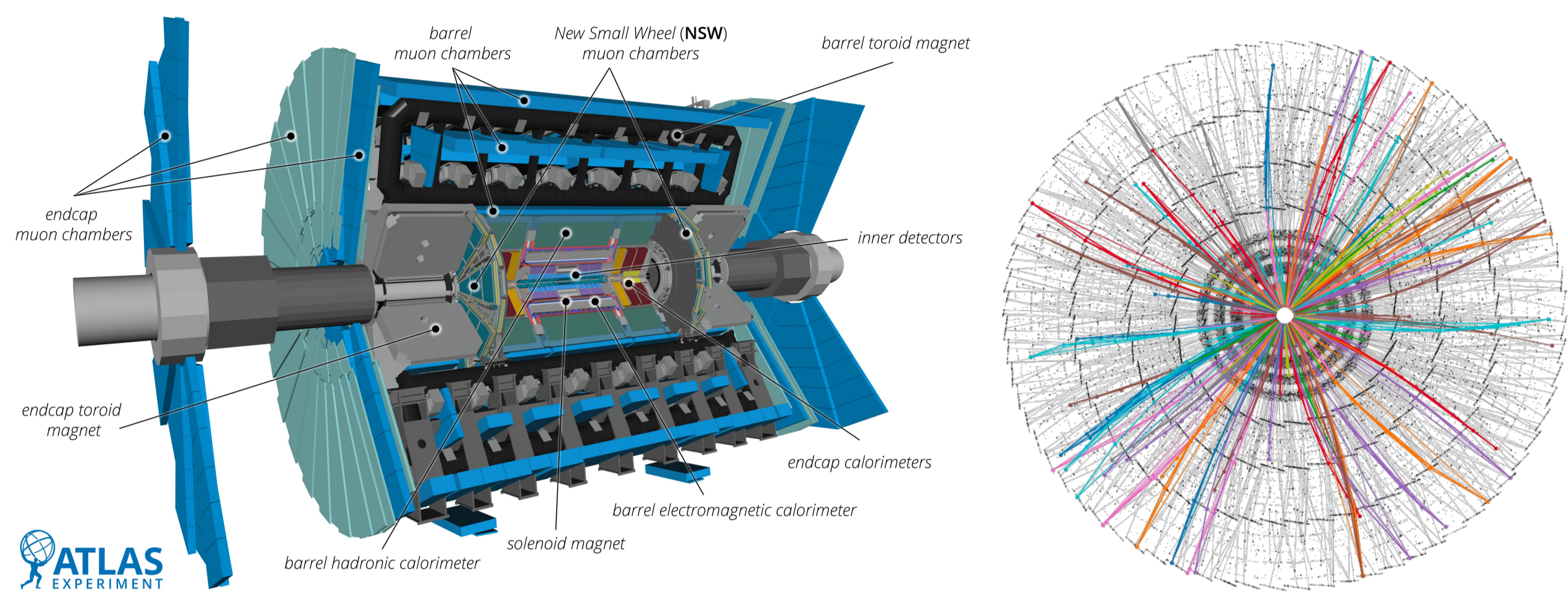


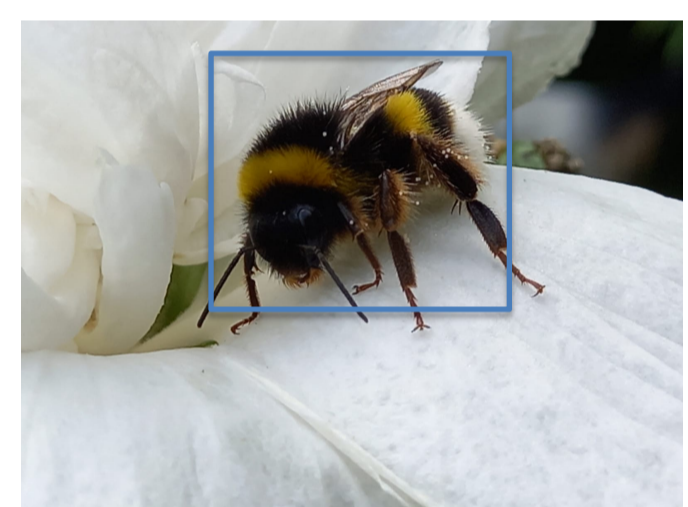
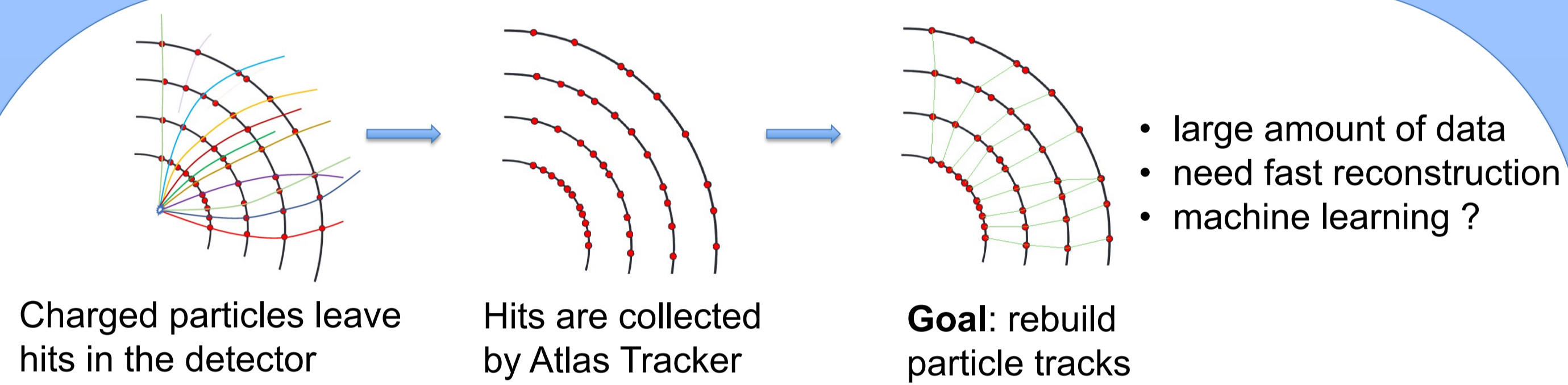
# NEW APPROACHES FOR FAST AND EFFICIENT GRAPH CONSTRUCTION ON CPU / GPU AND HETEROGENEOUS ARCHITECTURES FOR THE ATLAS EVENT RECONSTRUCTION

## ATLAS detector



18108 modules in the ATLAS ITK detector (for HL-LHC 2029)  
average number of space points / event: O(300k)

## Track reconstruction challenge for HL-LHC

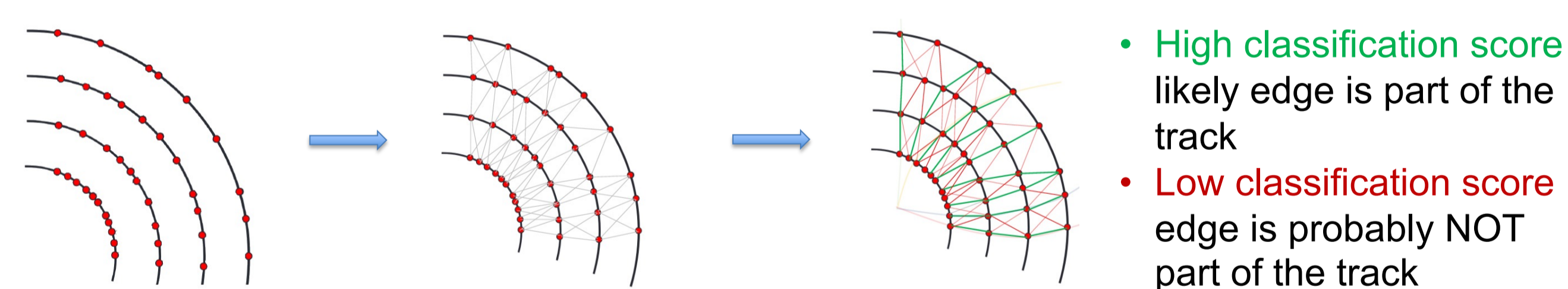


1600 \* 1200 pixels  
a large fraction of the image carries information of the subject

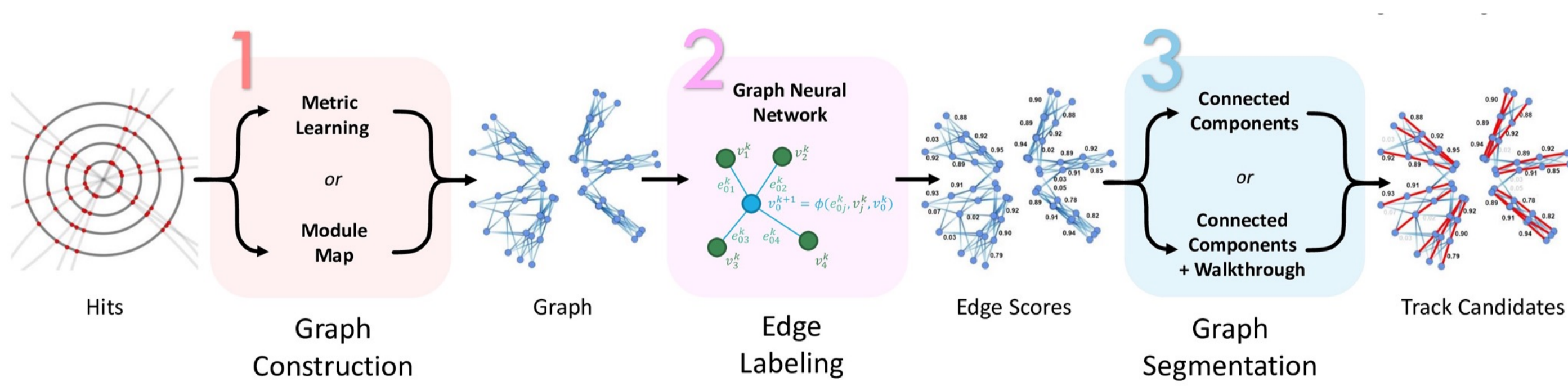
The ATLAS collaboration, "ATLAS Upgrade for the HL-LHC: meeting the challenges of a five-fold increase in collision rate", ATL-UPGRADE-PROC-2012-003, EPJ Web Conf. 28 (2012) 12069

Atlas Tracker for HL-LHC  
5 \* 10<sup>9</sup> readout channels  
~ 3 \* 10<sup>5</sup> 3D space points / event  
**data are sparse**

- Can't use the same tools
- How to represent tracking data with a Neural Network ?



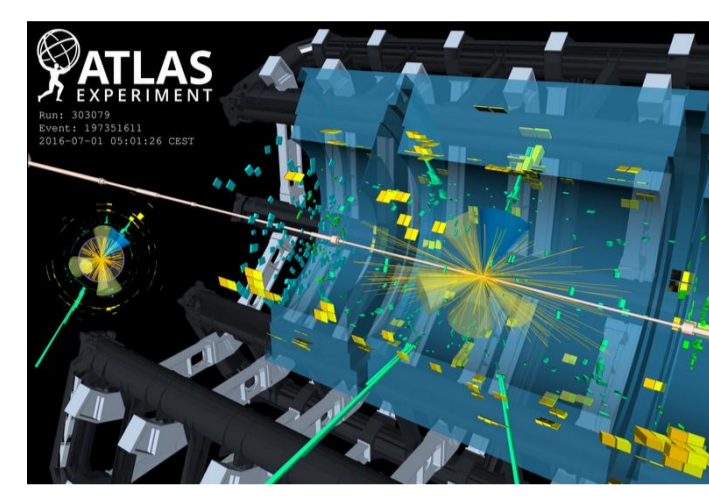
Hits are collected by Atlas Tracker → Represent data using a graph → Classify the edges of the graph



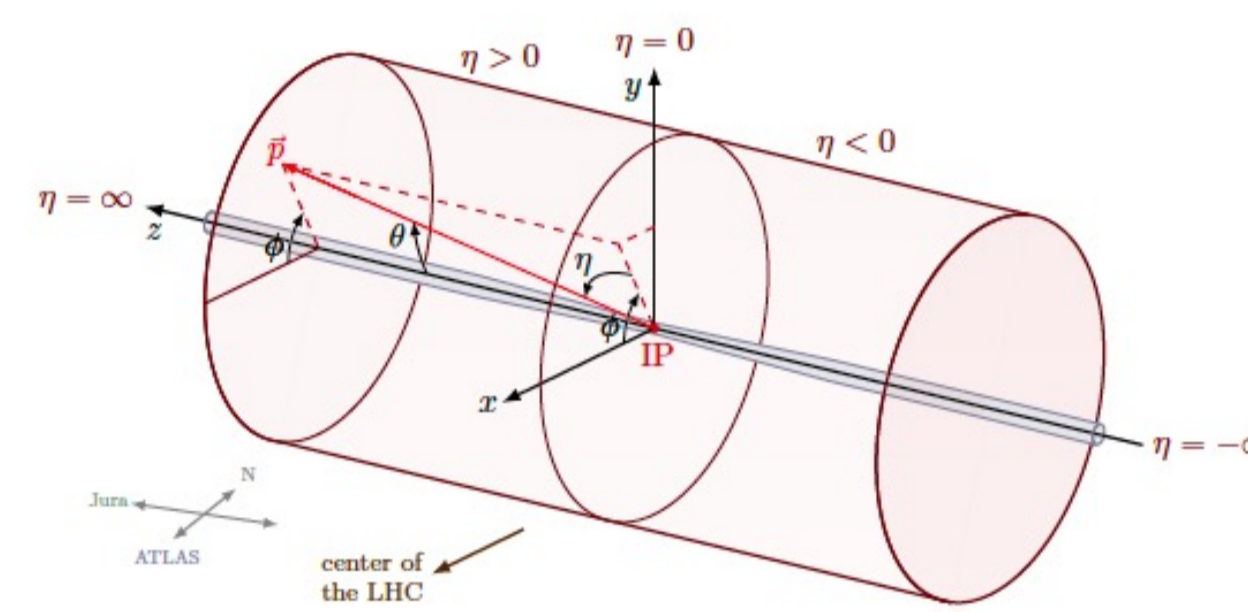
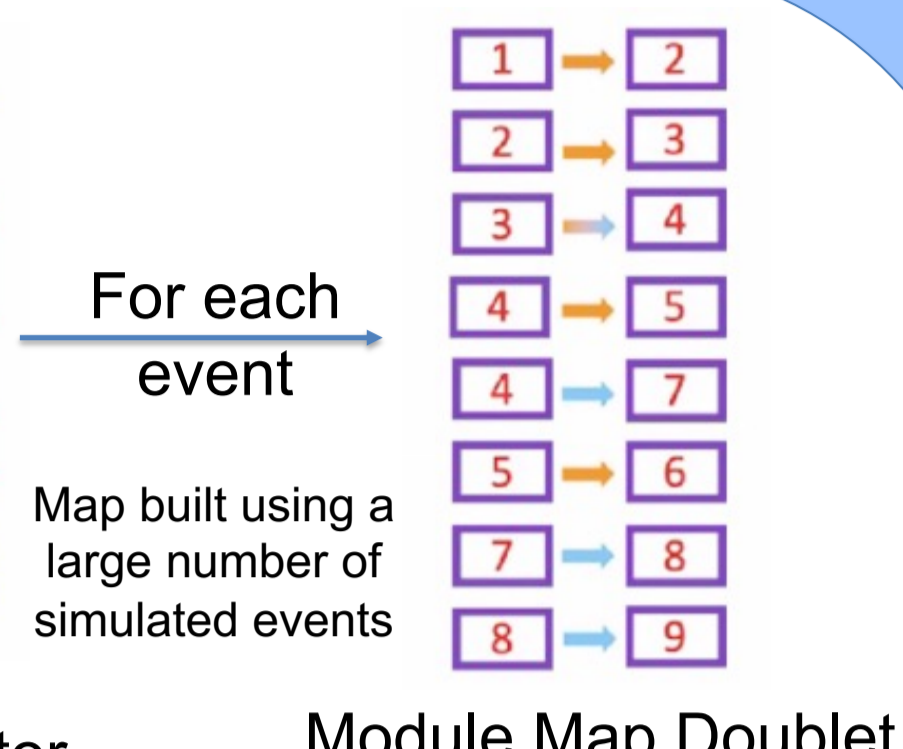
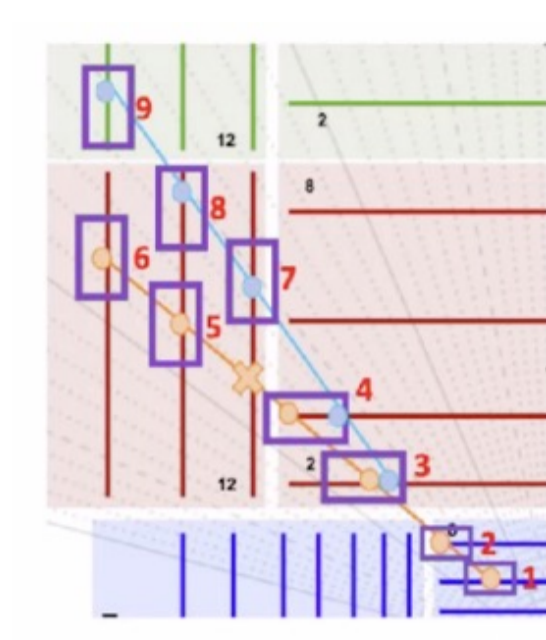
GNN4ITK Full pipeline

The ATLAS Collaboration, "Track finding performance plots for a Graph Neural Network pipeline on ATLAS ITK Simulated Data", IDTR-2022-01, <https://atlas.web.cern.ch/Atlas/GROUPS/PHYSICS/PLOTS/IDTR-2022-01/> (2022)

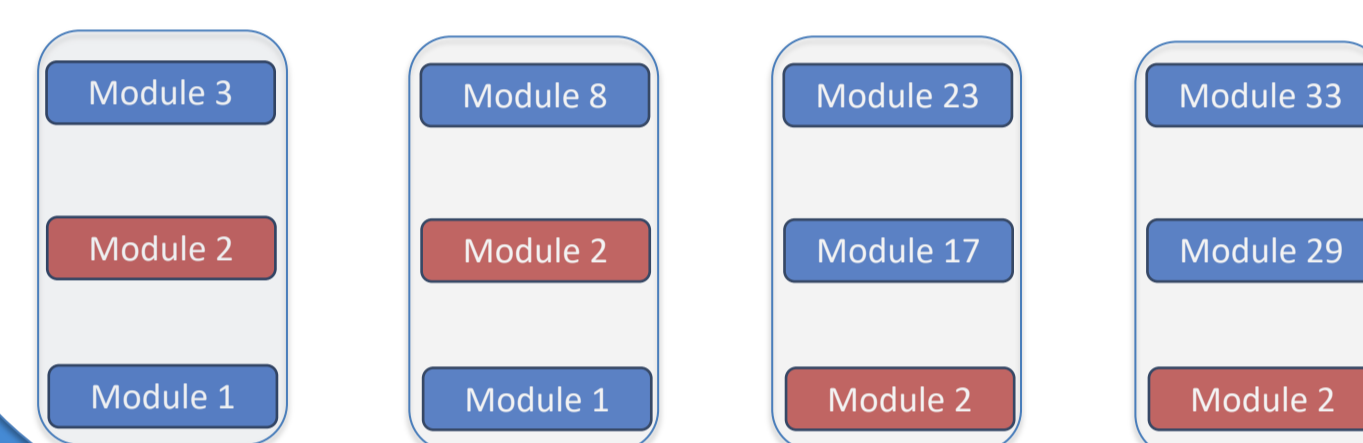
## Module Maps



Module Maps are computed from hits (energy deposits) left in the detector by particles based on several events



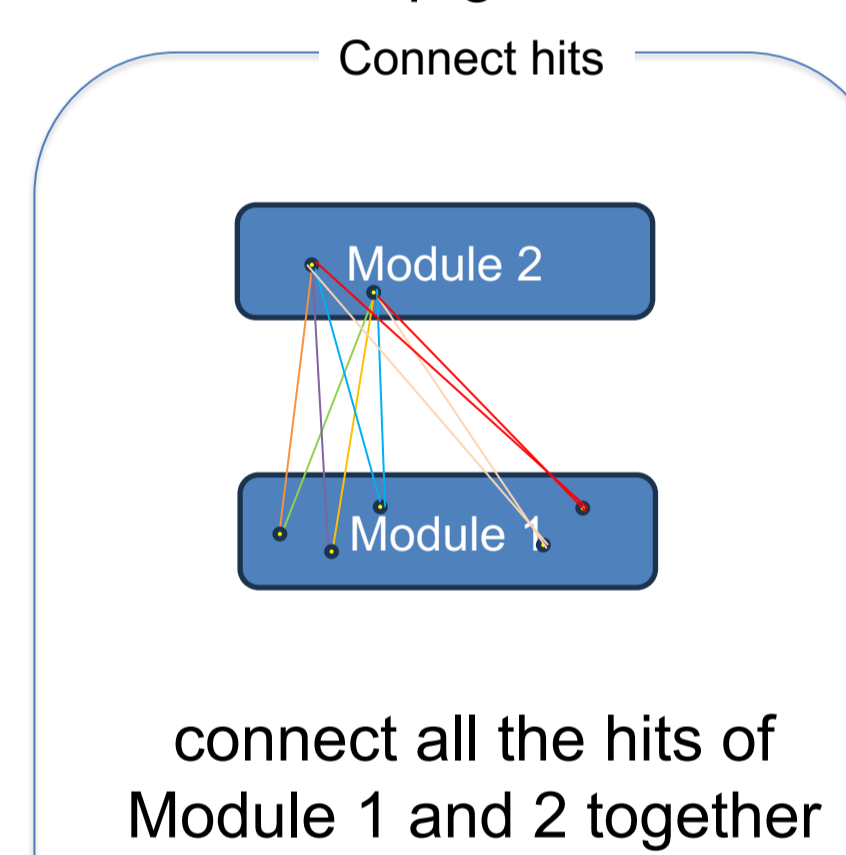
relation hits ↔ particles  
several events used to create a map of energy deposits in the detector  
For each triplet of modules calculate geometrical cuts



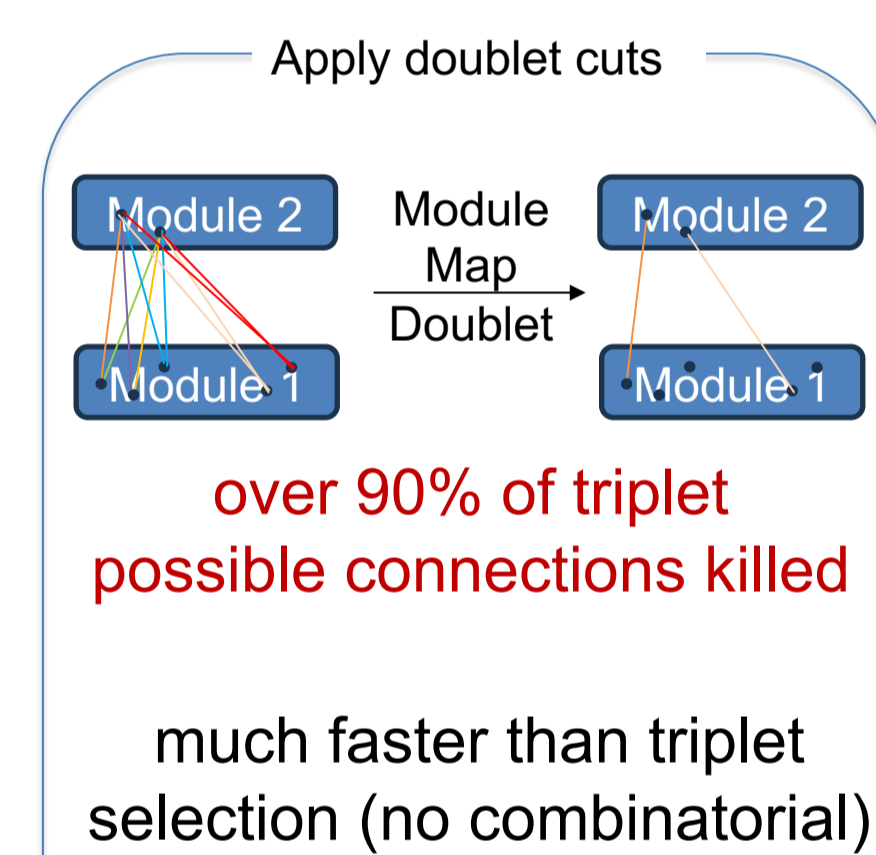
Module Map Doublet  
calculate geometrical cuts on doublets as min and max values of triplet cuts  
→ deduced from Module Map Triplet

## Graph creation

Module Map gives interconnected modules



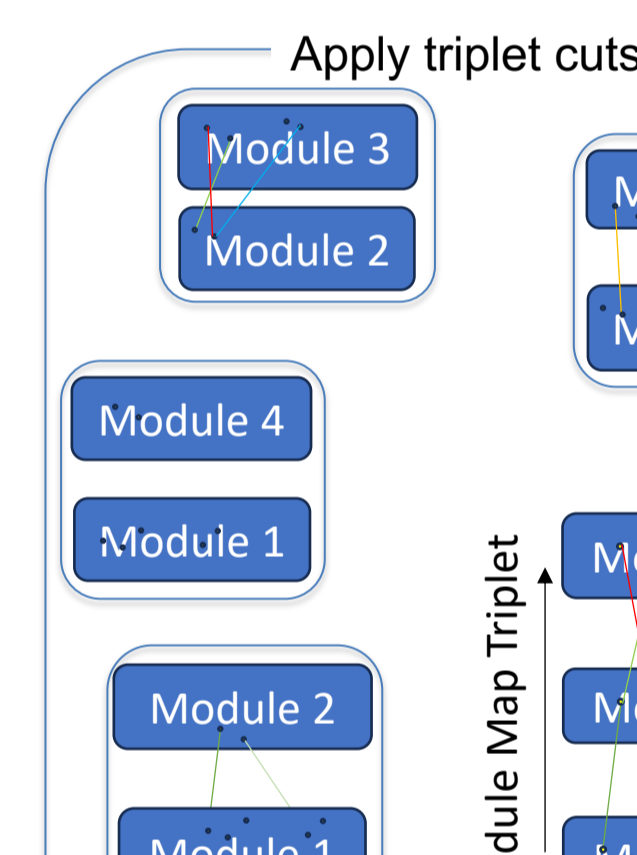
connect all the hits of Module 1 and 2 together



over 90% of triplet possible connections killed

much faster than triplet selection (no combinatorial)

allows to reduce triplet selection time



over 50% of remaining connections killed

## Graph creation : from CPU to GPU

Step 1: Doublet cuts

- send event data to GPU
- apply doublet cuts
- sort edges on output hit (quick sort)

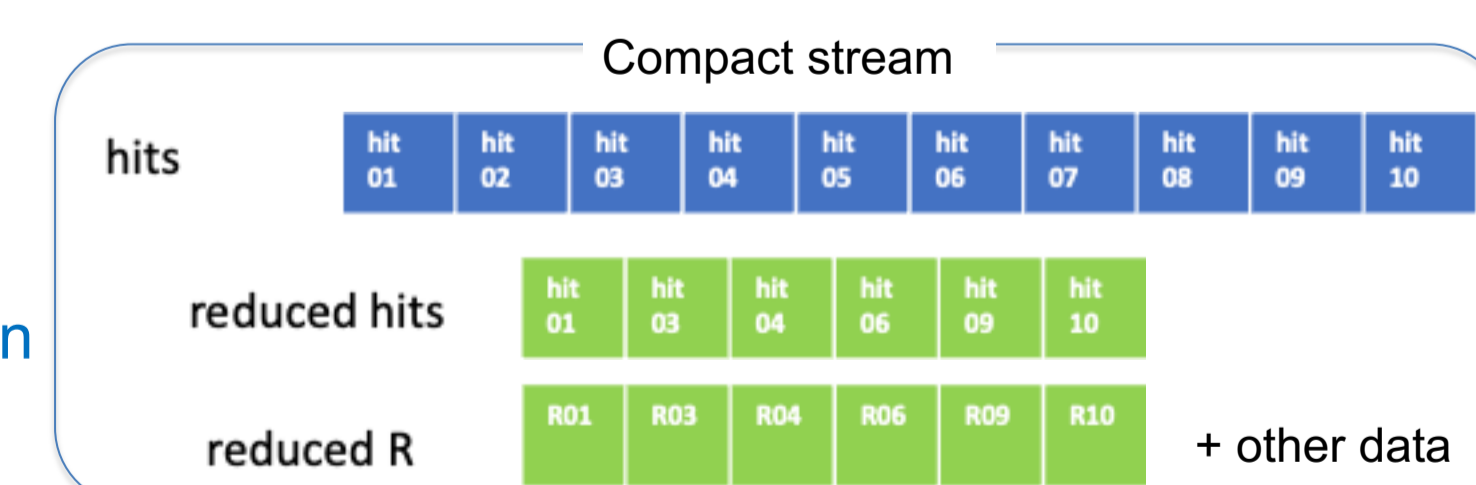
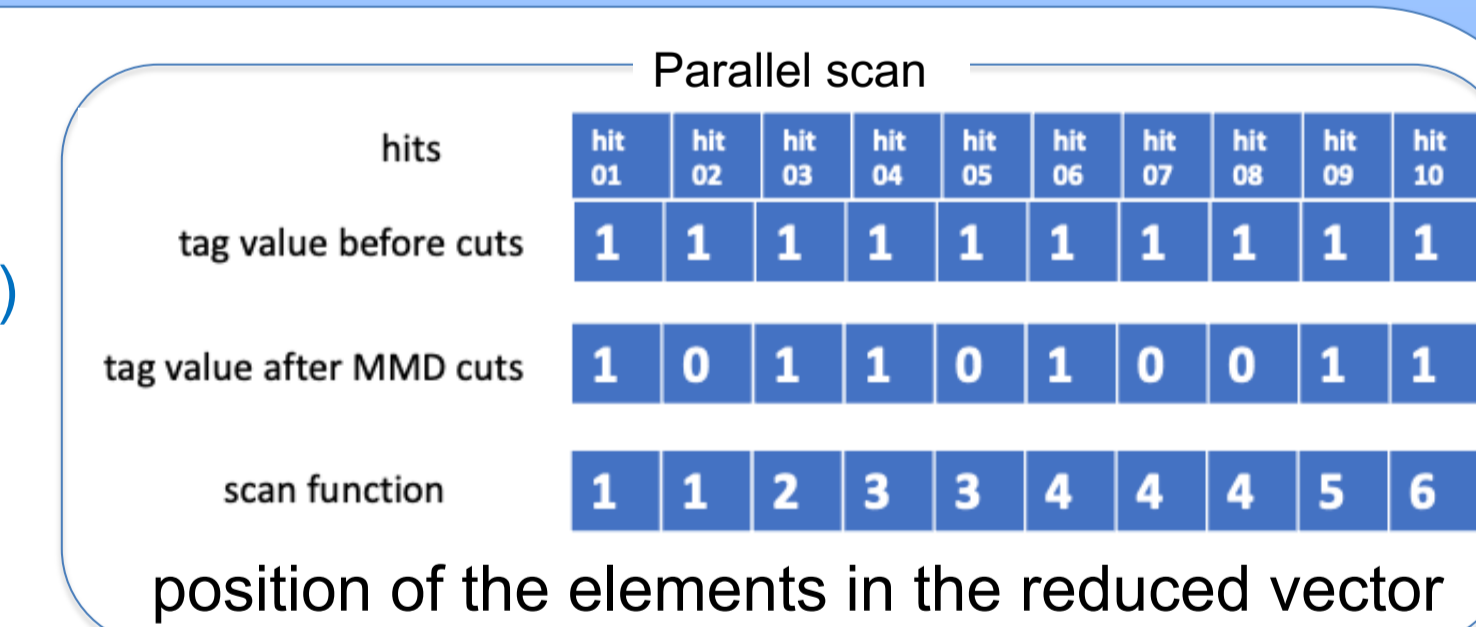
Step 2: Count edges after doublet cuts and reduce data

- parallel scan
- compact stream

Step 3: Triplet cuts

Step 4: Data reduction

- active edges only
- parallel scan + stream compaction
- copy reduced data to cpu
- free cuda memory



## GPU profiling and computational time

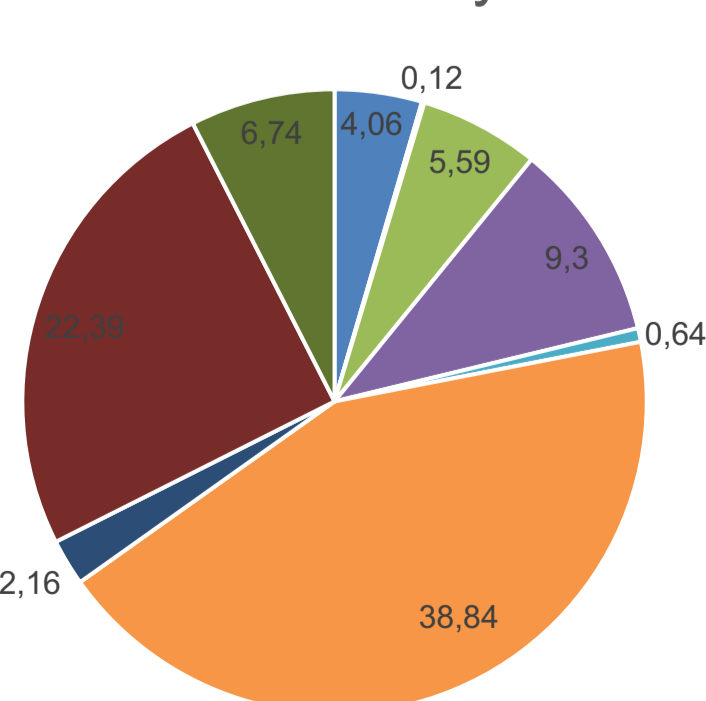
TTree\_hits\_constants  
Section: Occupancy

Metric Name	Metric Unit	Metric Value
Block Limit SM	block	16
Block Limit Registers	block	3
Block Limit Shared Mem	block	16
Block Limit Warps	block	2
Theoretical Active Warps per SM	warp	32
Theoretical Occupancy	%	100
Achieved Occupancy	%	91.23
Achieved Active Warps Per SM	warp	29.19

TTree\_hits\_constants

Metric Name	Metric Unit	Metric Value
DRAM Frequency	cycle/nsecond	6.02
SM Frequency	cycle/nsecond	1.15
Elapsed Cycles	cycle	783,443
Memory Throughput	%	1.15
DRAM Throughput	%	1.15
Duration	usecond	609.38
L1/TEX Cache Throughput	%	0.56
L2 Cache Throughput	%	0.71
SM Active Cycles	cycle	656,862.85
Compute (SM) Throughput	%	80.70

% GPU time chart of the 69ms by function



cp to device constants GPU alloc doublet cuts scan/compact  
triplet cuts reduction cp to host free memory

CUDA kernel for graph construction

Currently being integrated in ACTS

Can be integrated to Athena

Comparison CPU vs GPU

differences on 100 events data  
nb nodes 22 / 27,384,853  
nb of edges 446 / 177,067,905

Memory resources

10 to 11 Go / event

Average GPU time

69ms / event

based on 100 random events of CTD 2023 (simulated pile up of 200 tbar events) dataset running on Nvidia A100