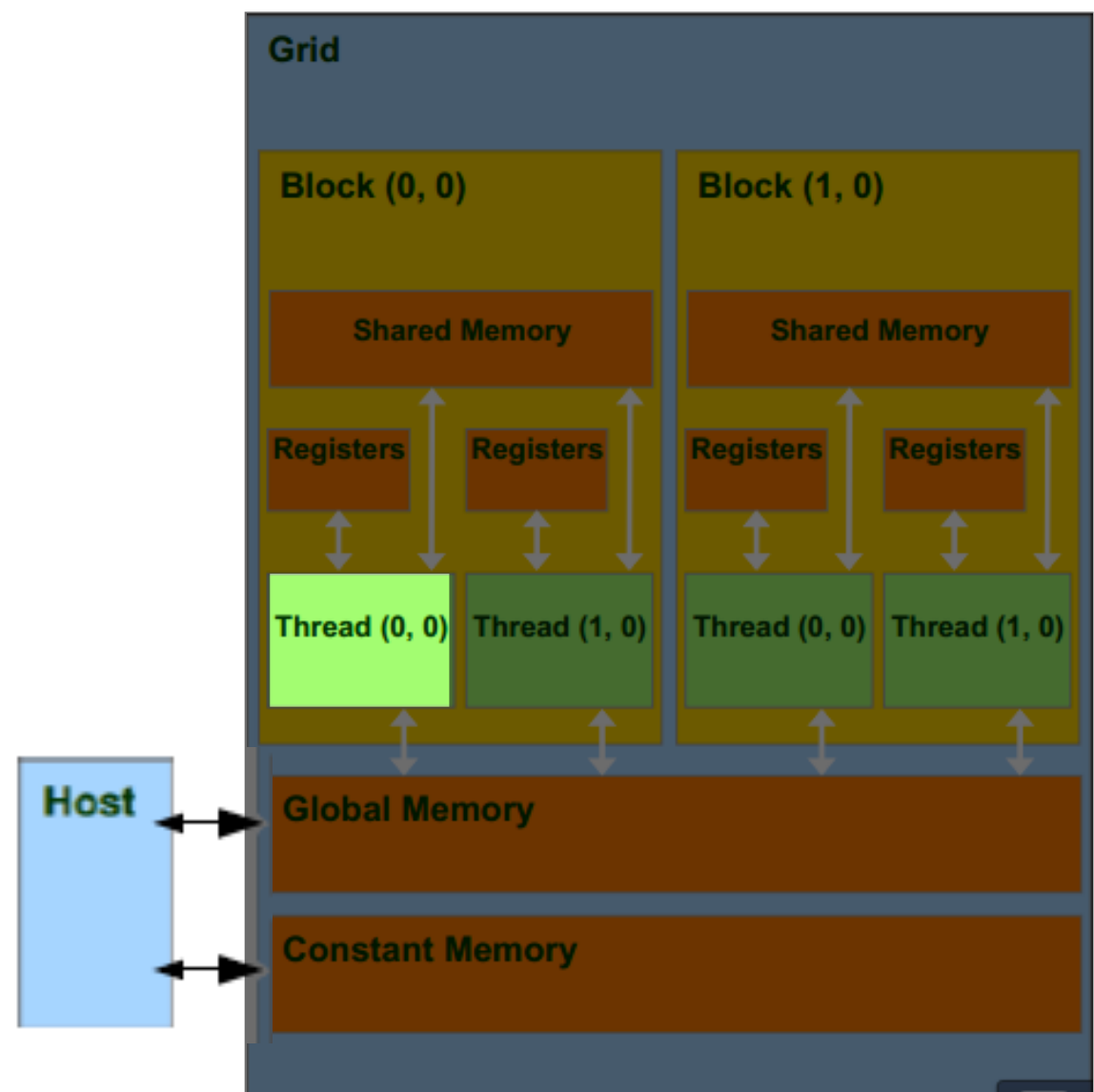
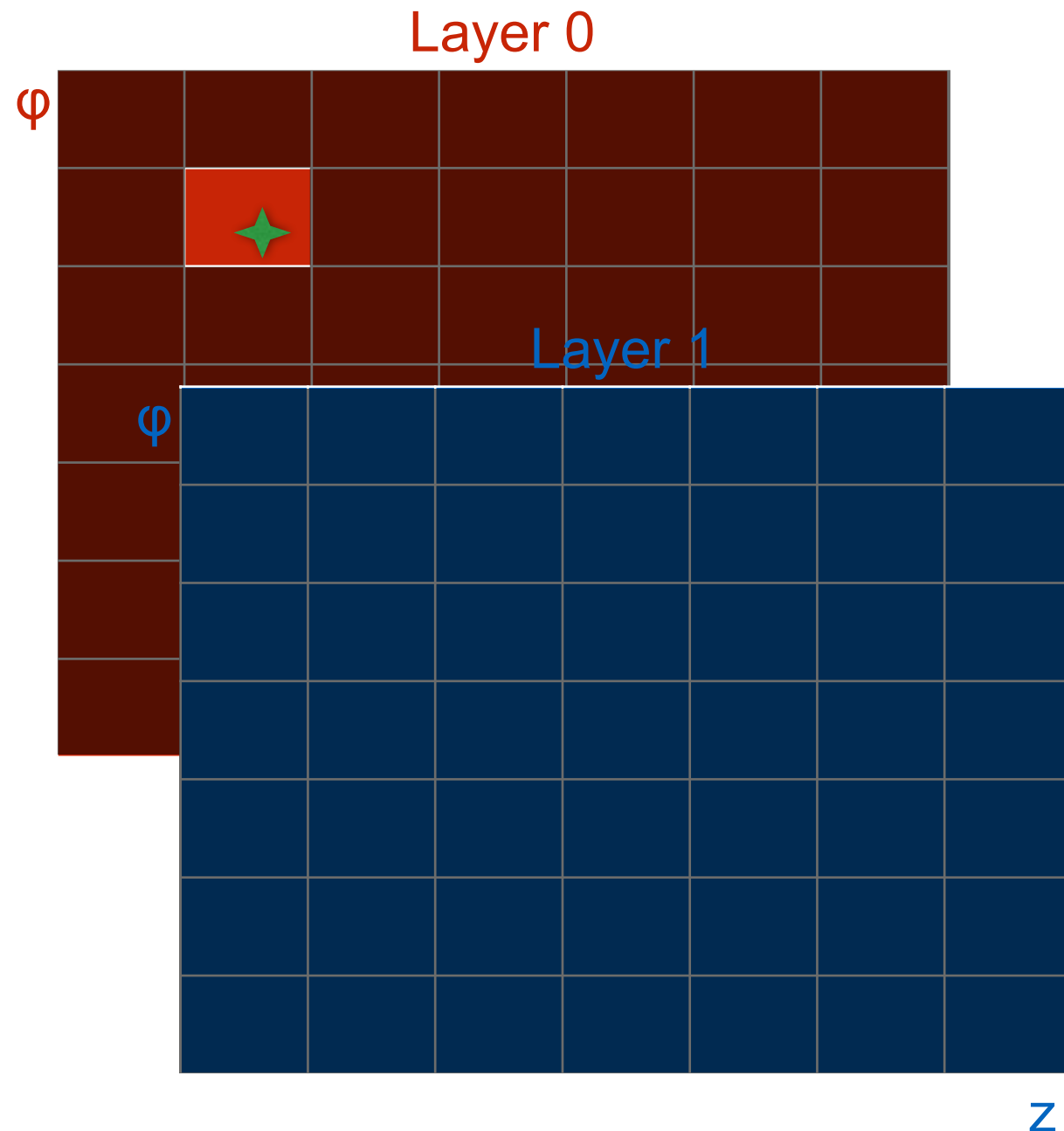
Transfer of the clusters and LUT on the GPU global (constant) memory

Porting on GPU: how to make doublets



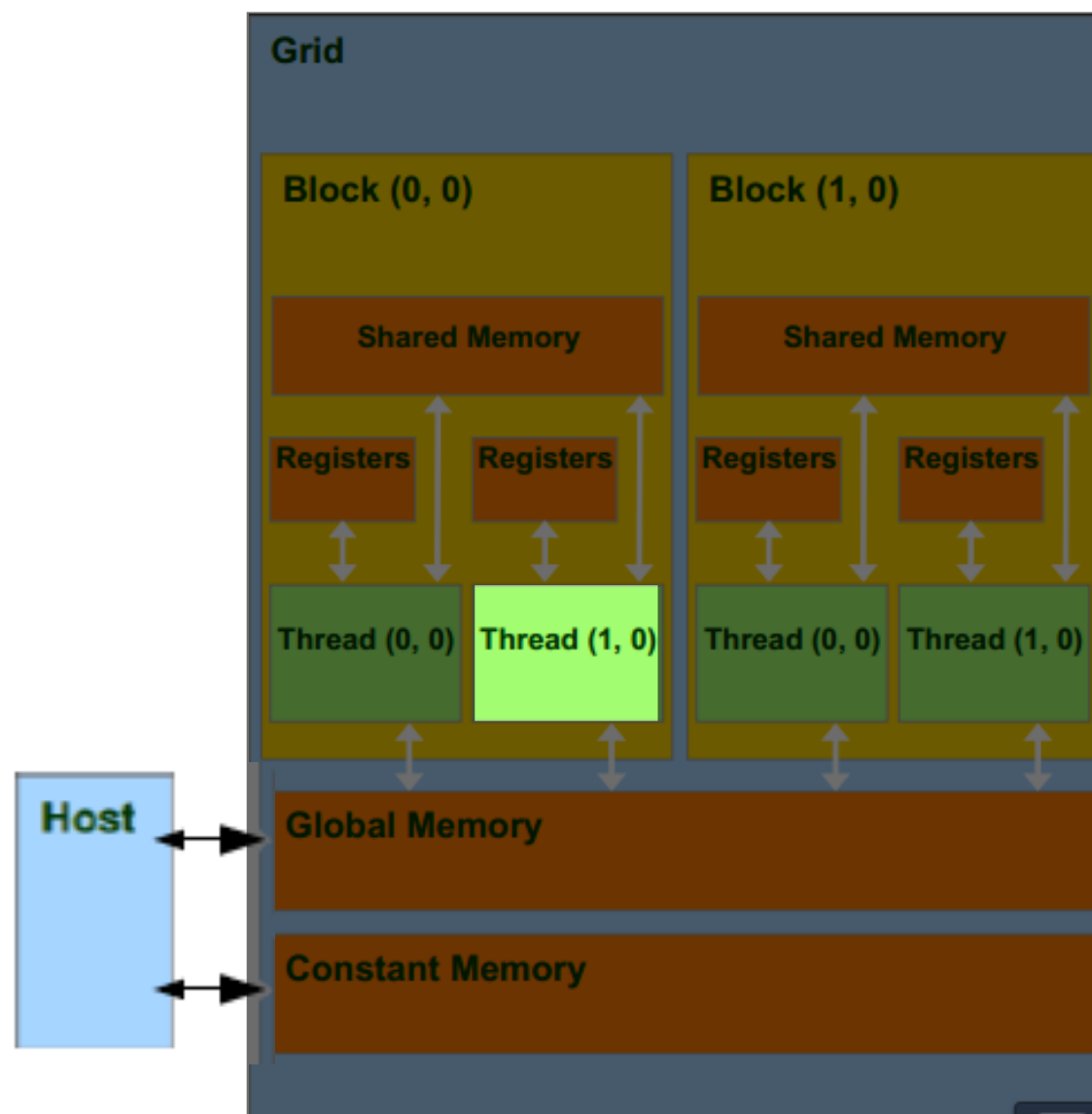
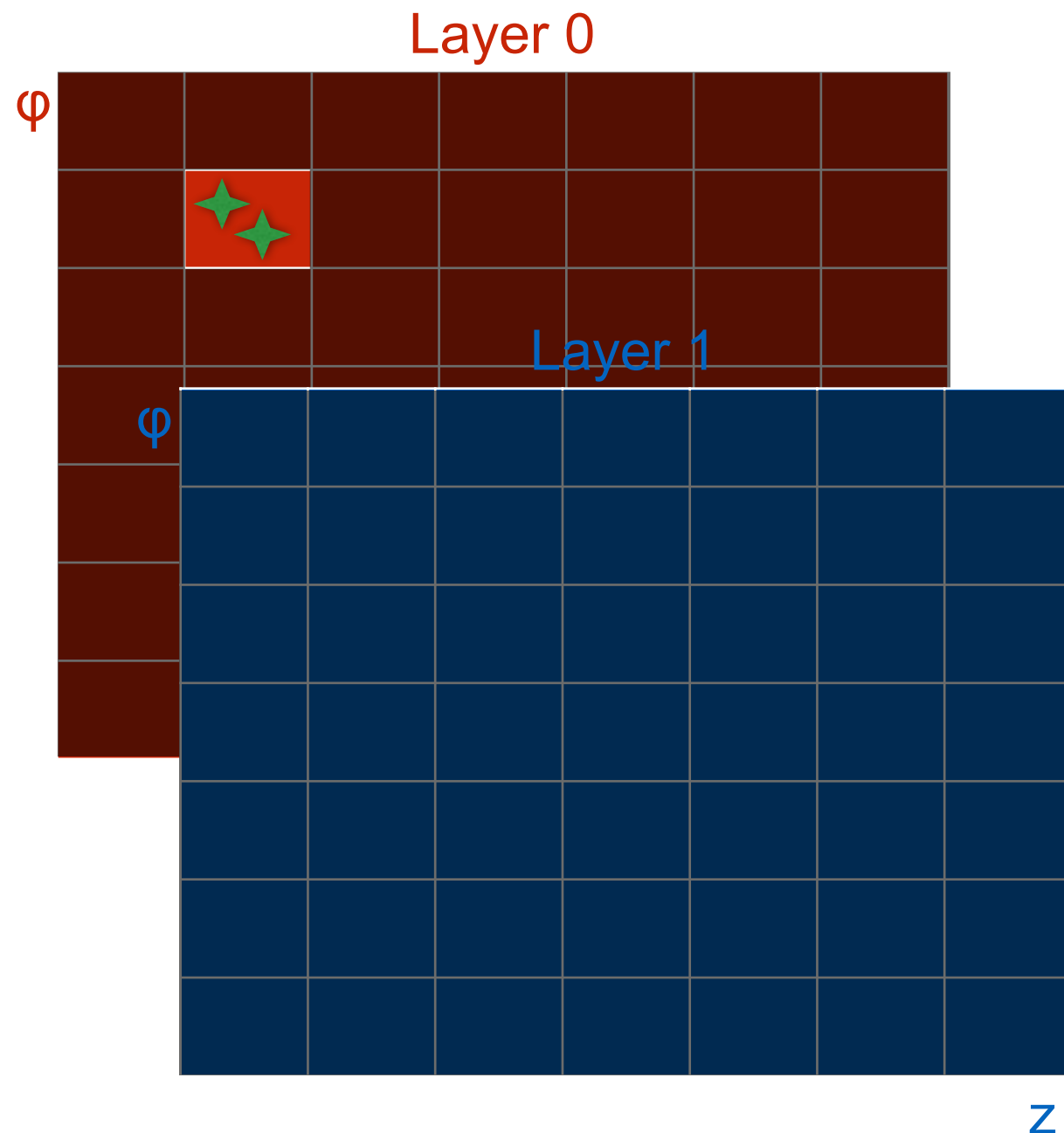
Each block of thread will process one bin on layer 0

Porting on GPU: how to make doublets



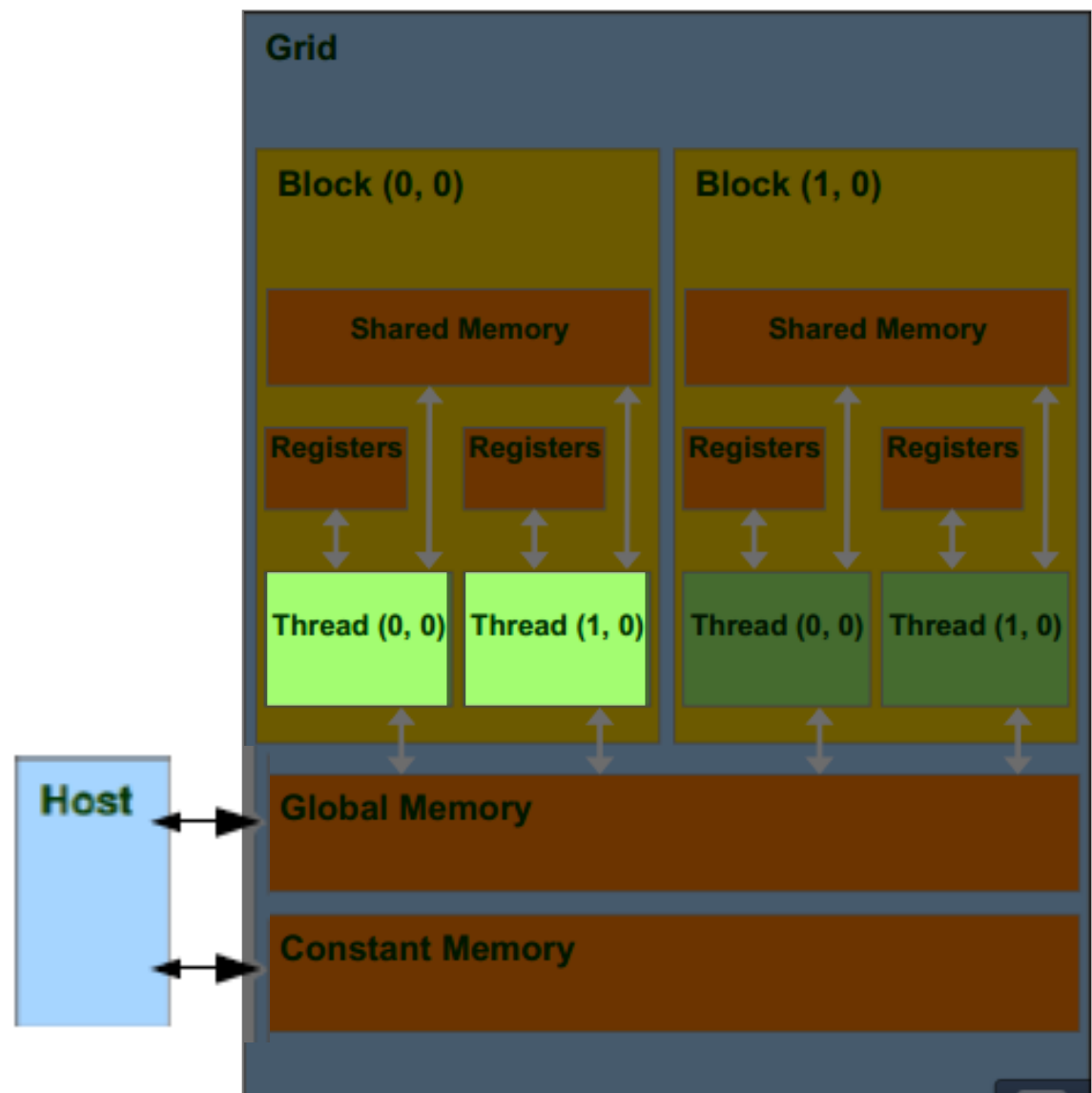
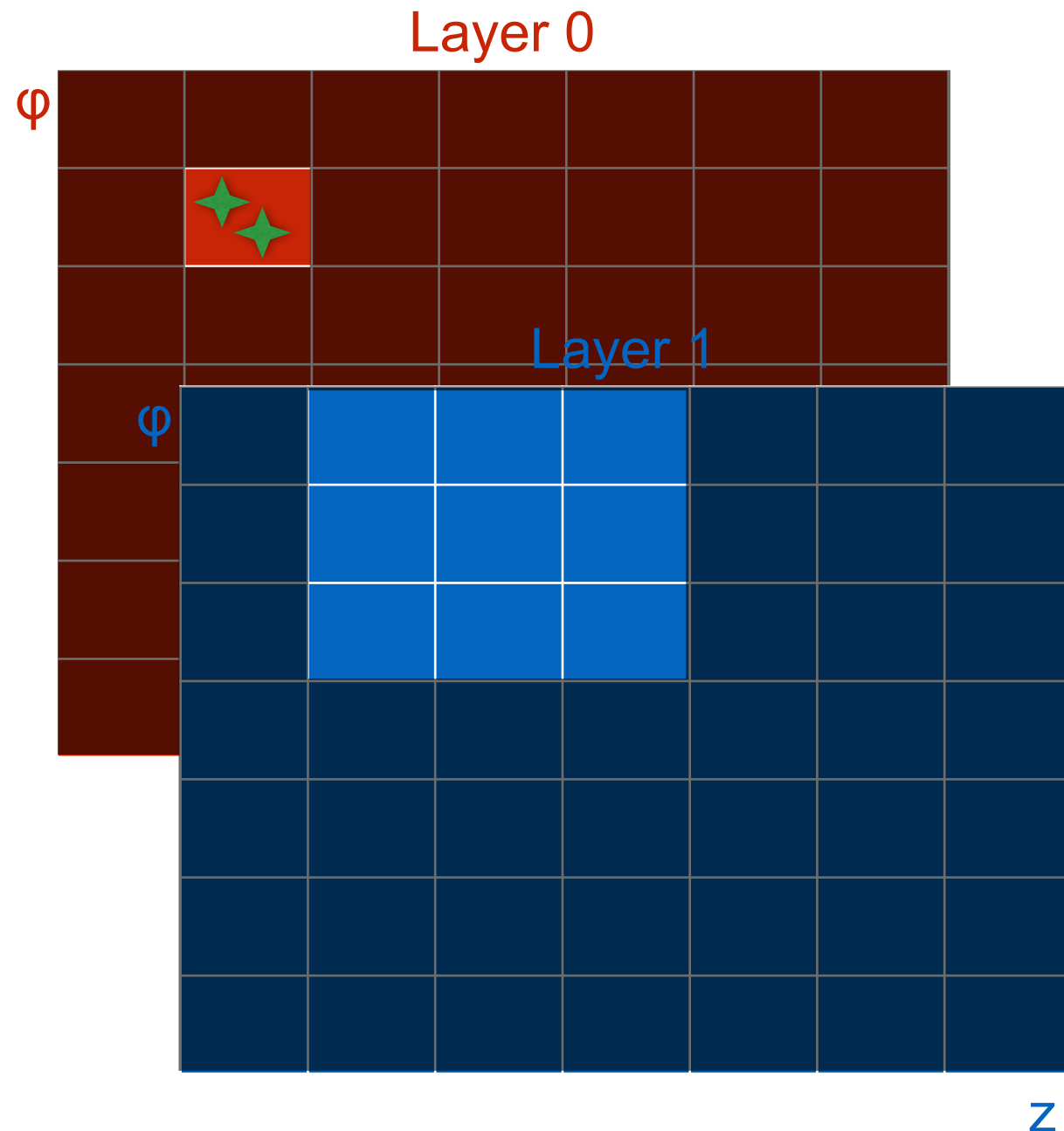
Each thread will take care of one cluster layer 0

Porting on GPU: how to make doublets



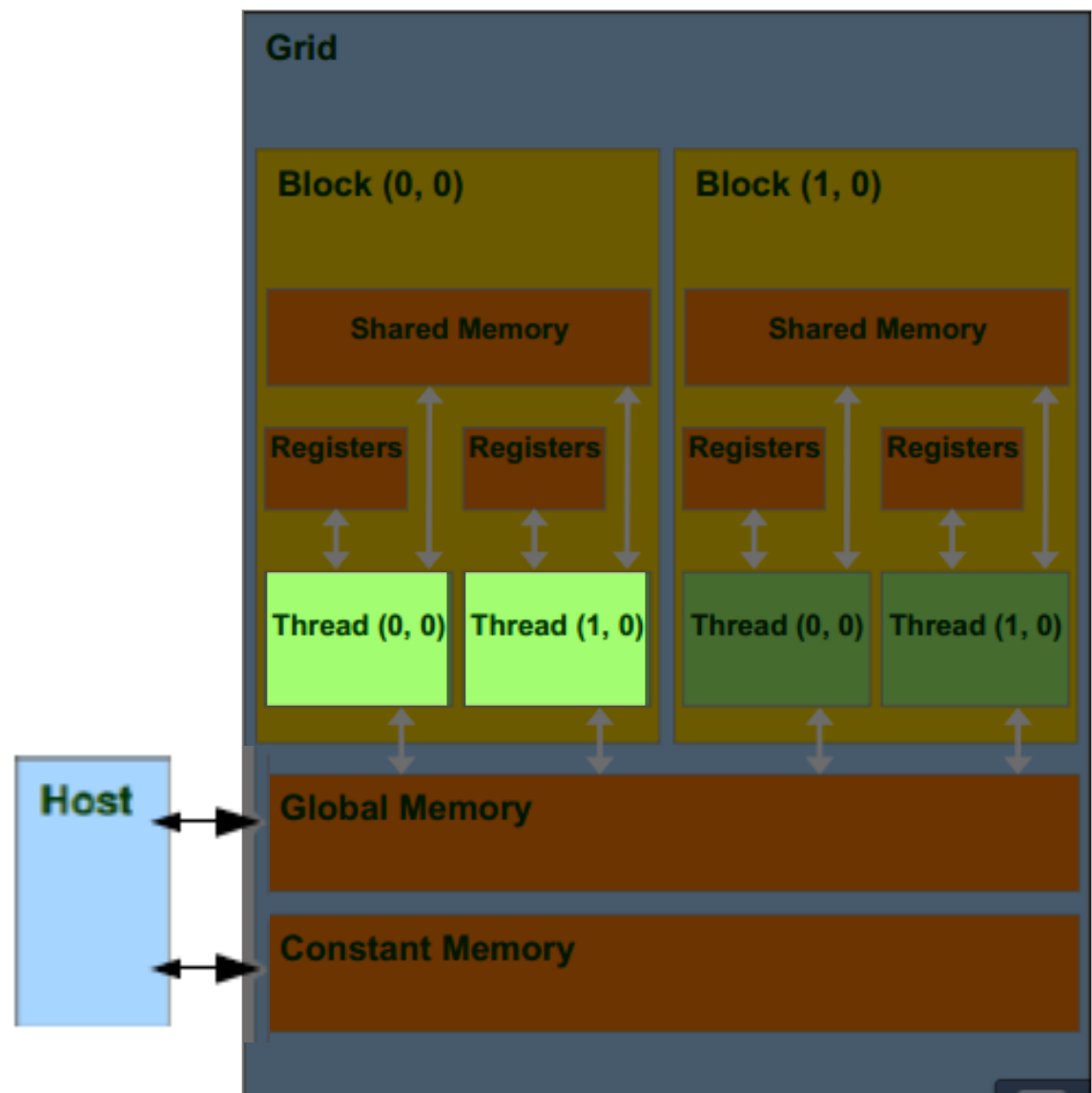
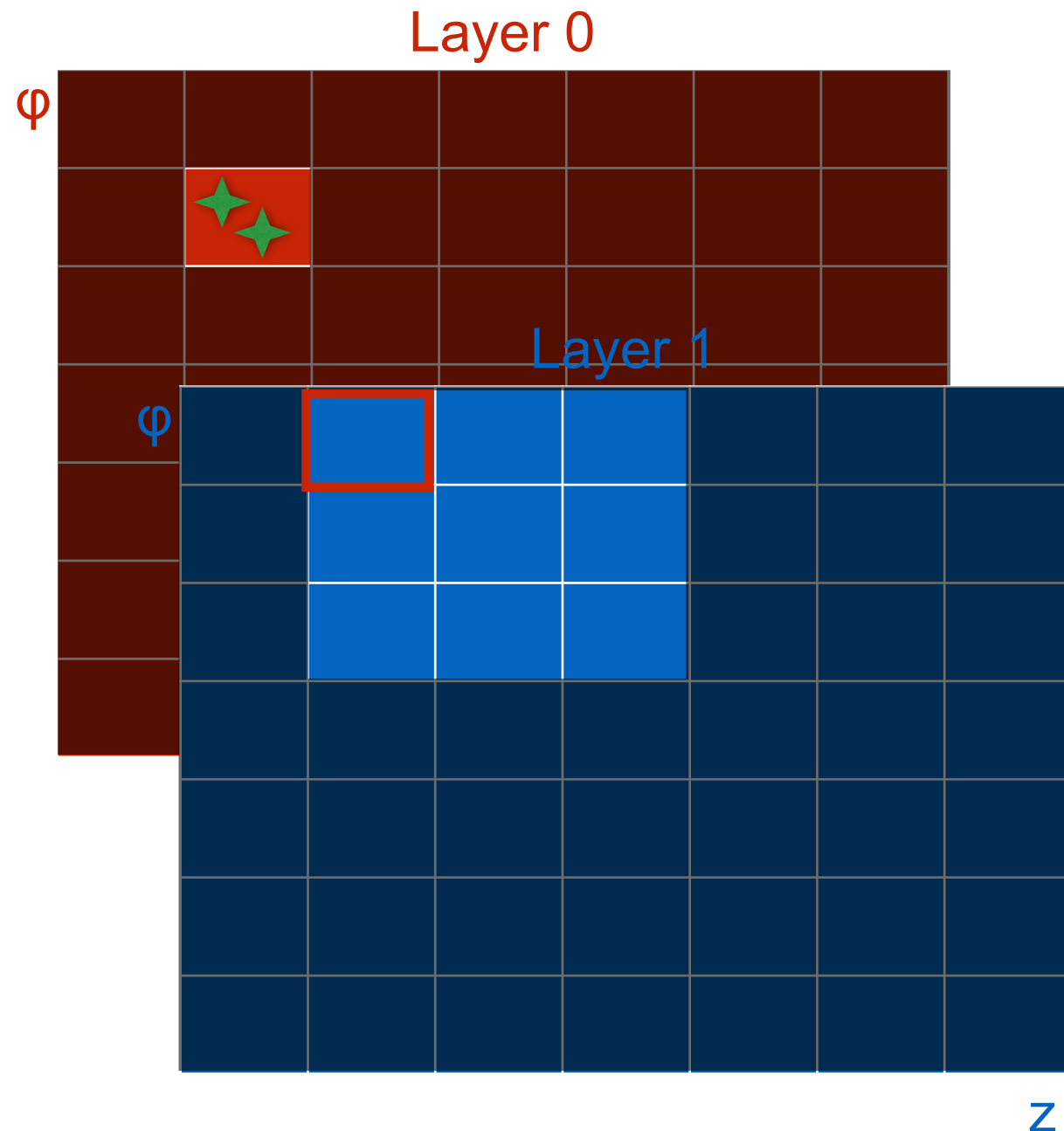
Each thread will take care of one cluster layer 0

Porting on GPU: how to make doublets



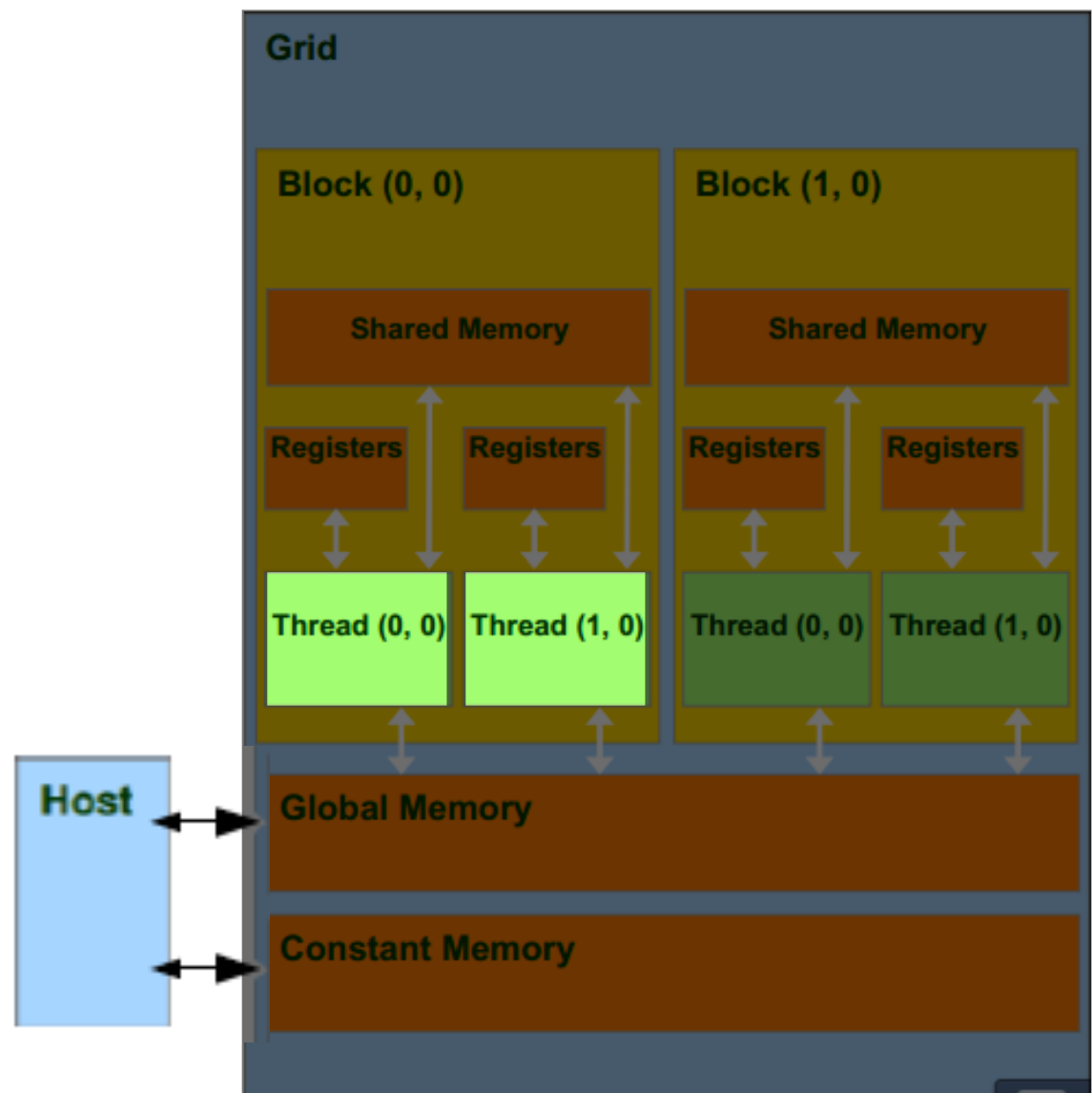
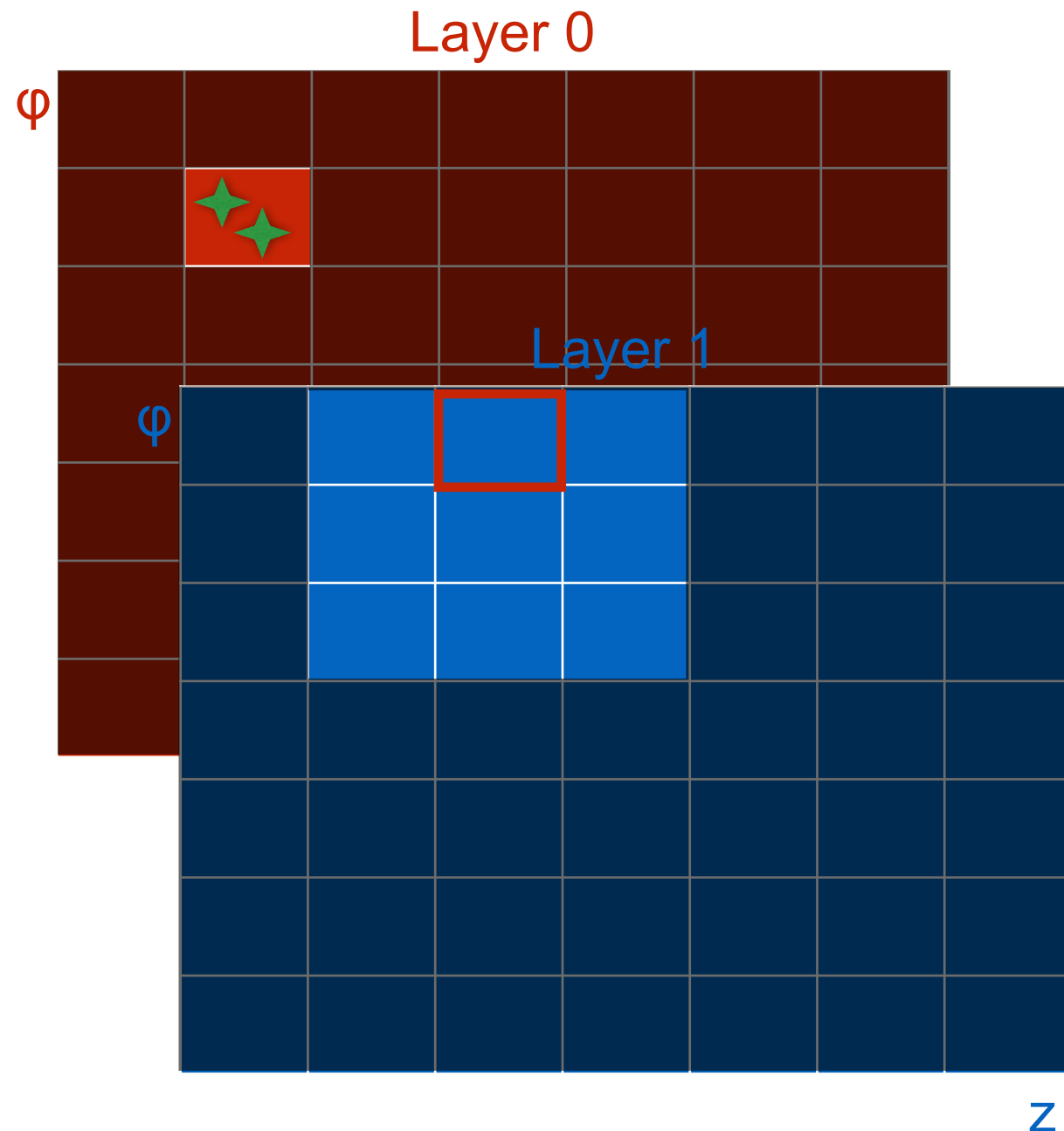
All the clusters of one bin on layer 1 are copied in the shared memory and a for loop starts over these clusters looking for good doublets

Porting on GPU: how to make doublets



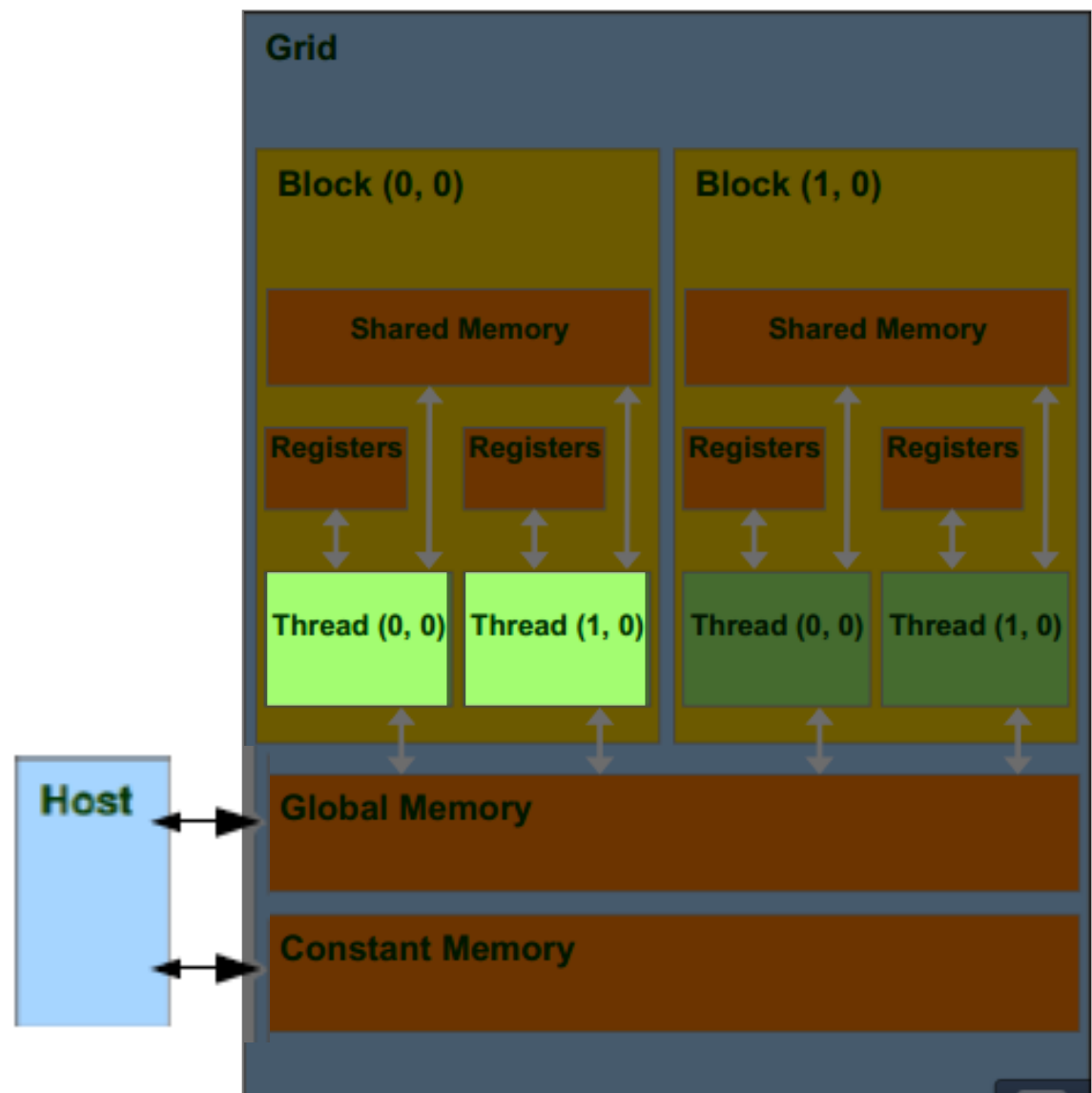
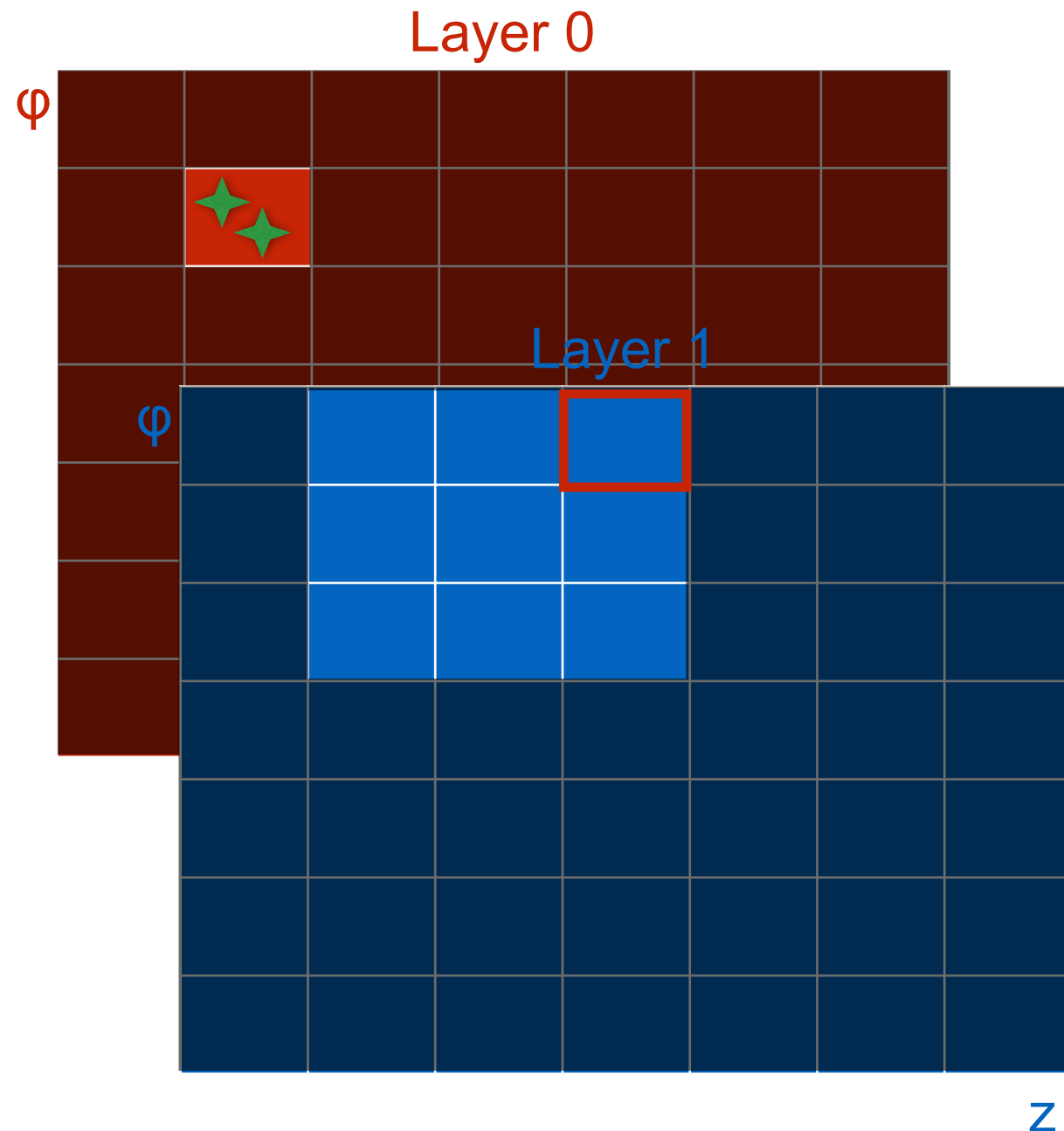
All the clusters of one bin on layer 1 are copied in the shared memory and a for loop starts over these clusters looking for good doublets

Porting on GPU: how to make doublets



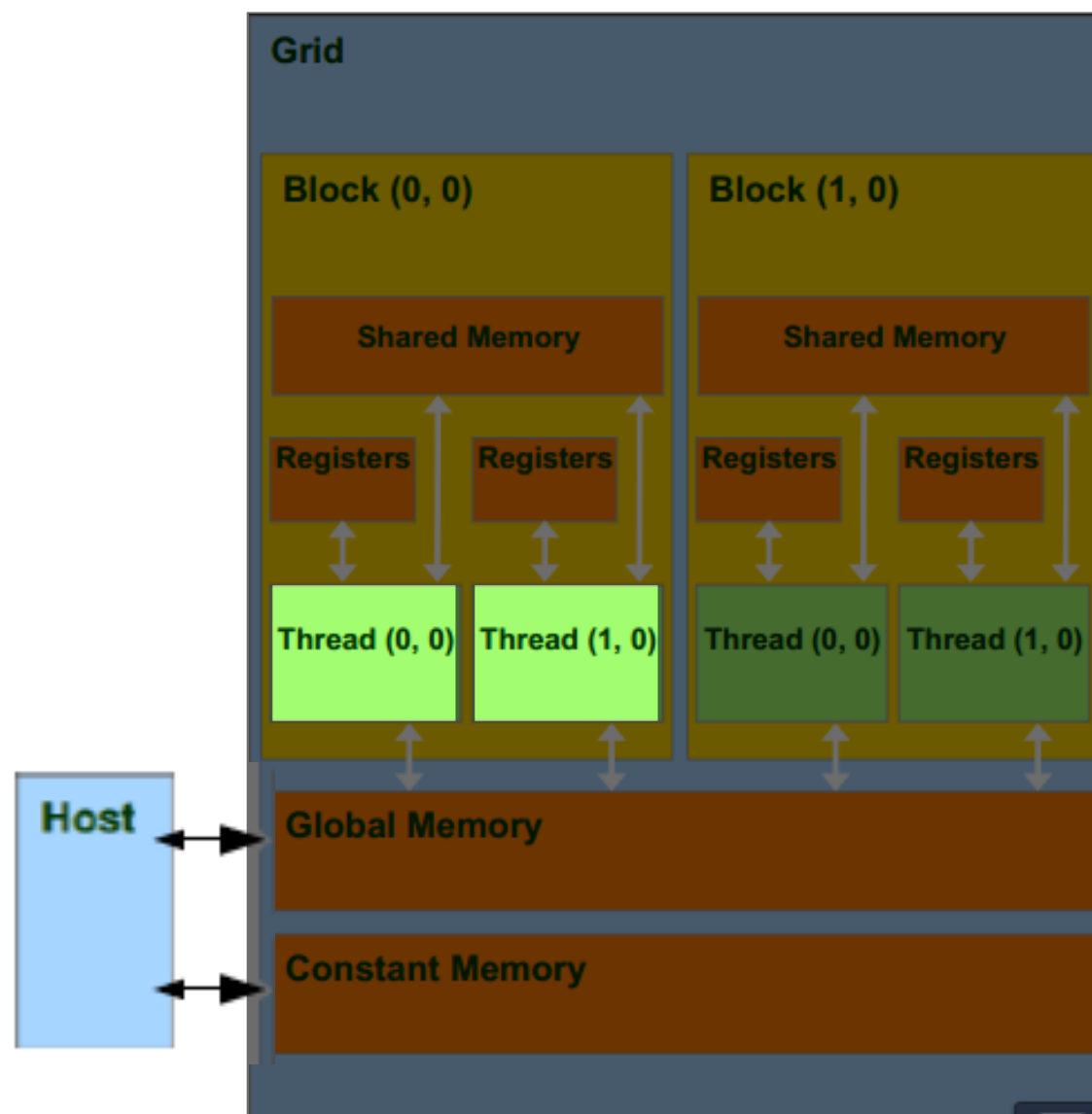
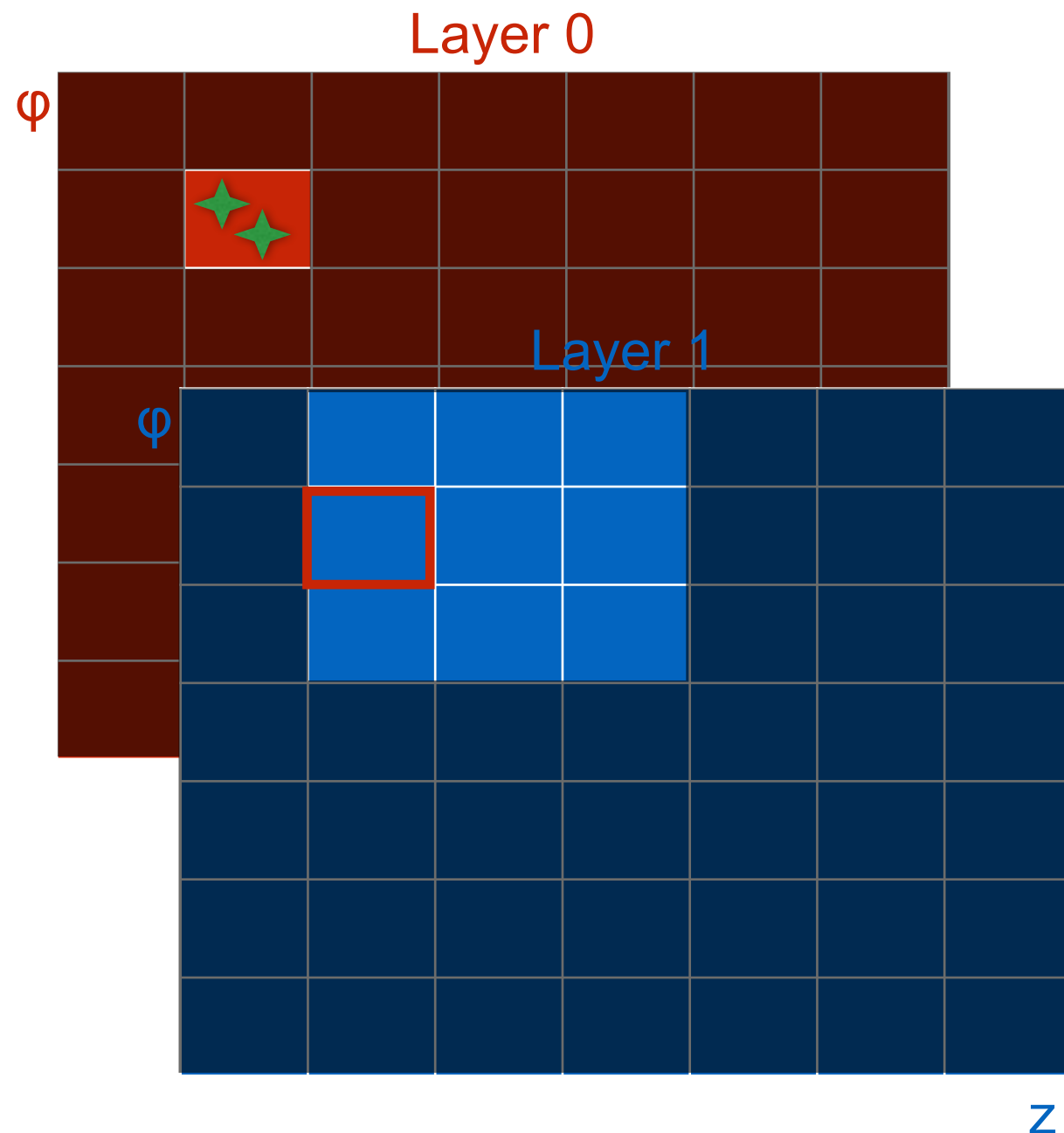
All the clusters of one bin on layer 1 are copied in the shared memory and a for loop starts over these clusters looking for good doublets

Porting on GPU: how to make doublets



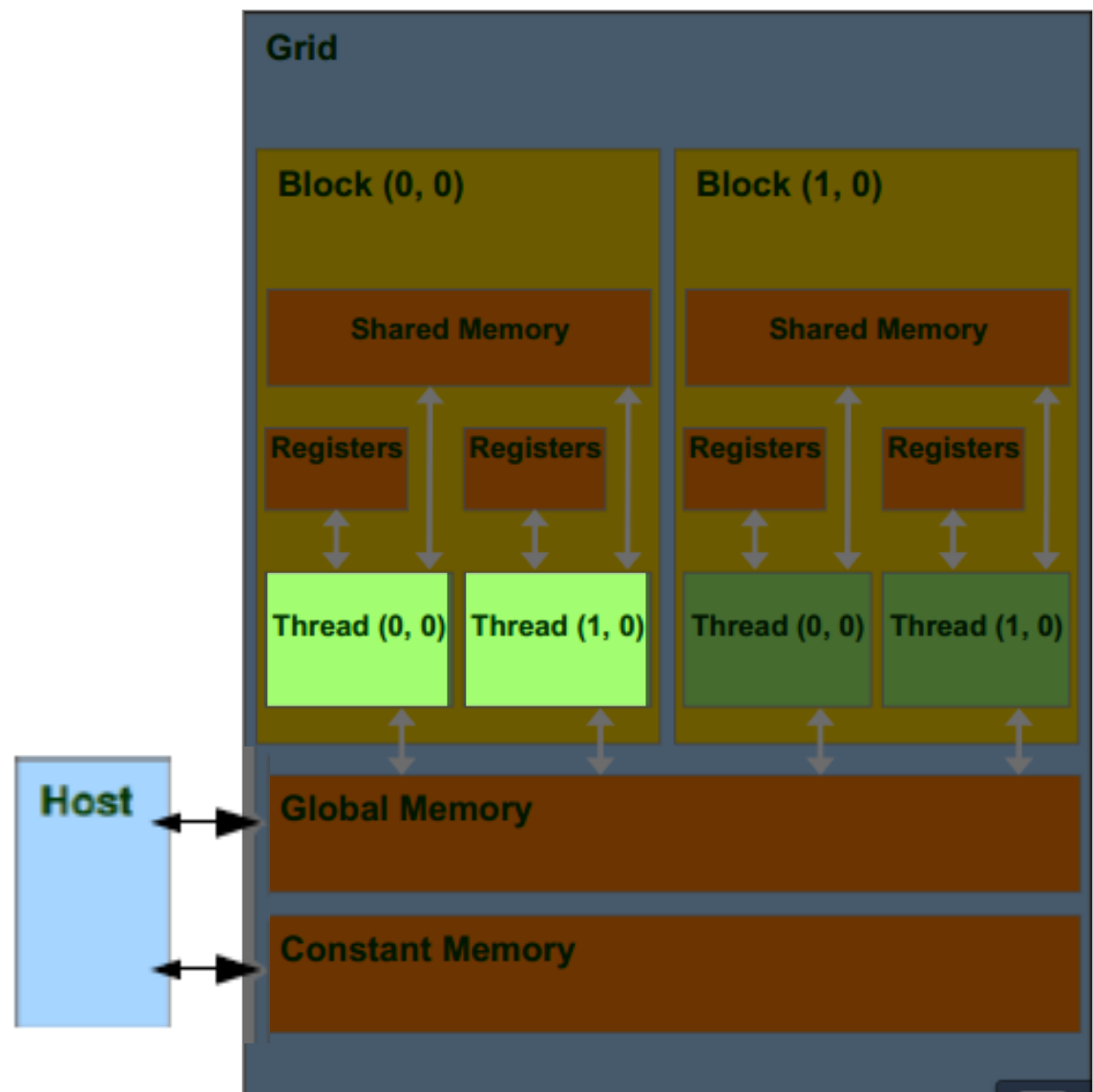
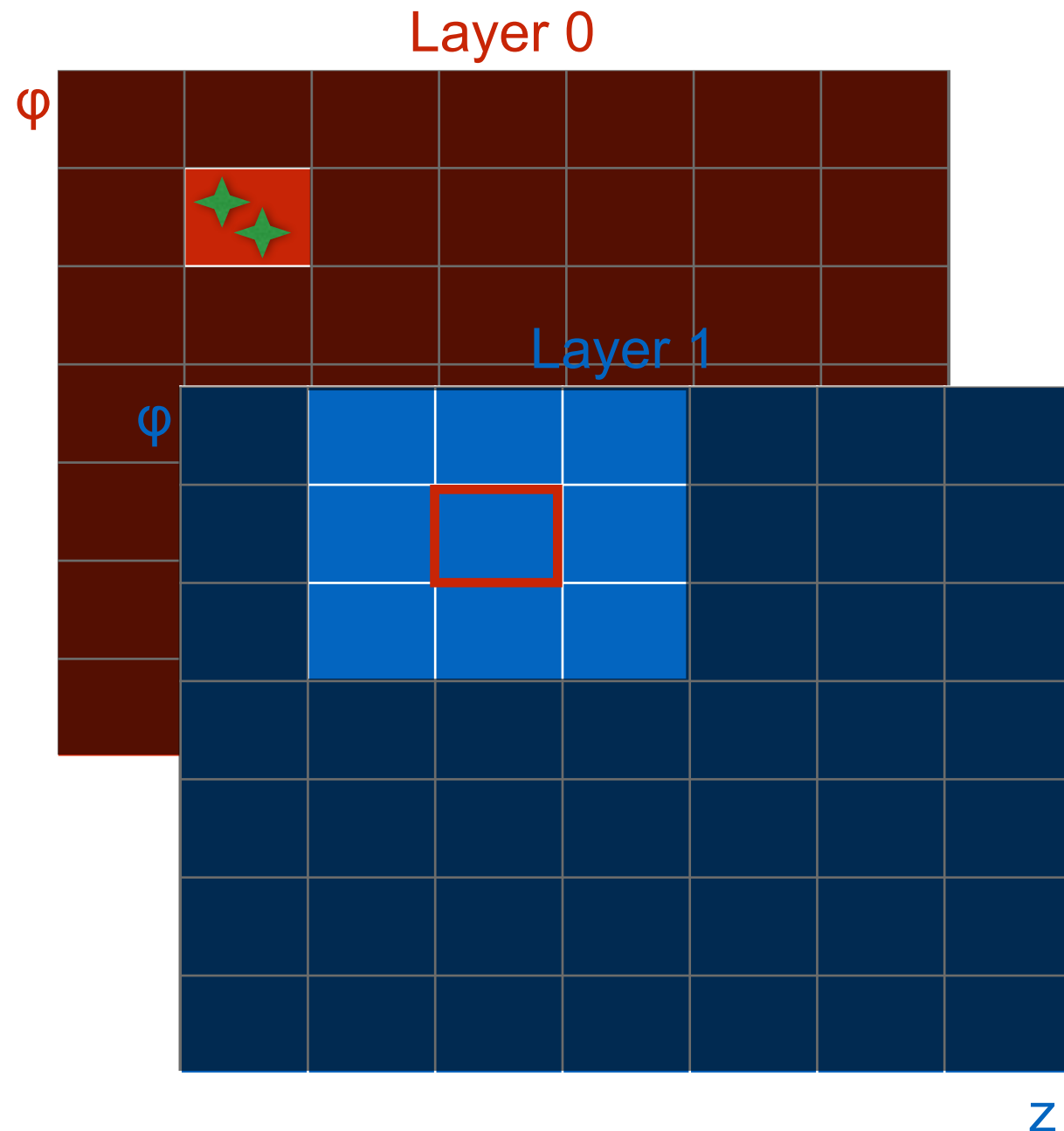
All the clusters of one bin on layer 1 are copied in the shared memory and a for loop starts over these clusters looking for good doublets

Porting on GPU: how to make doublets



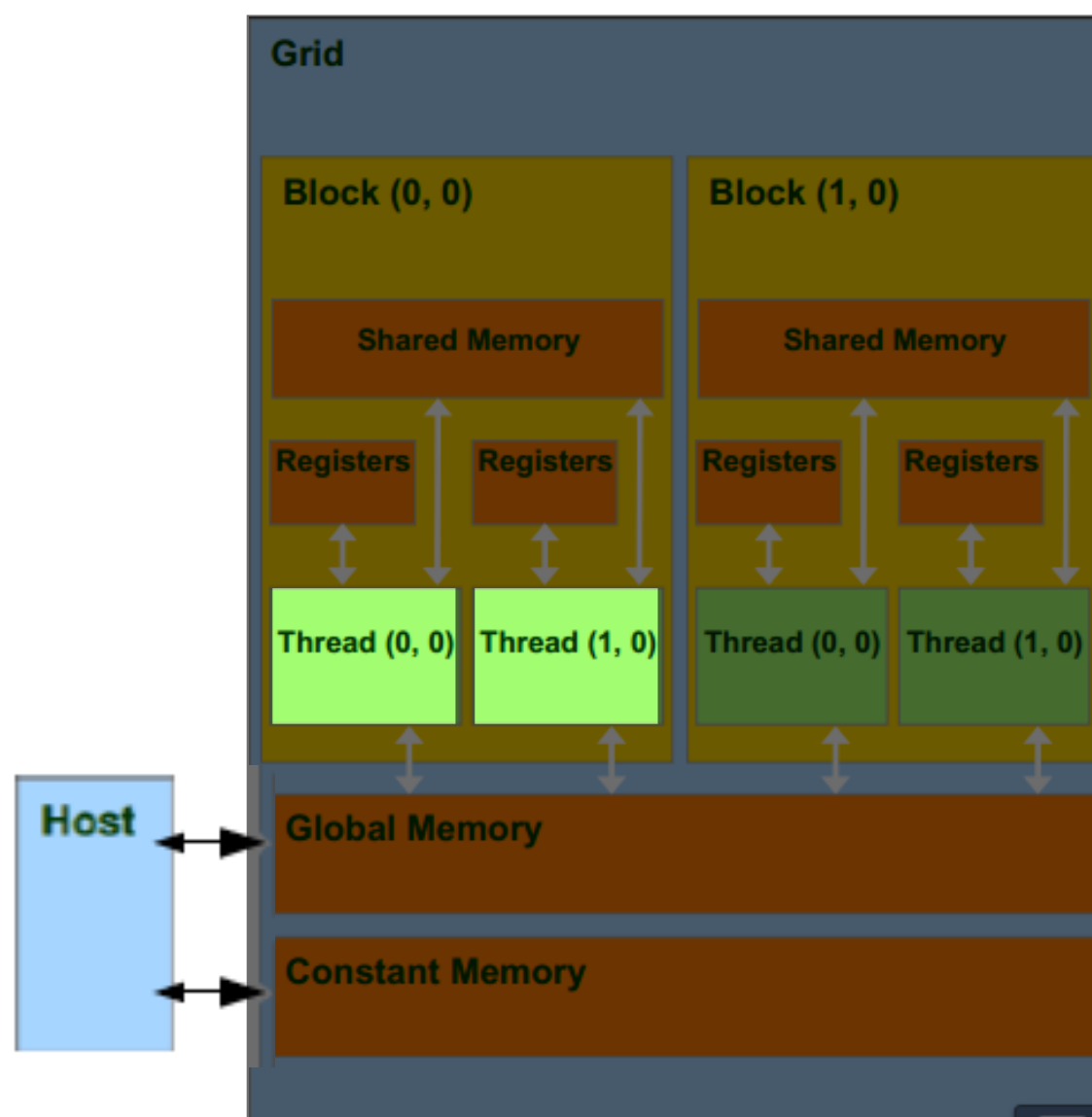
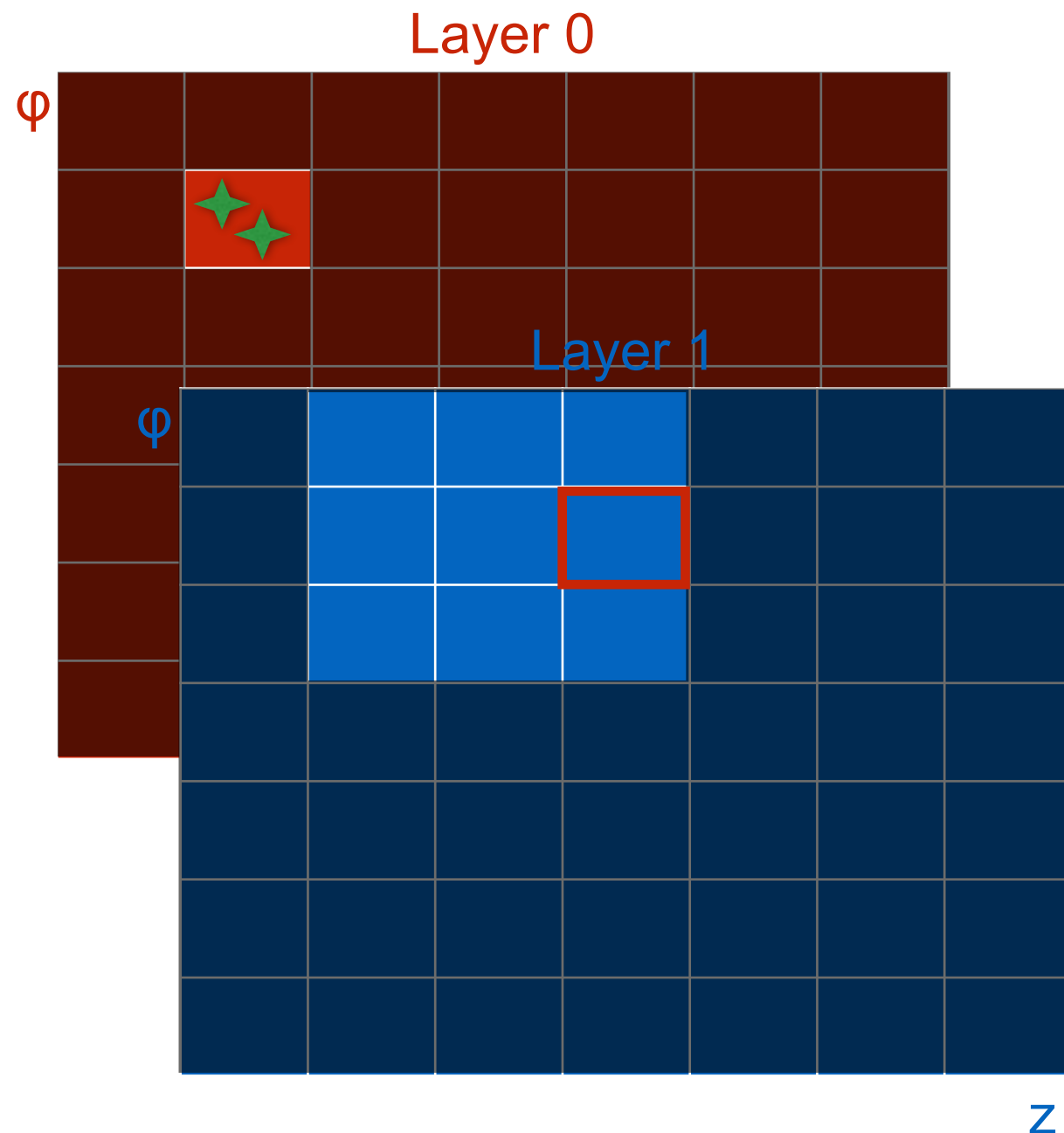
All the clusters of one bin on layer 1 are copied in the shared memory and a for loop starts over these clusters looking for good doublets

Porting on GPU: how to make doublets



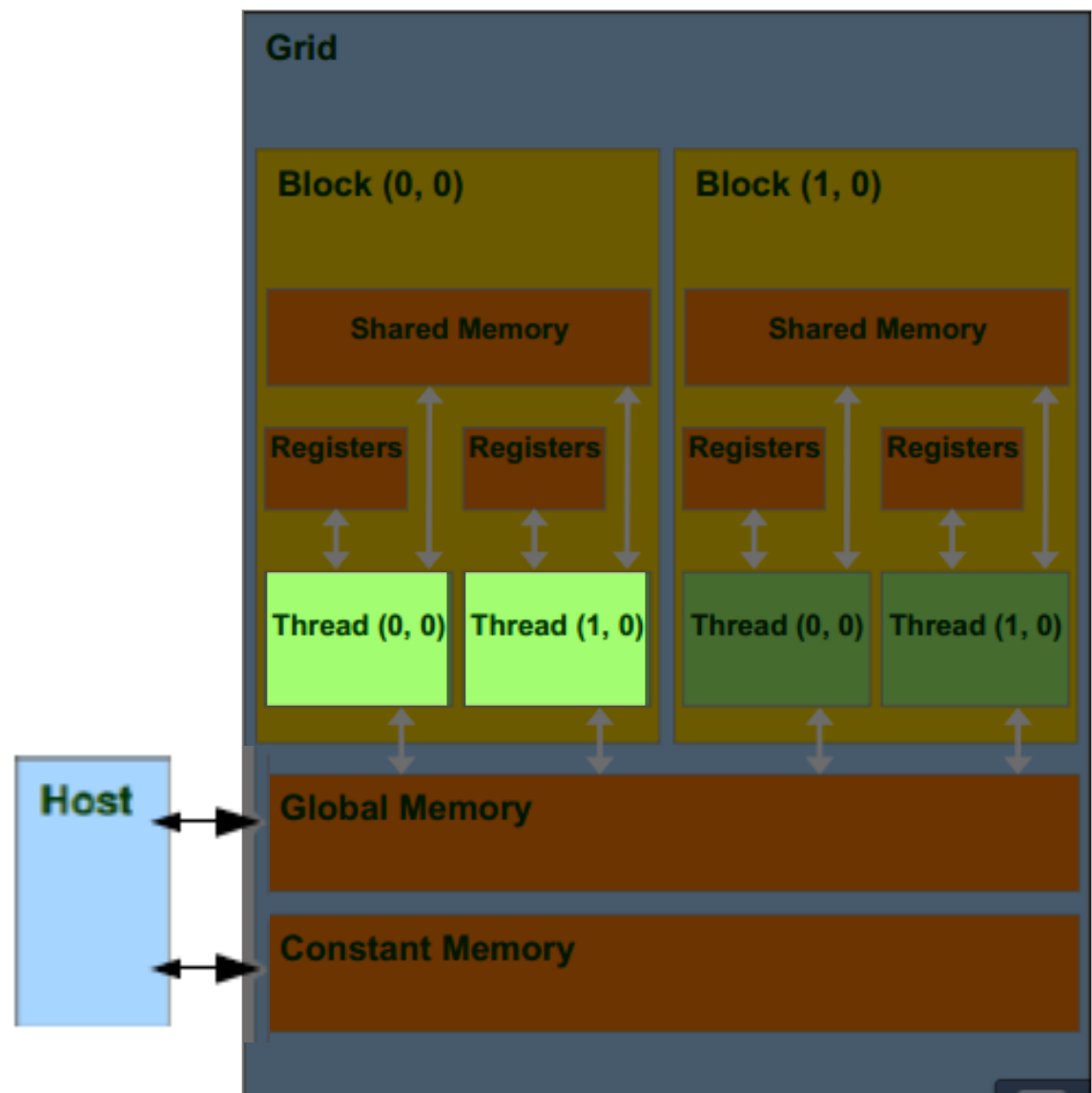
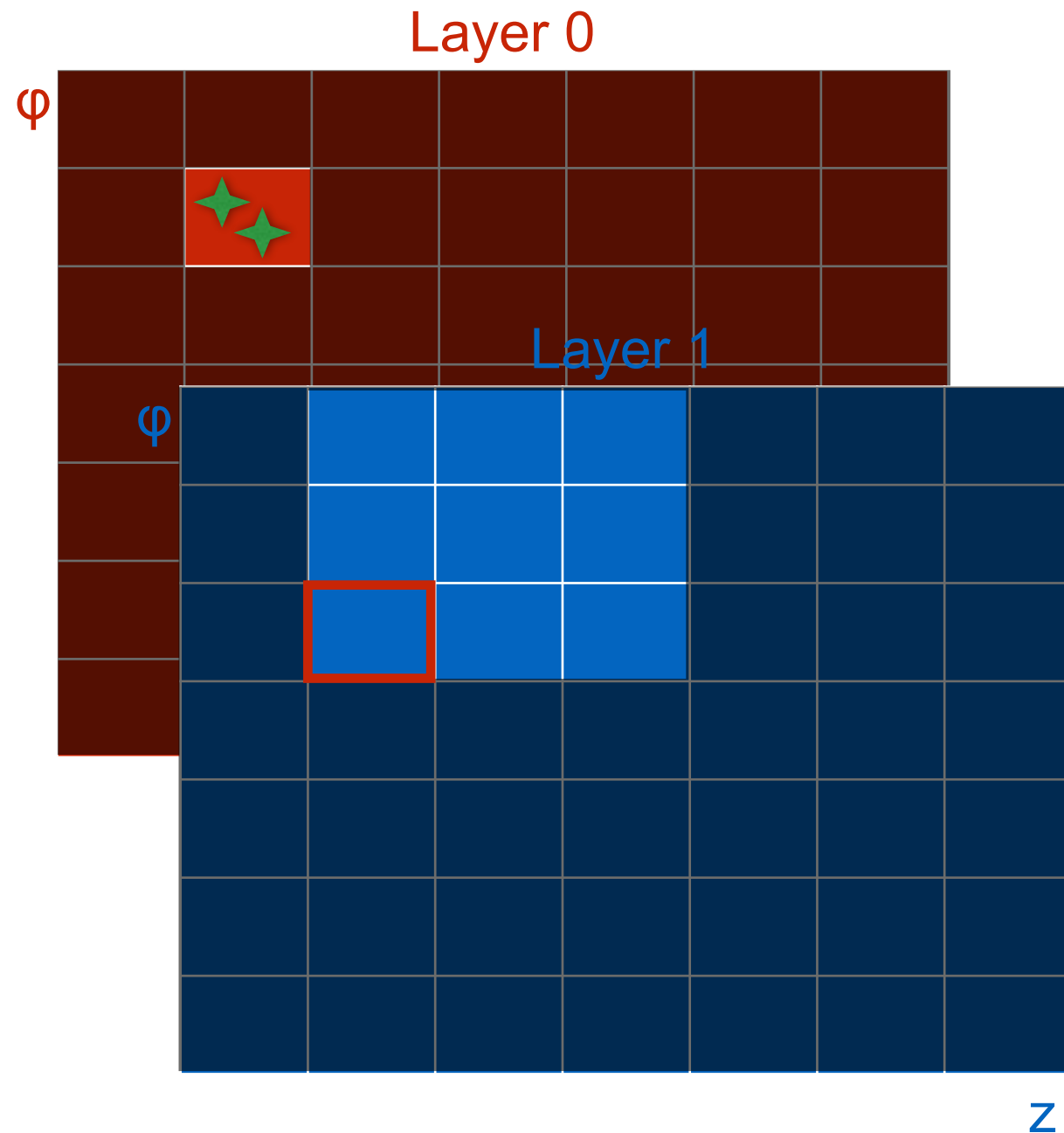
All the clusters of one bin on layer 1 are copied in the shared memory and a for loop starts over these clusters looking for good doublets

Porting on GPU: how to make doublets



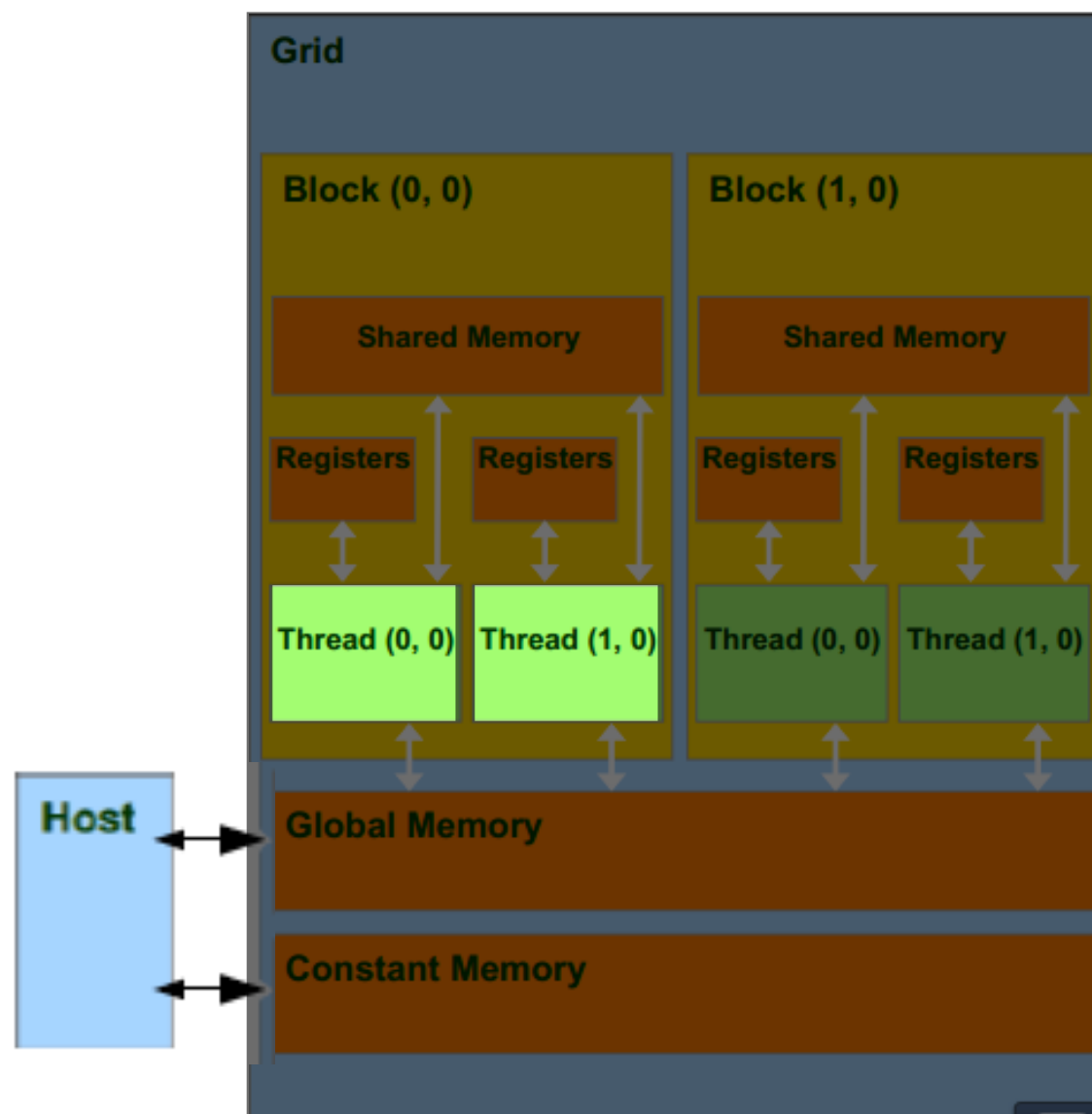
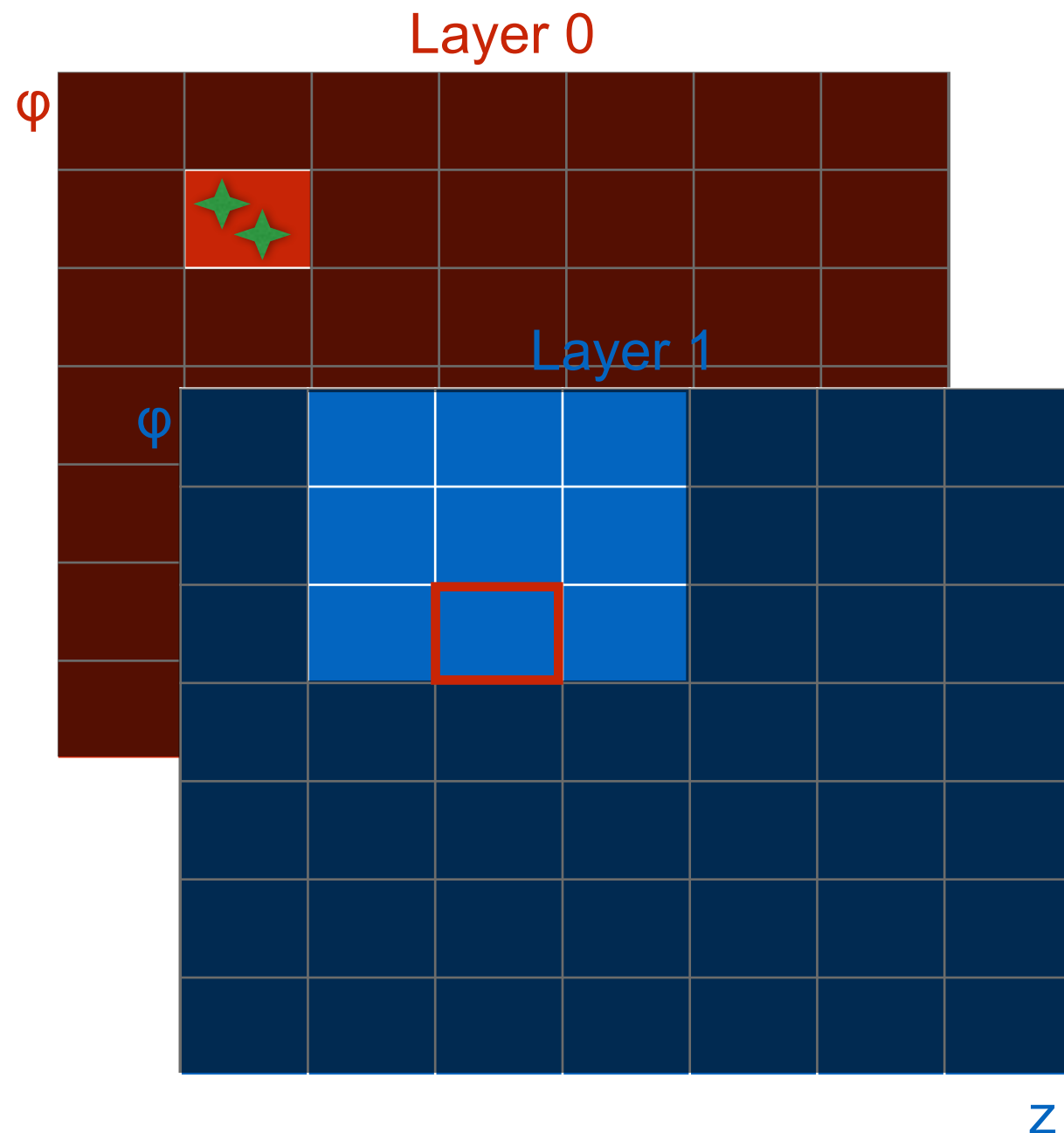
All the clusters of one bin on layer 1 are copied in the shared memory and a for loop starts over these clusters looking for good doublets

Porting on GPU: how to make doublets



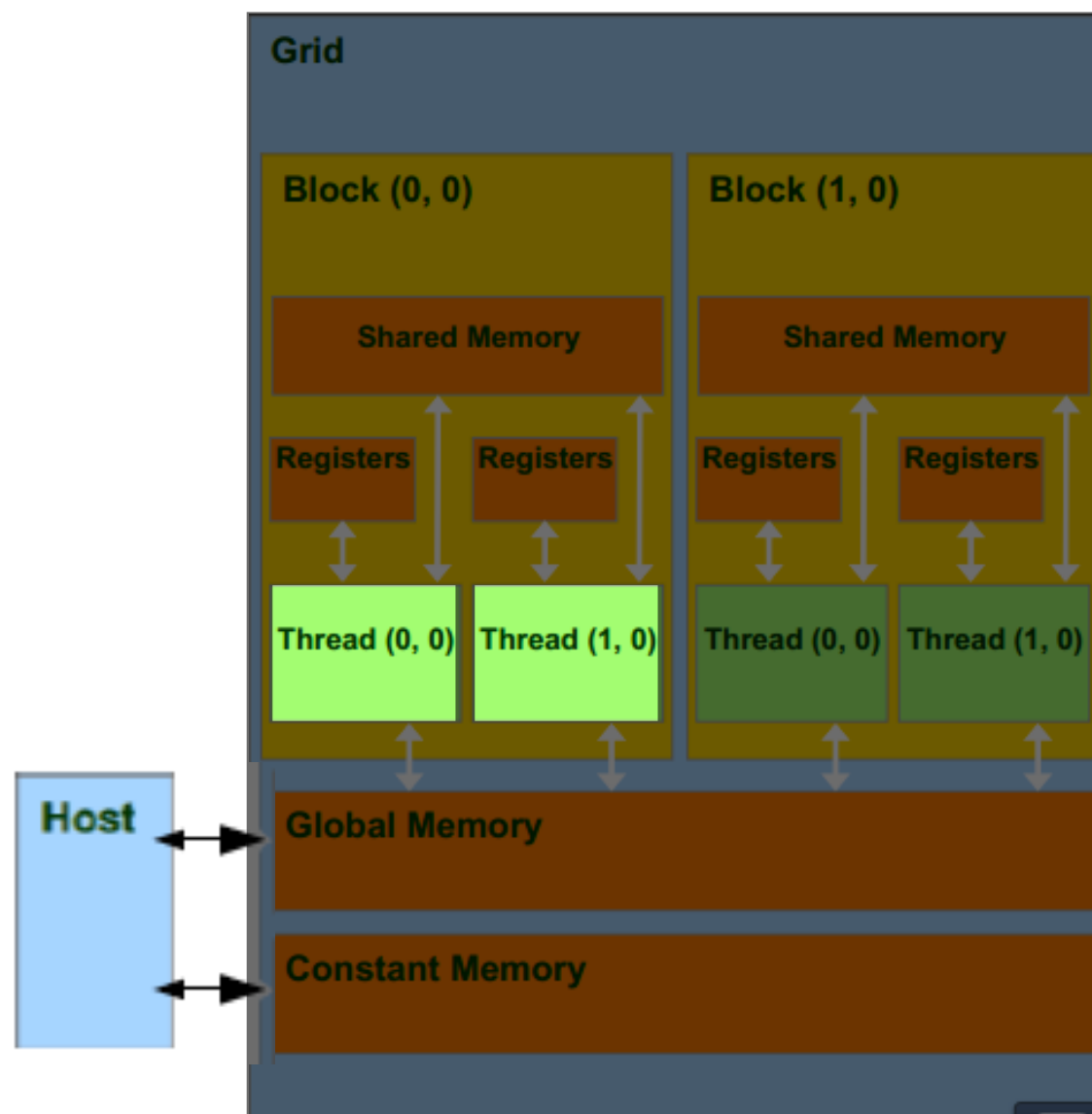
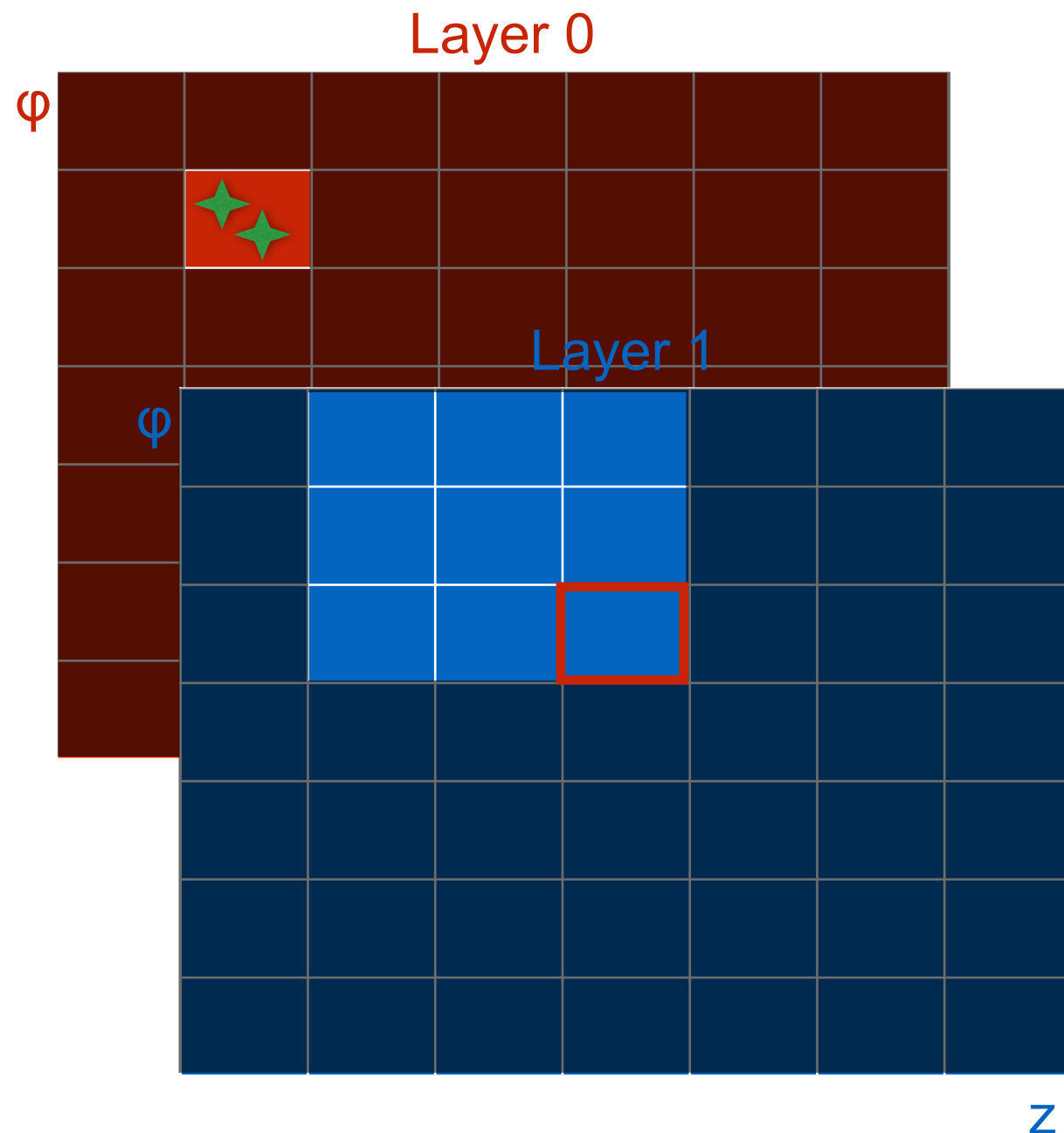
All the clusters of one bin on layer 1 are copied in the shared memory and a for loop starts over these clusters looking for good doublets

Porting on GPU: how to make doublets



All the clusters of one bin on layer 1 are copied in the shared memory and a for loop starts over these clusters looking for good doublets

Porting on GPU: how to make doublets



All the clusters of one bin on layer 1 are copied in the shared memory and a for loop starts over these clusters looking for good doublets