

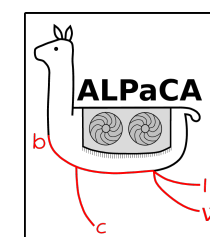
The Real Time Analysis (RTA) Upgrade II project at LHCb

Dorothea vom Bruch

CPPM, Aix-Marseille Université, Marseille, CNRS/IN2P3

Conseil Scientifique de l'IN2P3

February 13th 2025

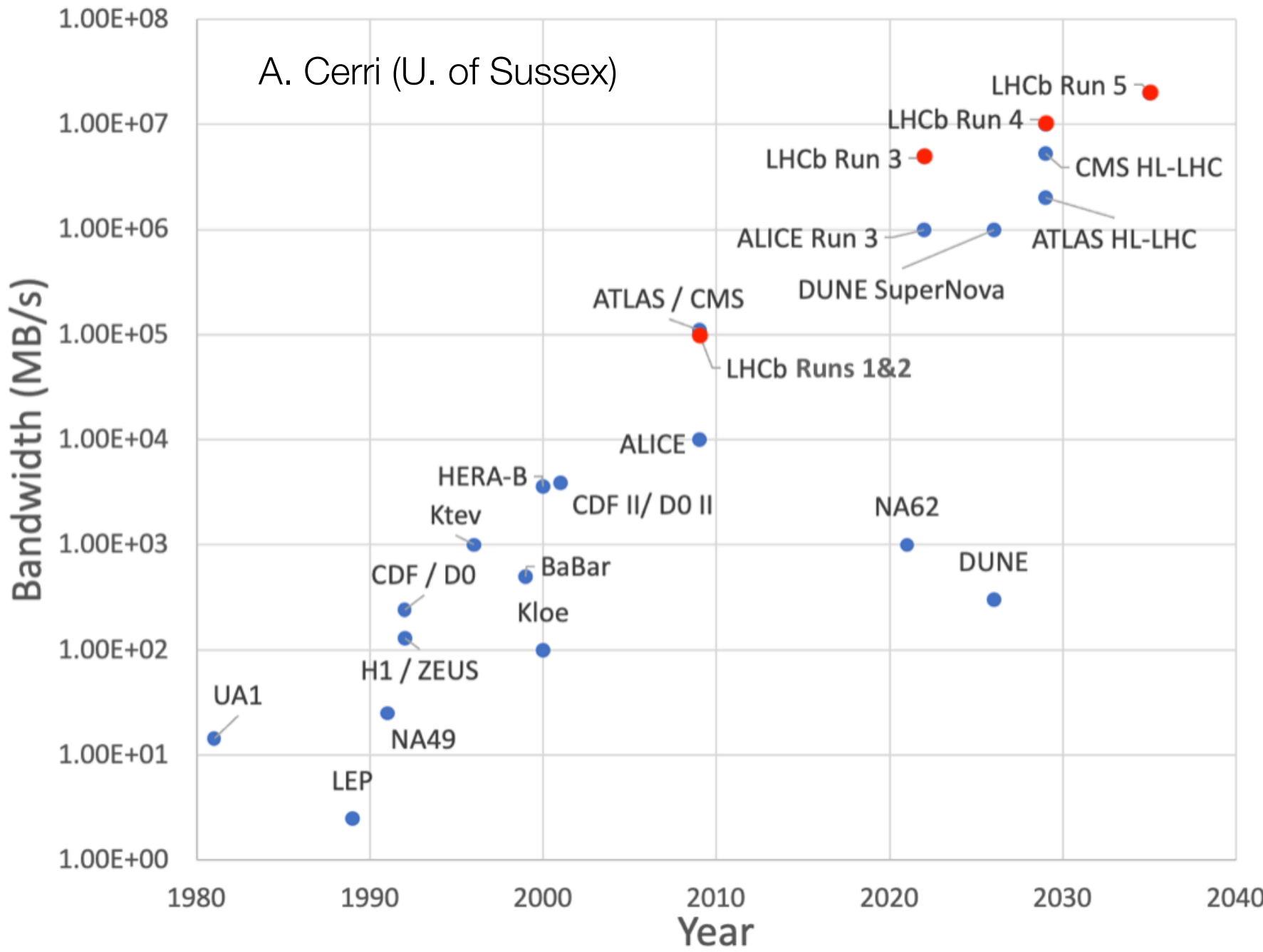
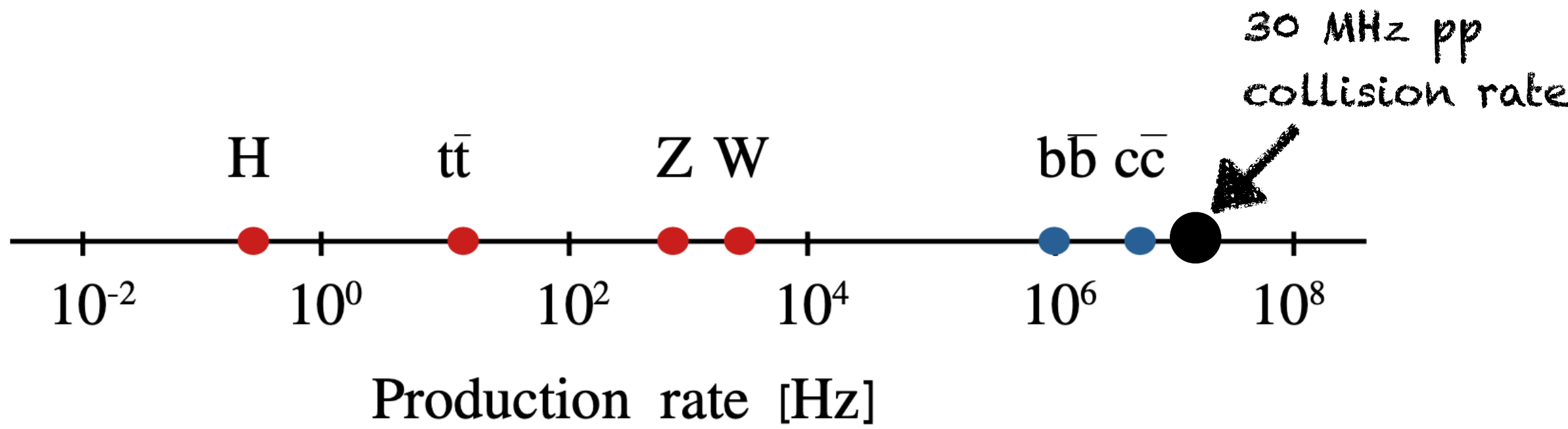


The LHCb trigger challenge

$\mathcal{L} = 2 \times 10^{34} \text{ cm}^{-2}\text{s}^{-1}$ (ATLAS/CMS)

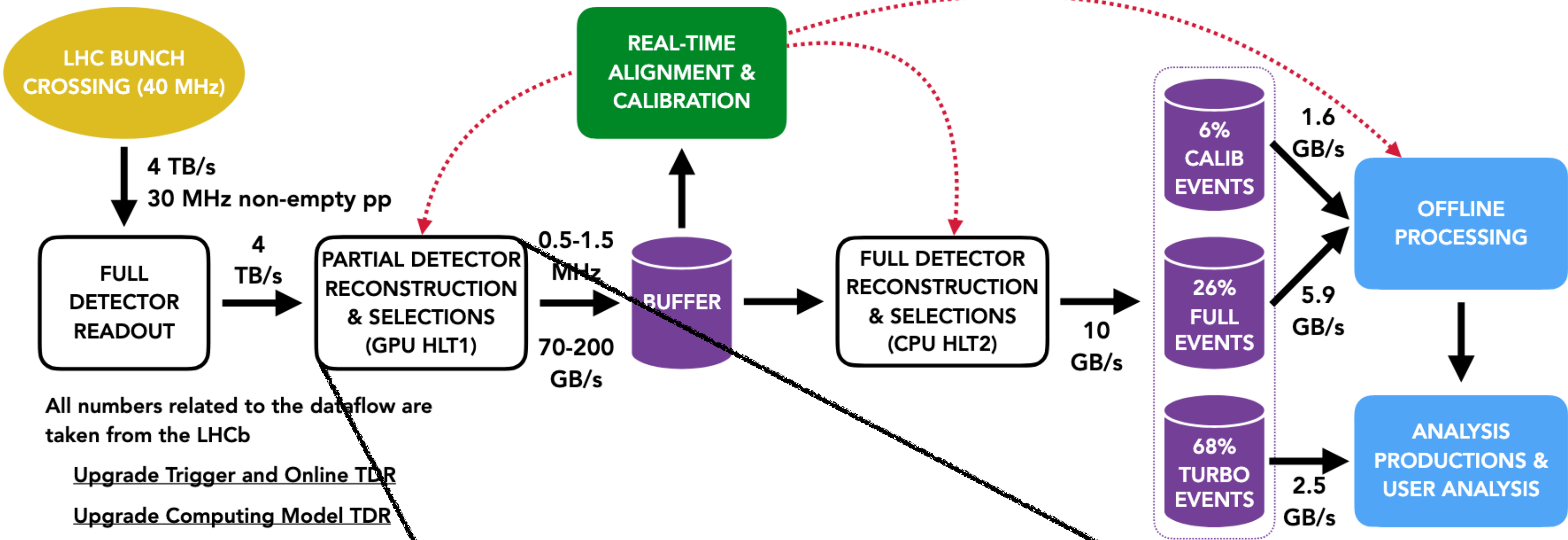
$\mathcal{L} = 2 \times 10^{33} \text{ cm}^{-2}\text{s}^{-1}$ (LHCb)

$\sqrt{s} = 13.6 \text{ TeV}$



- Key signature in LHCb is a secondary vertex with significant transverse momentum and displacement from the pp collision
- Charged particle reconstruction at 30 MHz in the full detector is necessary

The LHCb trigger in LHC Runs 3 & 4



All numbers related to the dataflow are taken from the LHCb

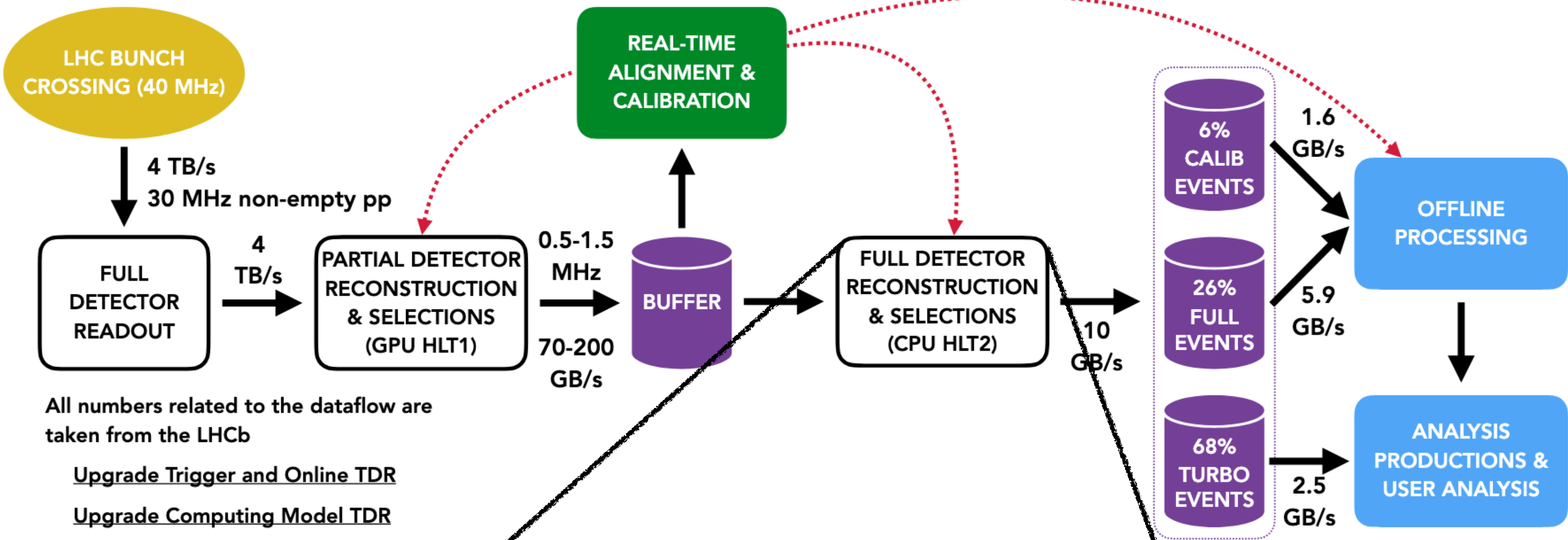
Upgrade Trigger and Online TDR

Upgrade Computing Model TDR

High Level Trigger 1 (HLT1)

- 30 MHz input
- Not latency bound
- O(100) algorithms to maintain
- O(10) developers
- 125k lines of code
- Processed by Allen software on O(500) GPUs

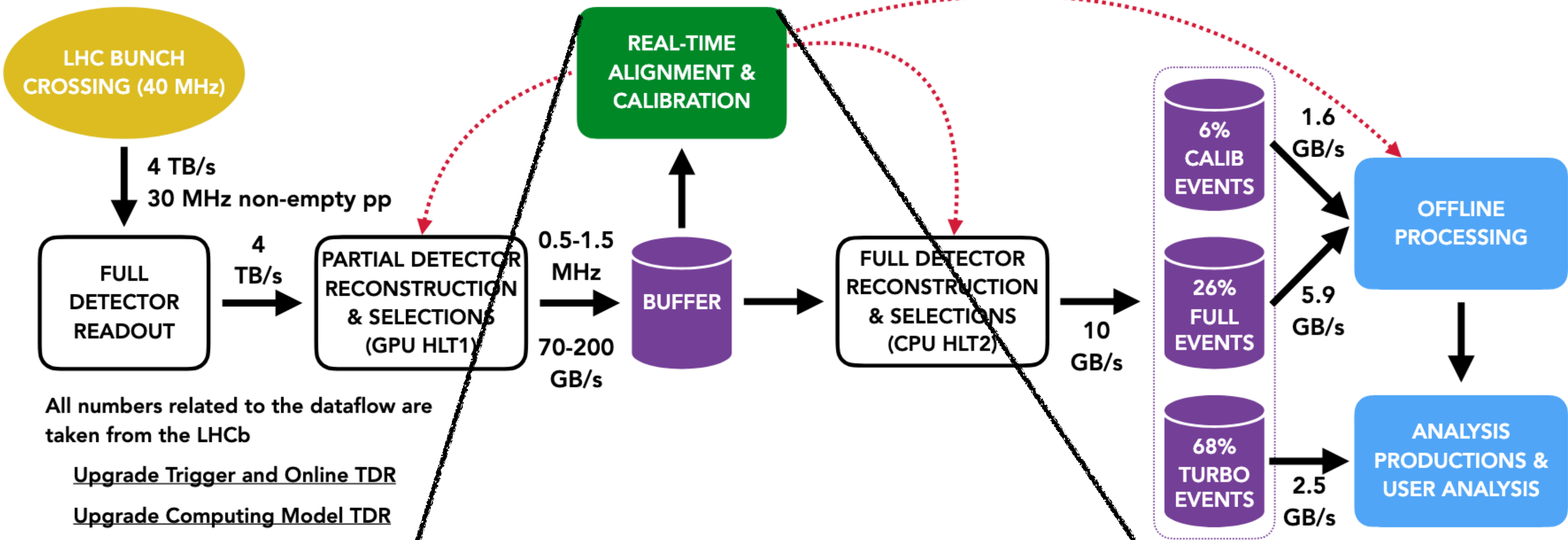
The LHCb trigger in LHC Runs 3 & 4



High Level Trigger 2 (HLT2)

- 1 MHz input
- Not latency bound
- O(2000) algorithms to maintain
- O(100) developers
- 2M lines of code
- Processed by Moore software on CPUs

The LHCb trigger in LHC Runs 3 & 4



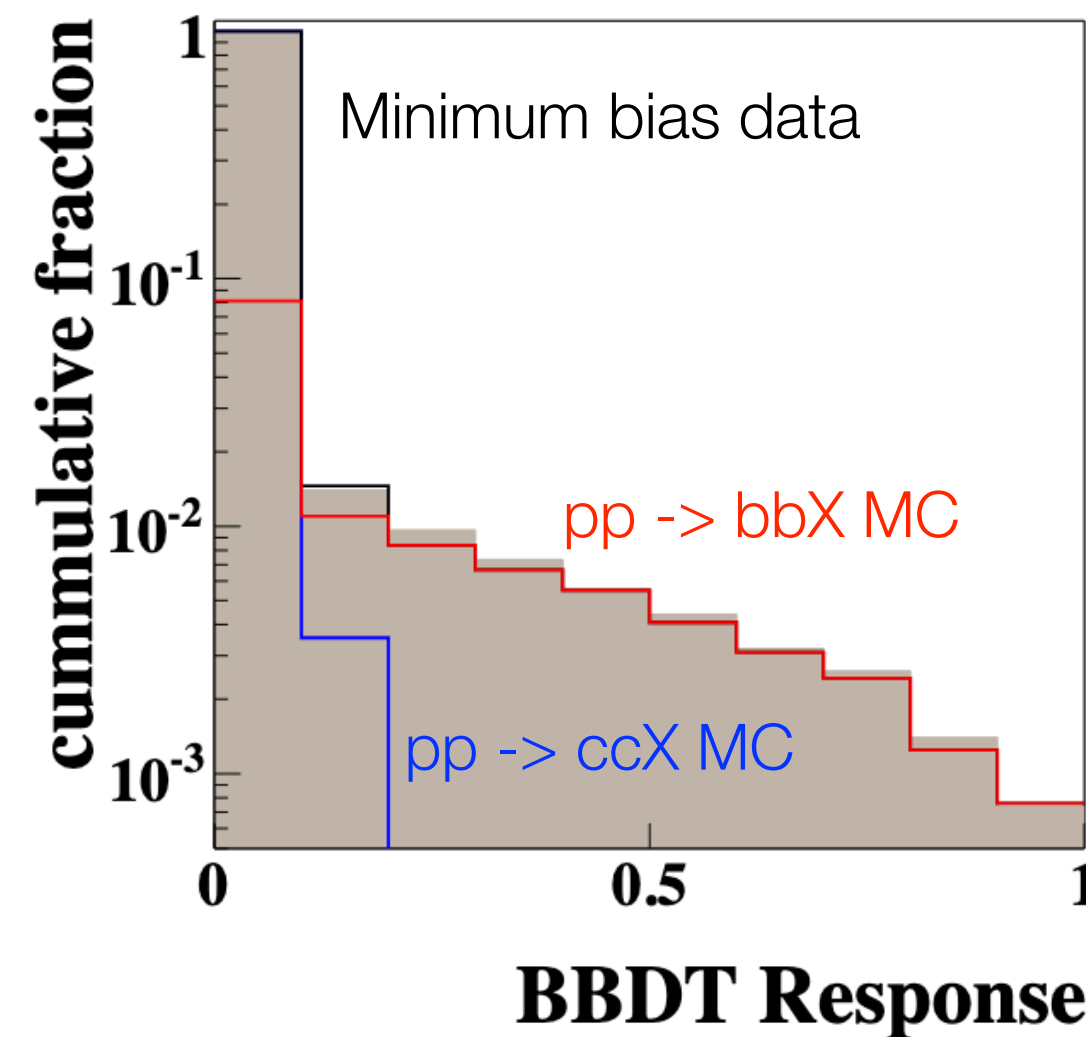
Alignment & Calibration

- Determine detector position based on reconstructed quantities
- RICH mirror alignment & refractive index calibration
- Calibrate calorimeter

Machine learning and Artificial Intelligence in LHCb RTA

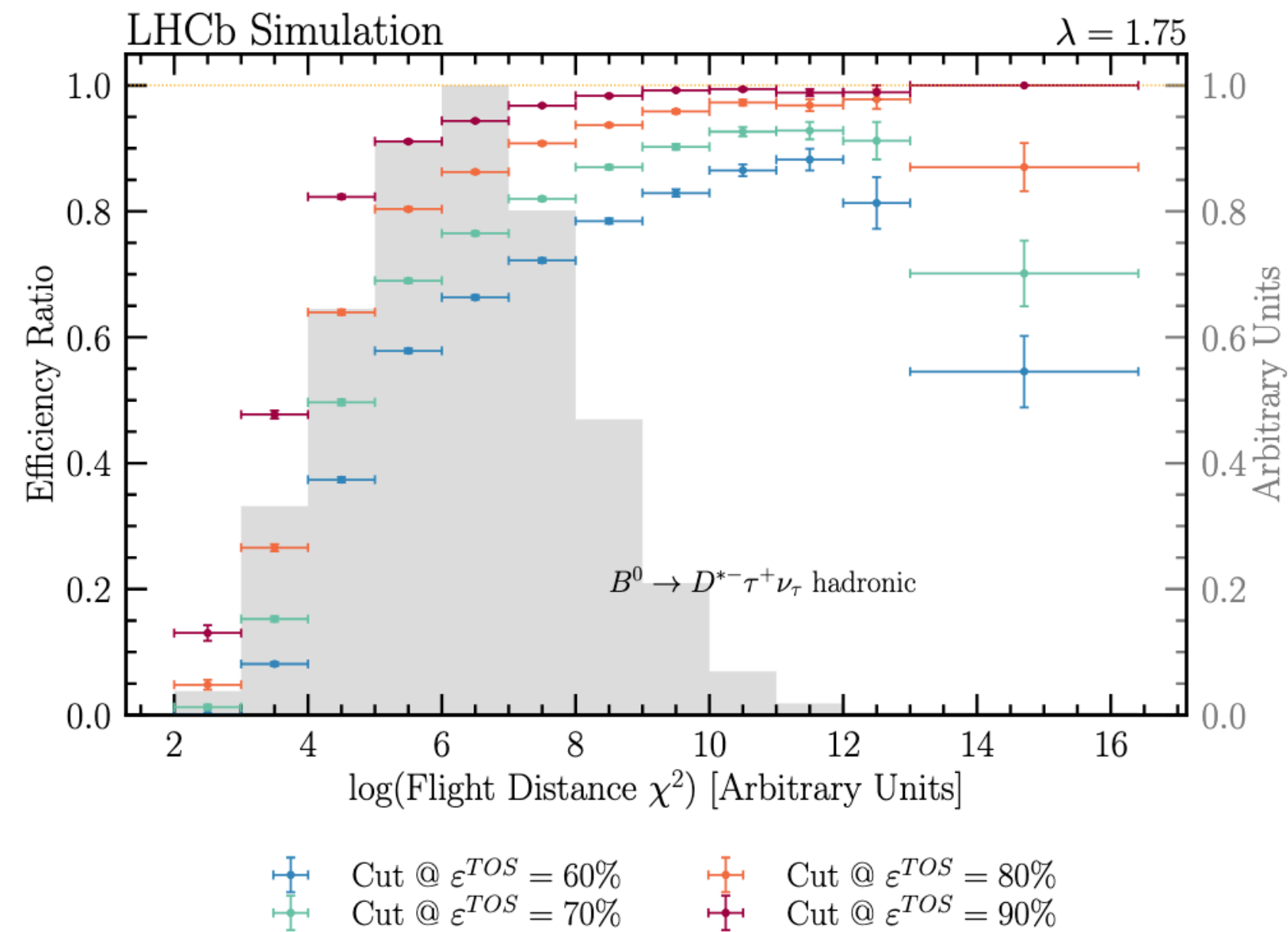
LHCb has pioneered AI/ML techniques, using resource-aware models since the start of Run 1

Bonsai Boosted Decision Tree (BBDT) in Run 1



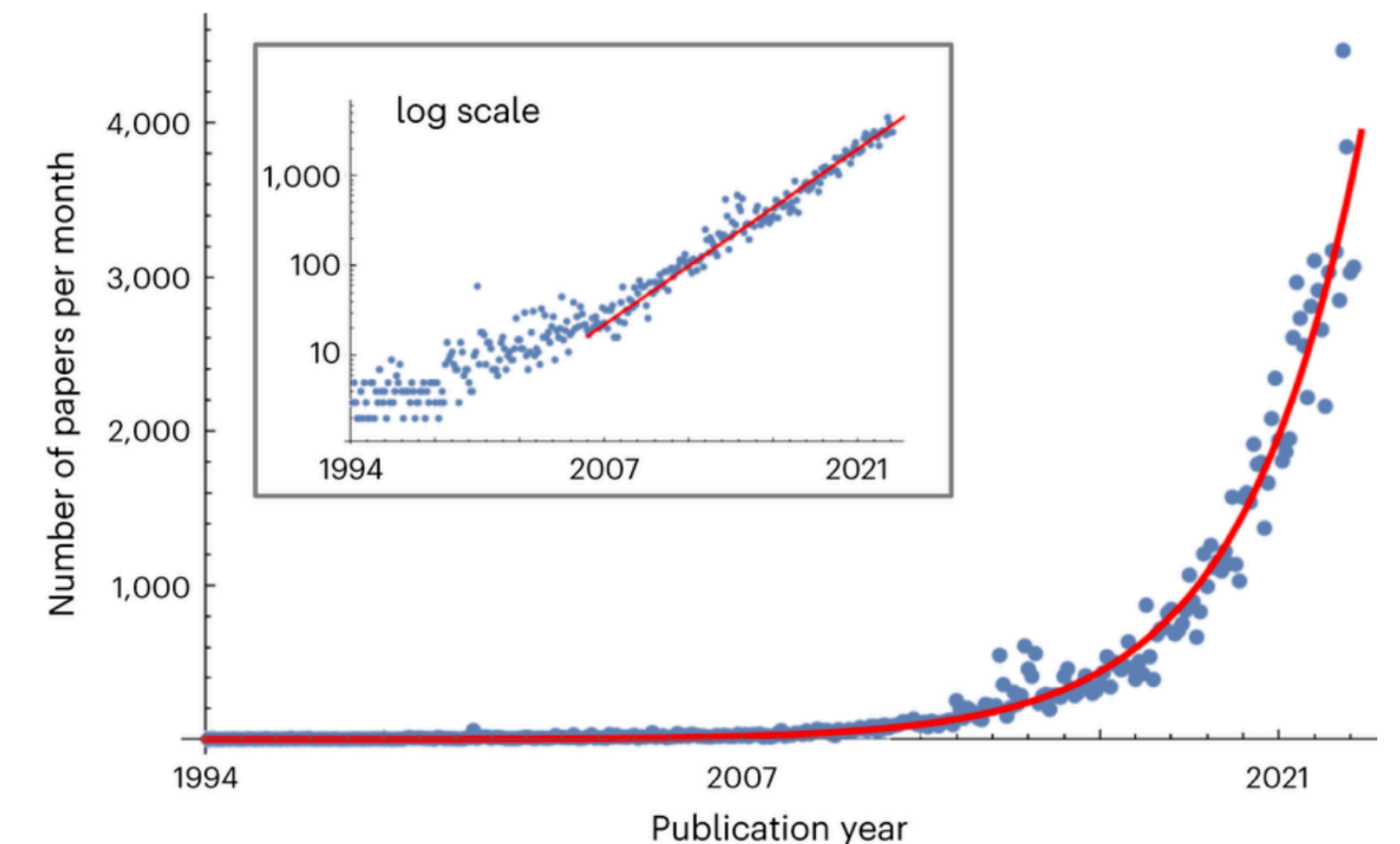
[LHCb-PUB-2011-016](#)

Lipschitz Neural Networks in Run 3



[arXiv:2312.14265](#)

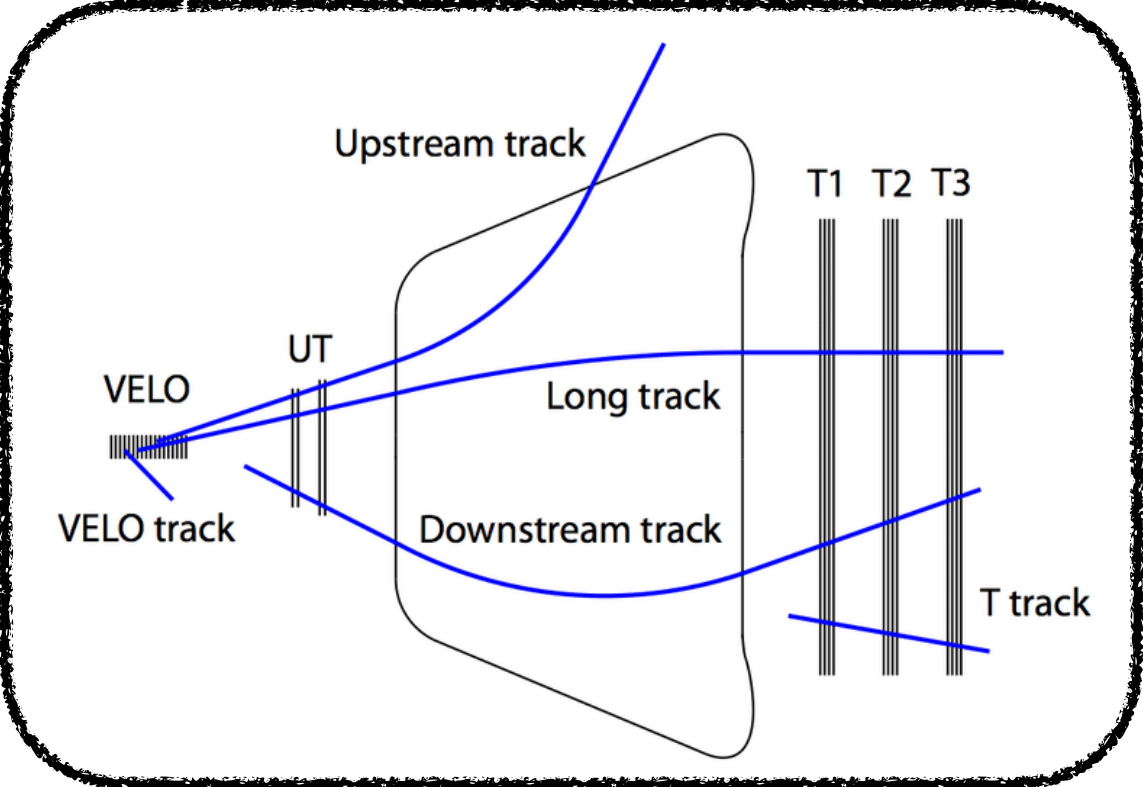
Explosion of ML /AI methods in industry In recent years



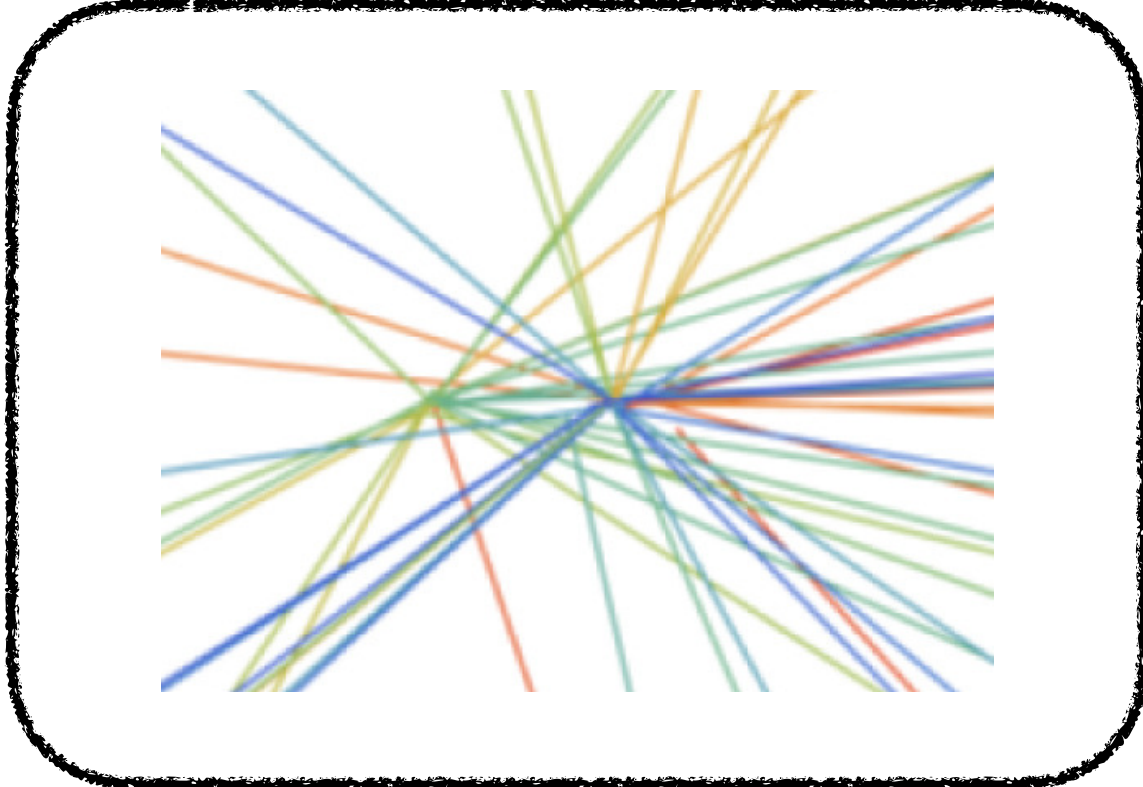
Nature Machine Intelligence, Volume 5, November 2023, 1326-1335

What do we reconstruct in the High Level Trigger (HLT)?

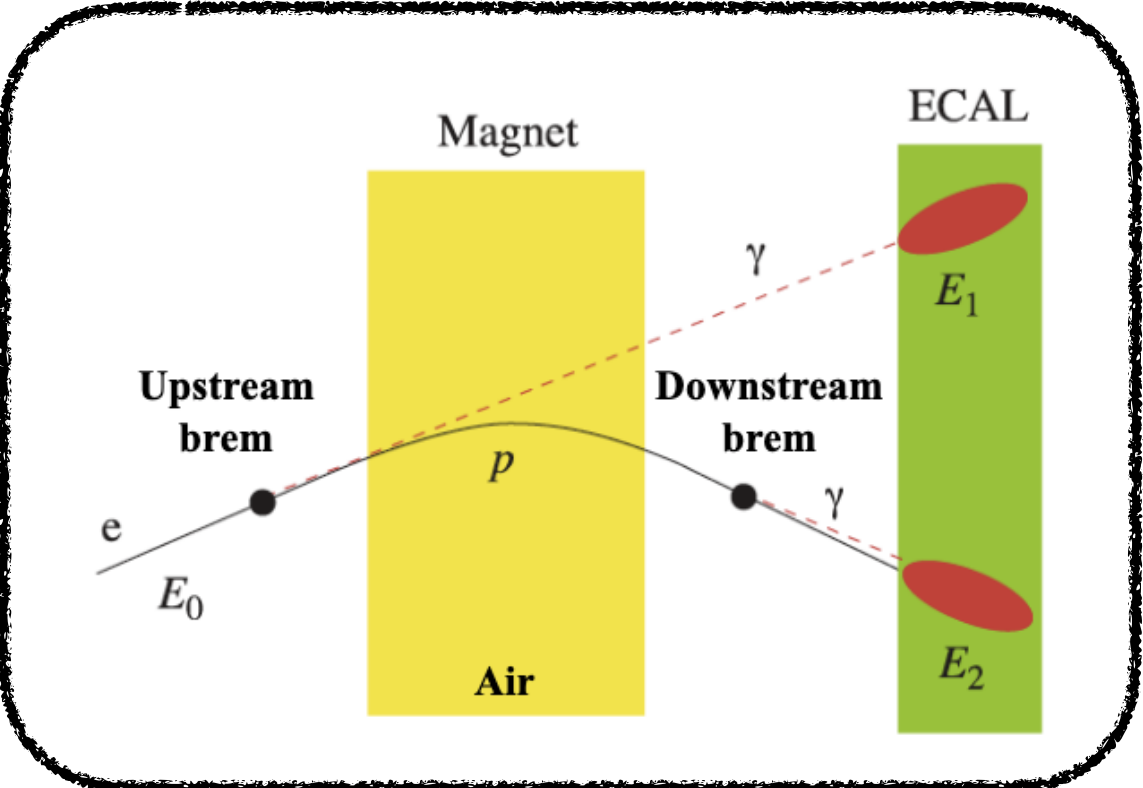
Tracks



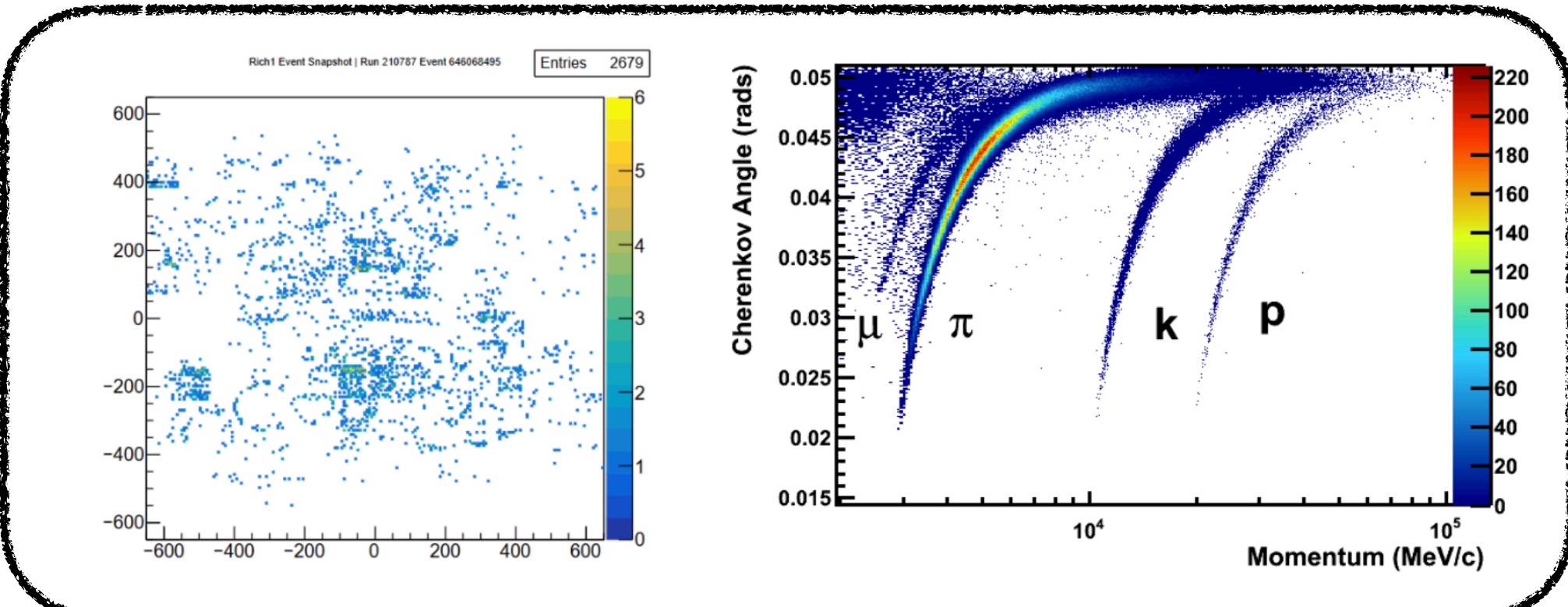
Vertices



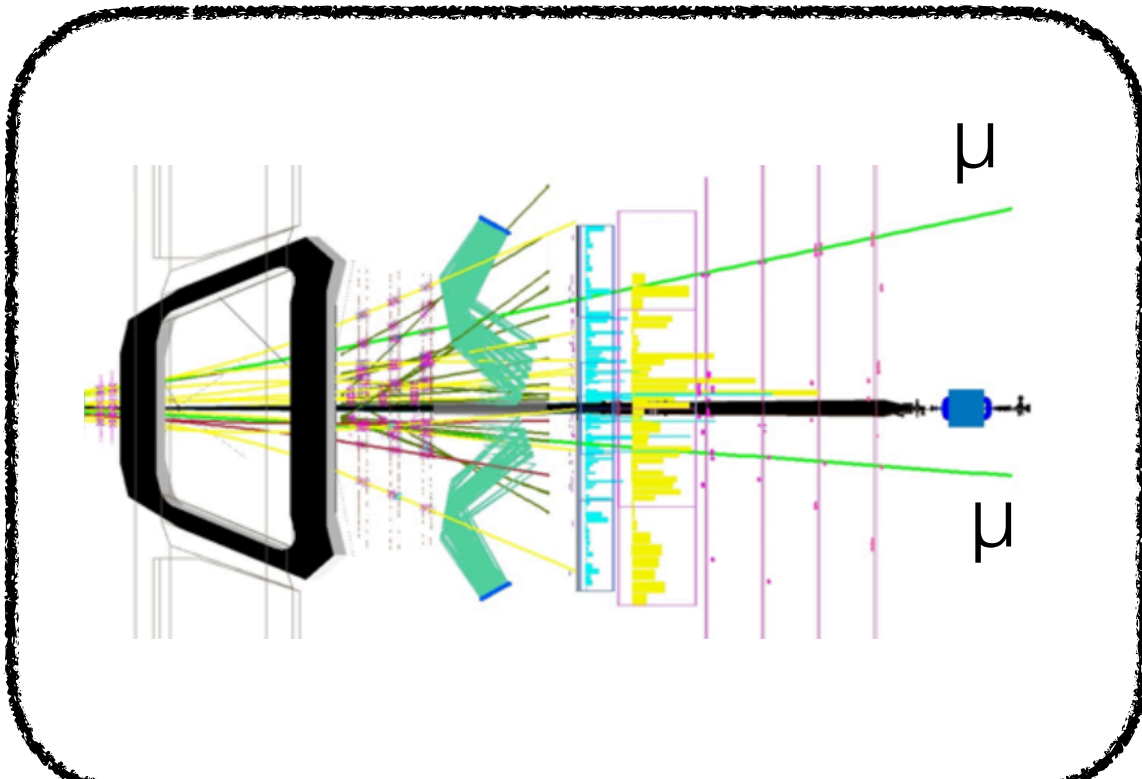
Calorimeter objects



Cherenkov rings



Muon hits



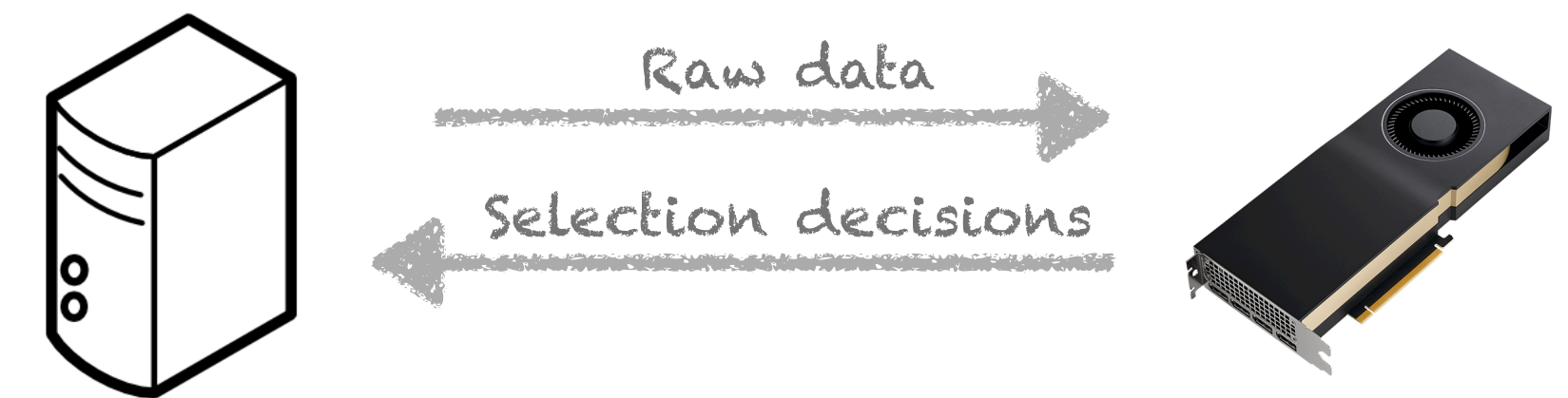
How does the HLT map to GPUs?

Characteristics of LHCb's HLT	Characteristics of GPUs
Intrinsically parallel problem: Process collisions and objects in parallel	Suited for: <ul style="list-style-type: none">• Data-intensive parallelizable applications• High throughput applications
Huge compute load	Many TFLOPs available
Full data stream is read out → No stringent latency requirements	Higher latency than CPUs Not as predictable as FPGAs
Small raw event data (~100 kB)	<ul style="list-style-type: none">• Connection via PCIe → limited I/O bandwidth• Thousands of events fit into O(10) GB of memory

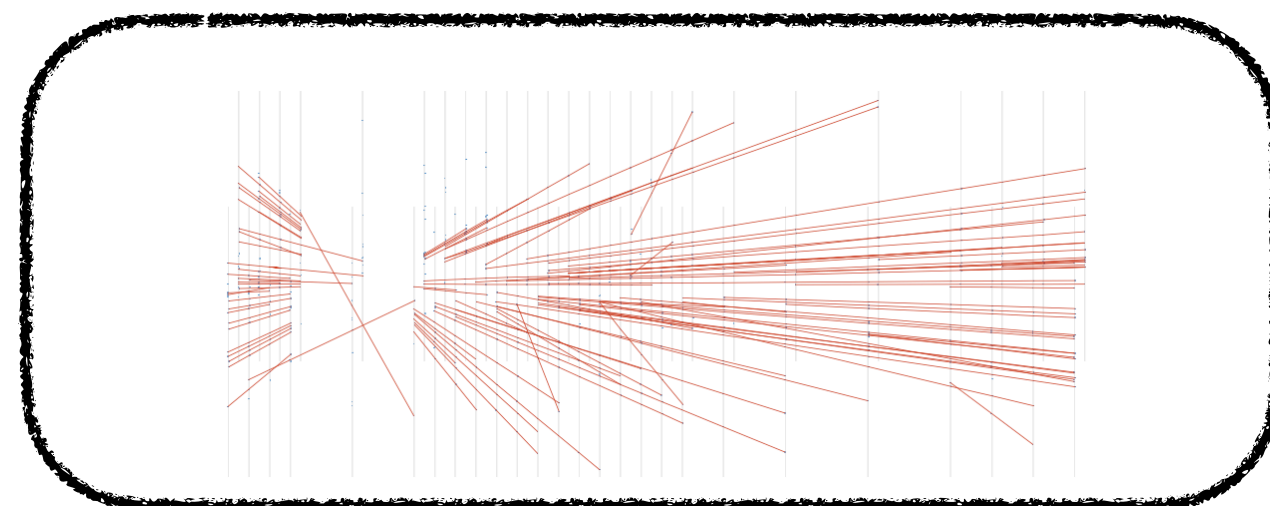
Design of the Allen GPU HLT1

*Mostly designed and built
in France*

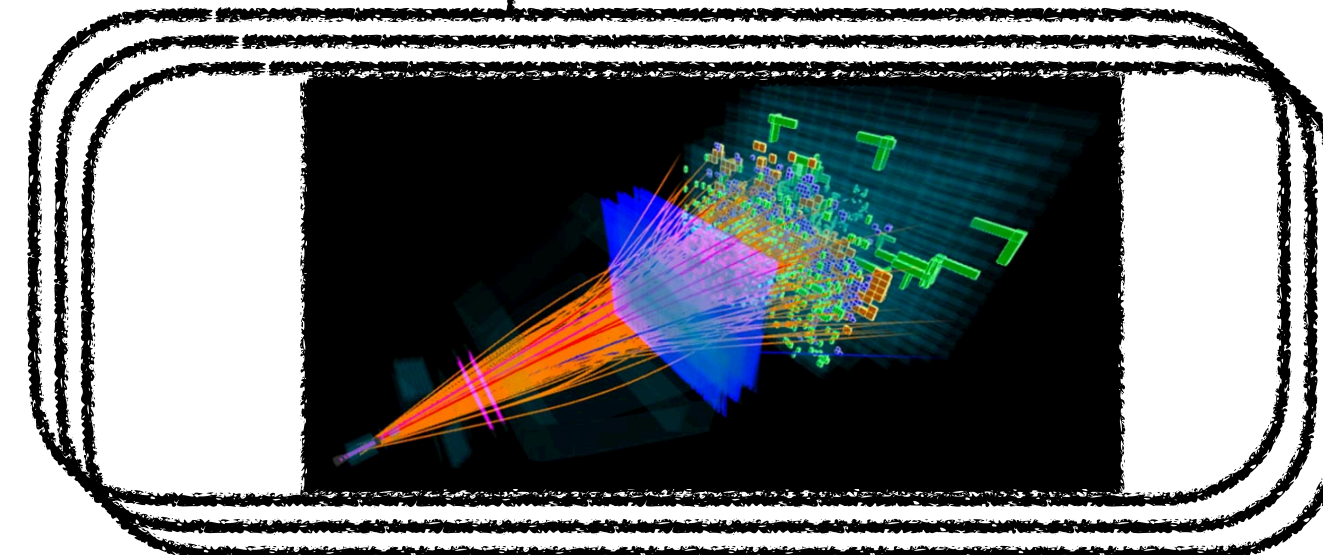
- Do all work on the GPU
 - Minimise copies to/from the GPU
- Parallelise on multiple levels
- Maximise GPU algorithm performance
 - Design software framework to optimise algorithm throughput performance
 - Interleave ML/AI and classical algorithms
 - Single precision only
- Execute on multiple compute architectures
- Simple event model



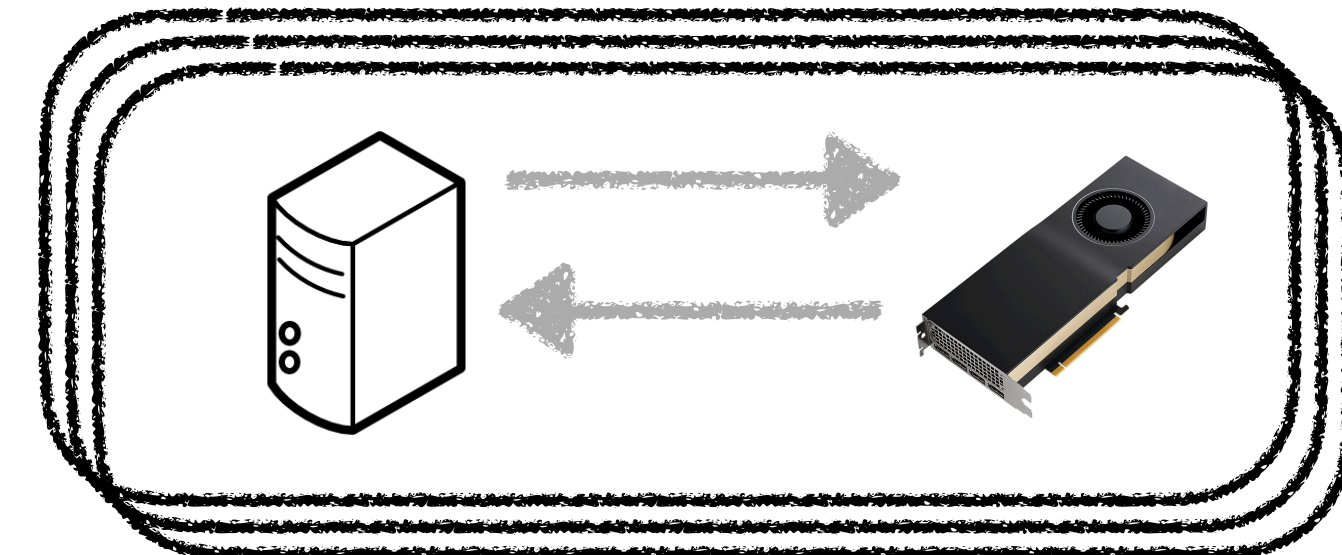
Intra collision: tracks, vertices, ...



Proton-proton collisions



Collision batches

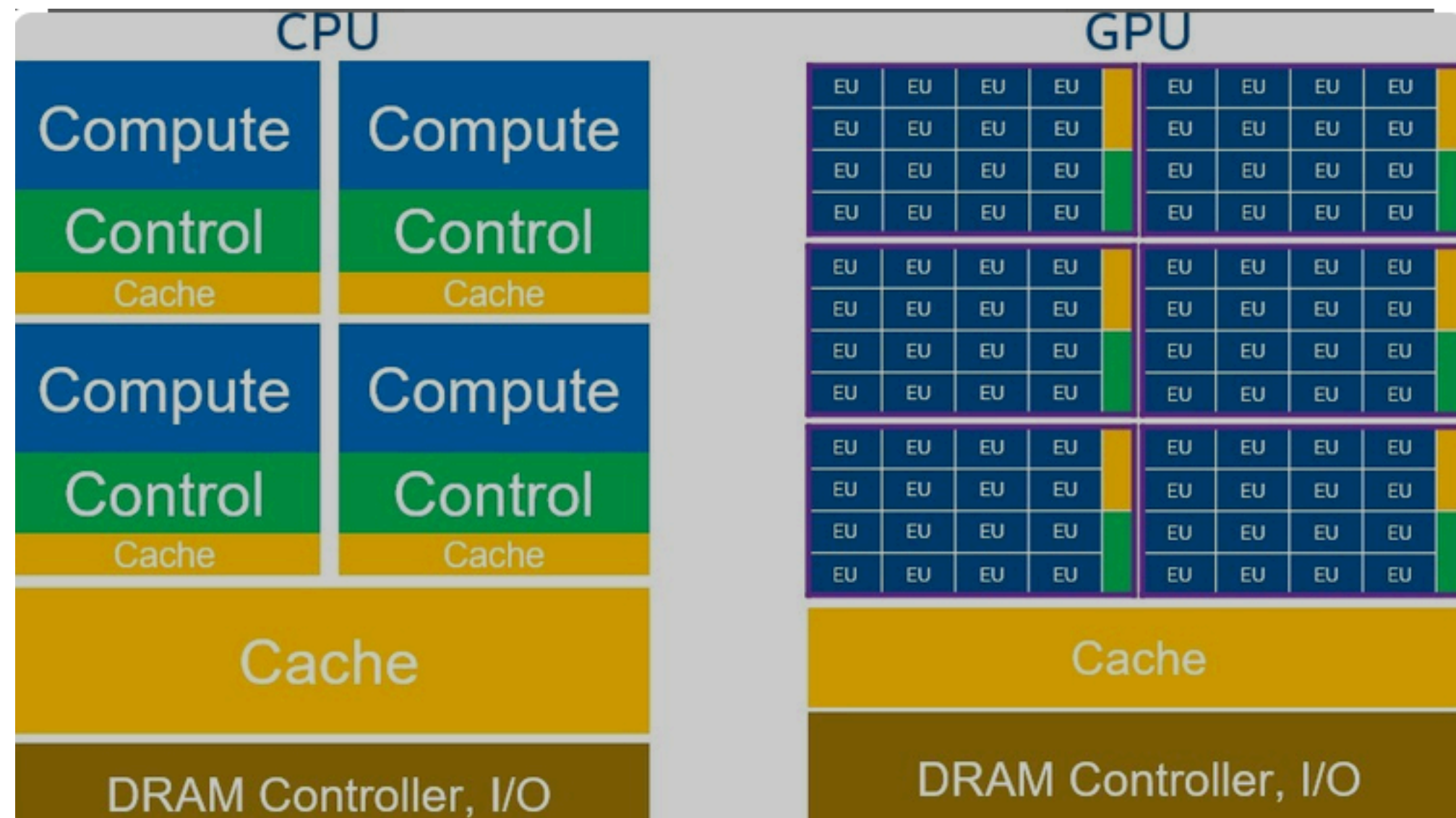


The Allen software framework

- Named after Frances E. Allen
- Hosted on giplab: <https://gitlab.cern.ch/lhcb/Allen>
- Documentation pages: <https://allen-doc.docs.cern.ch/index.html>
- Built with CMake
- Single source code, runs on CPU and GPU (Nvidia and AMD)
 - Portability between architectures provided by macros and few simple coding guide lines
- Standalone build and integrated with LHCb software stack
- Memory manager
- Multi-event scheduler
- Configuration via python
- Monitoring for development and data-taking
- Geometry loaded from DD4Hep, converted to simple structs for easy use on GPU



Allen: Memory management and algorithm scheduling



CPU

- Few strong cores
- Suited for serial workloads
- Quick access to large system memory

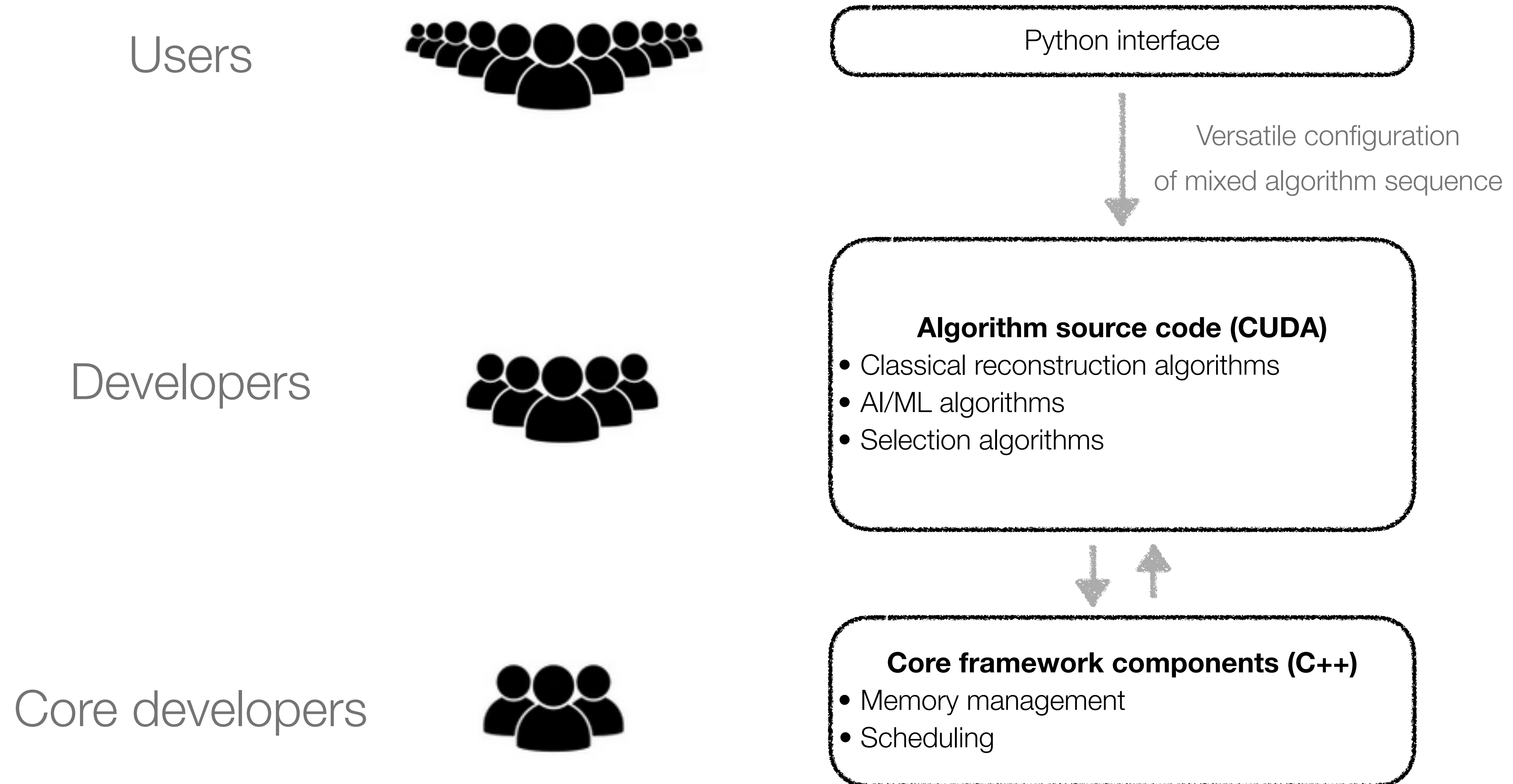
GPU

- Many weaker cores
- Suited for parallel workloads
- Limited memory on card

GPU challenges

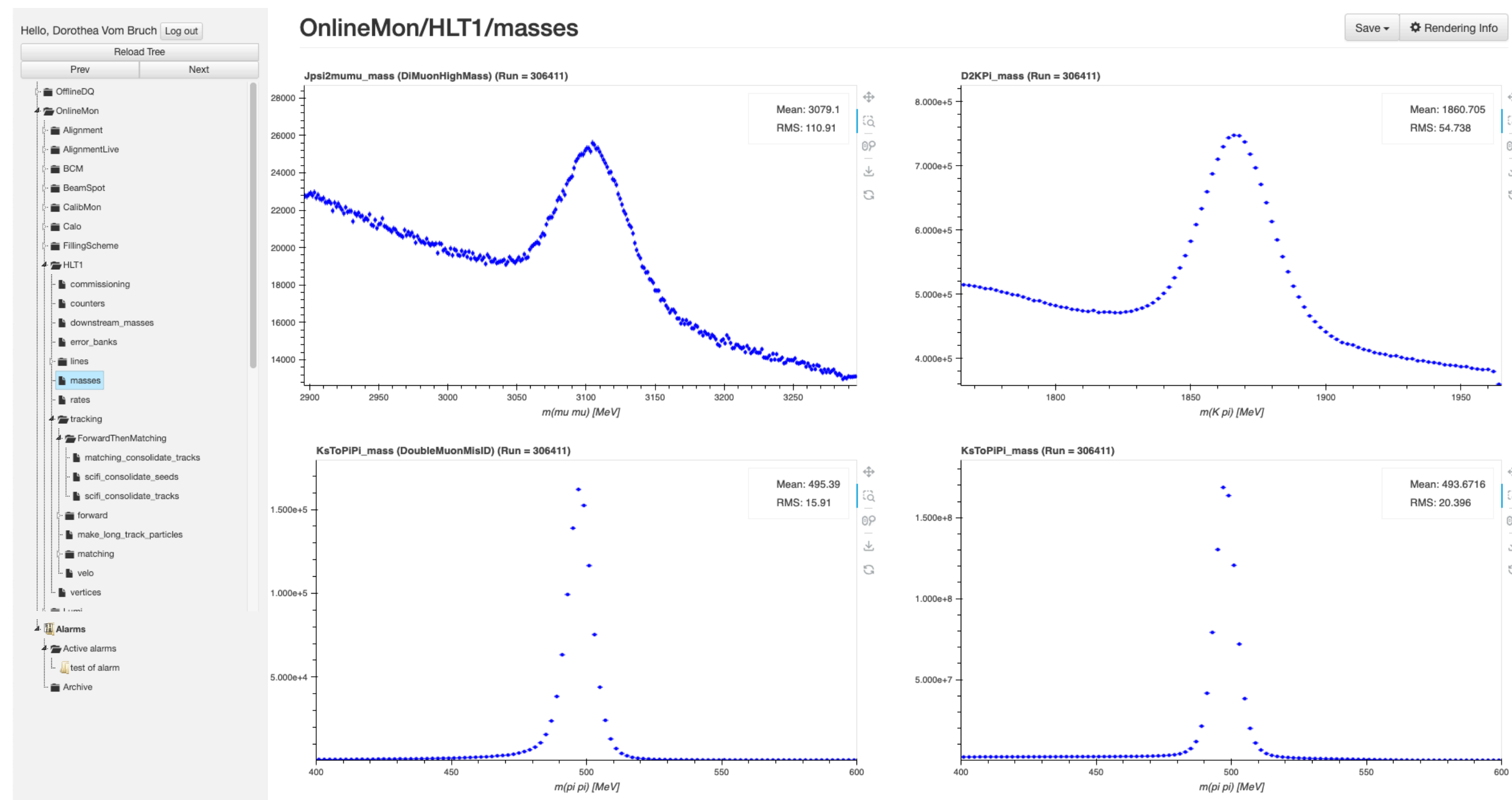
- Fit workload into limited memory resources
 - Allen uses count first - write later model
 - Only allocate as much memory as needed
 - Allen provides memory manager and algorithm scheduler
 - Only keep objects in memory for as long as needed
- Sufficient parallelisable work to utilise compute cores
 - In Allen, every algorithm processes many collision events in parallel
 - « Multi-event scheduler » in Allen generates static sequence of algorithms to be processed using event masks

Allen: Versatile and scalable user interface



Allen: Monitoring

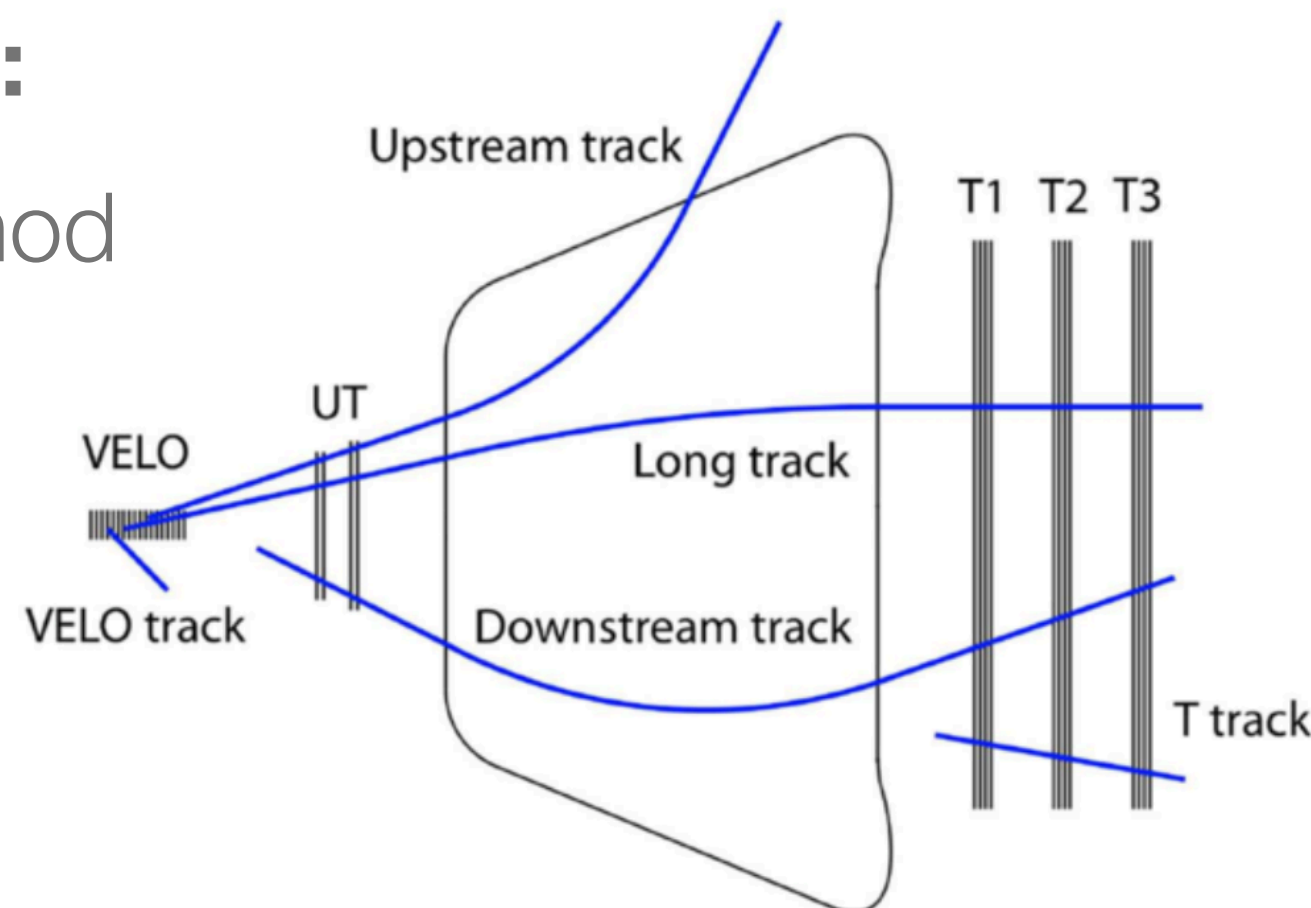
- Ntuple writing for algorithm development
- Histogram and counter filling for monitoring
 - Interfaced with LHCb's monitoring infrastructure « Monet », largely developed in France



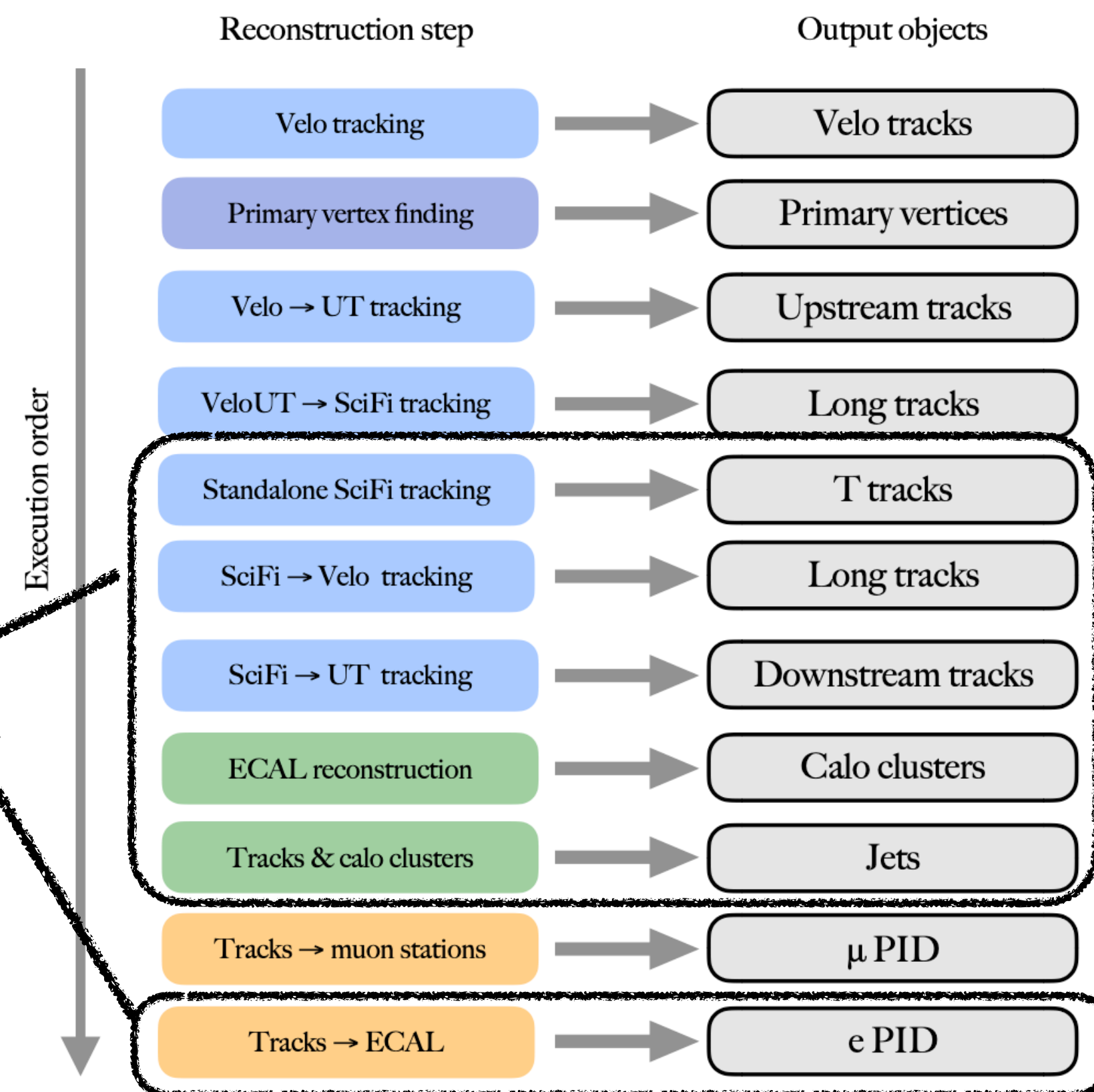
Allen: HLT1 reconstruction in Run 3

Many algorithms beyond TDR design:

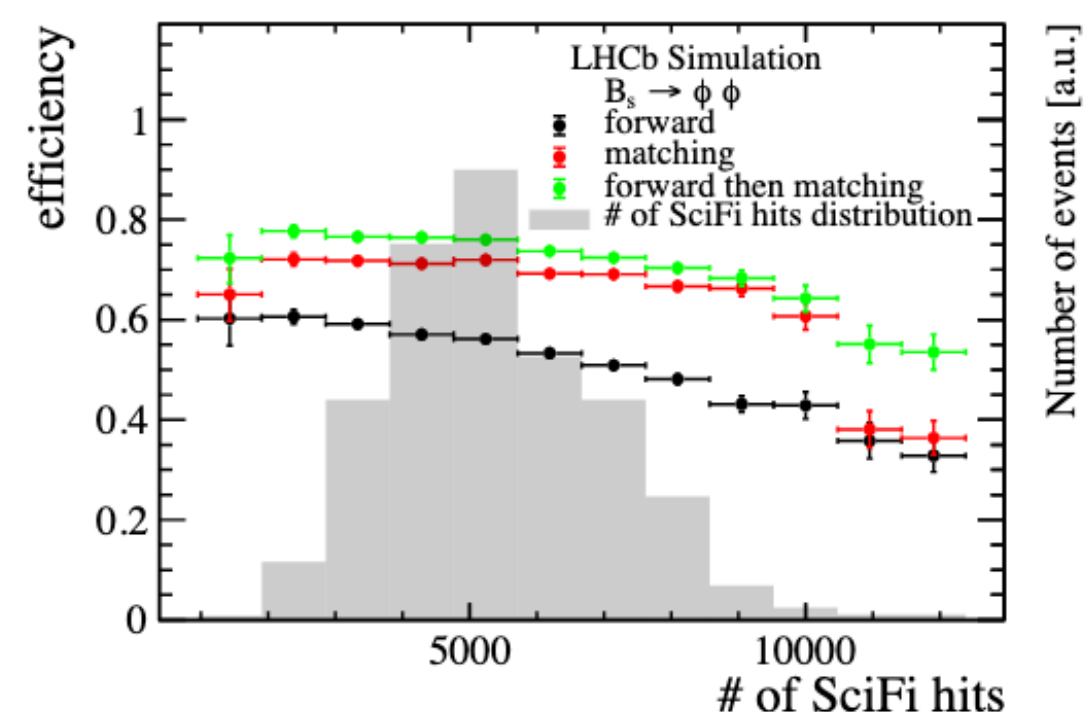
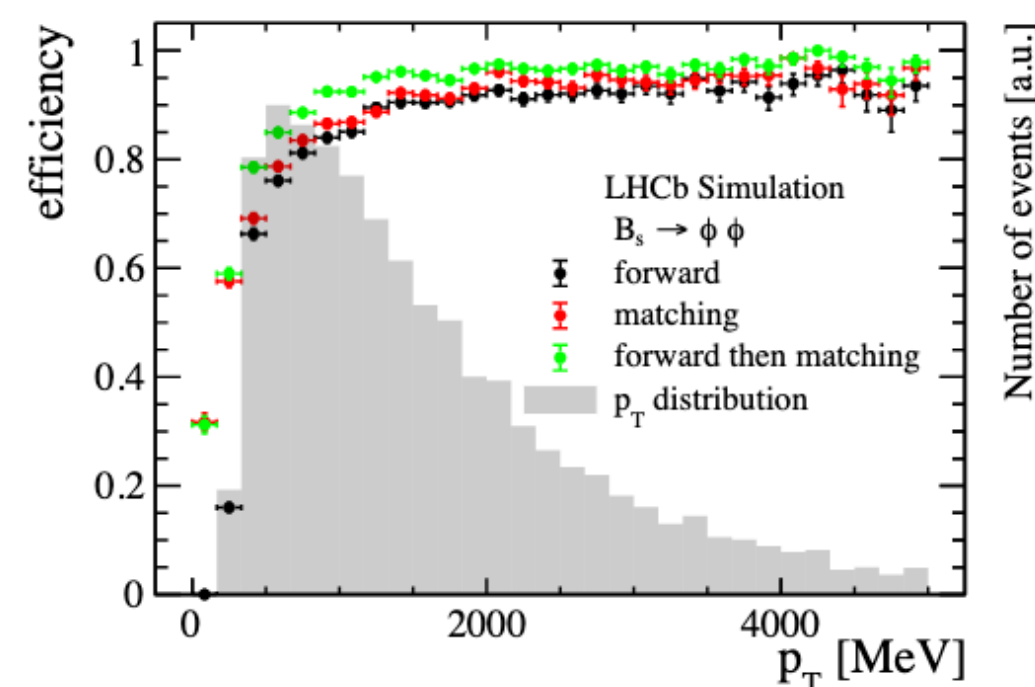
- Additional long track reconstruction method
 - Gained 20% efficiency at $p_T = 500$ MeV
 - Robust with respect to detector occupancy
- Downstream track reconstruction
- Ecal reconstruction
- Jet reconstruction
- Additional PID methods for μ and e



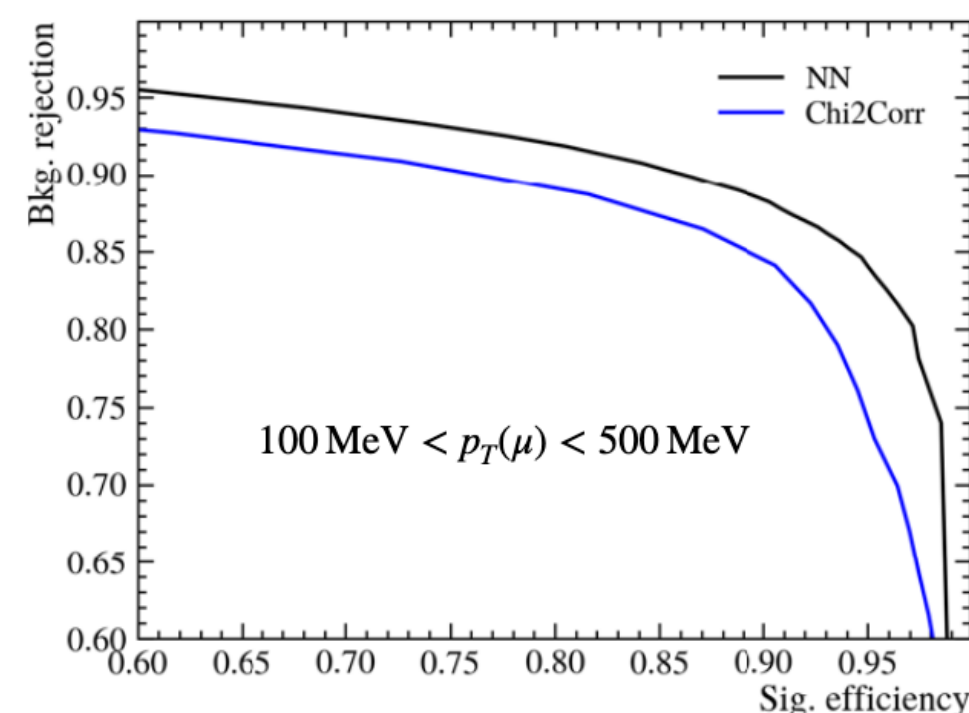
Beyond TDR reconstruction



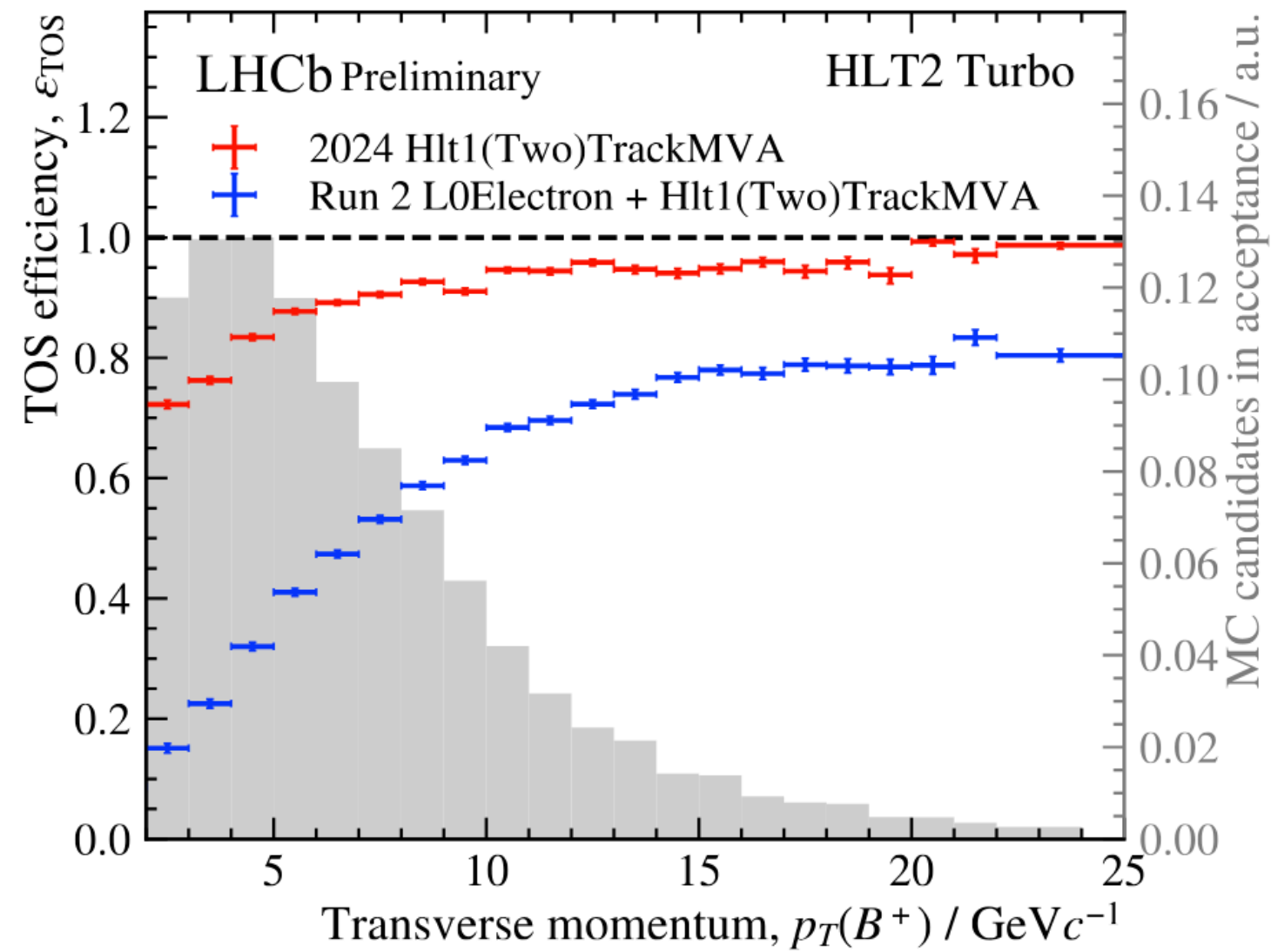
Long tracks form b-hadrons in HLT1



HLT1 Muon PID performance

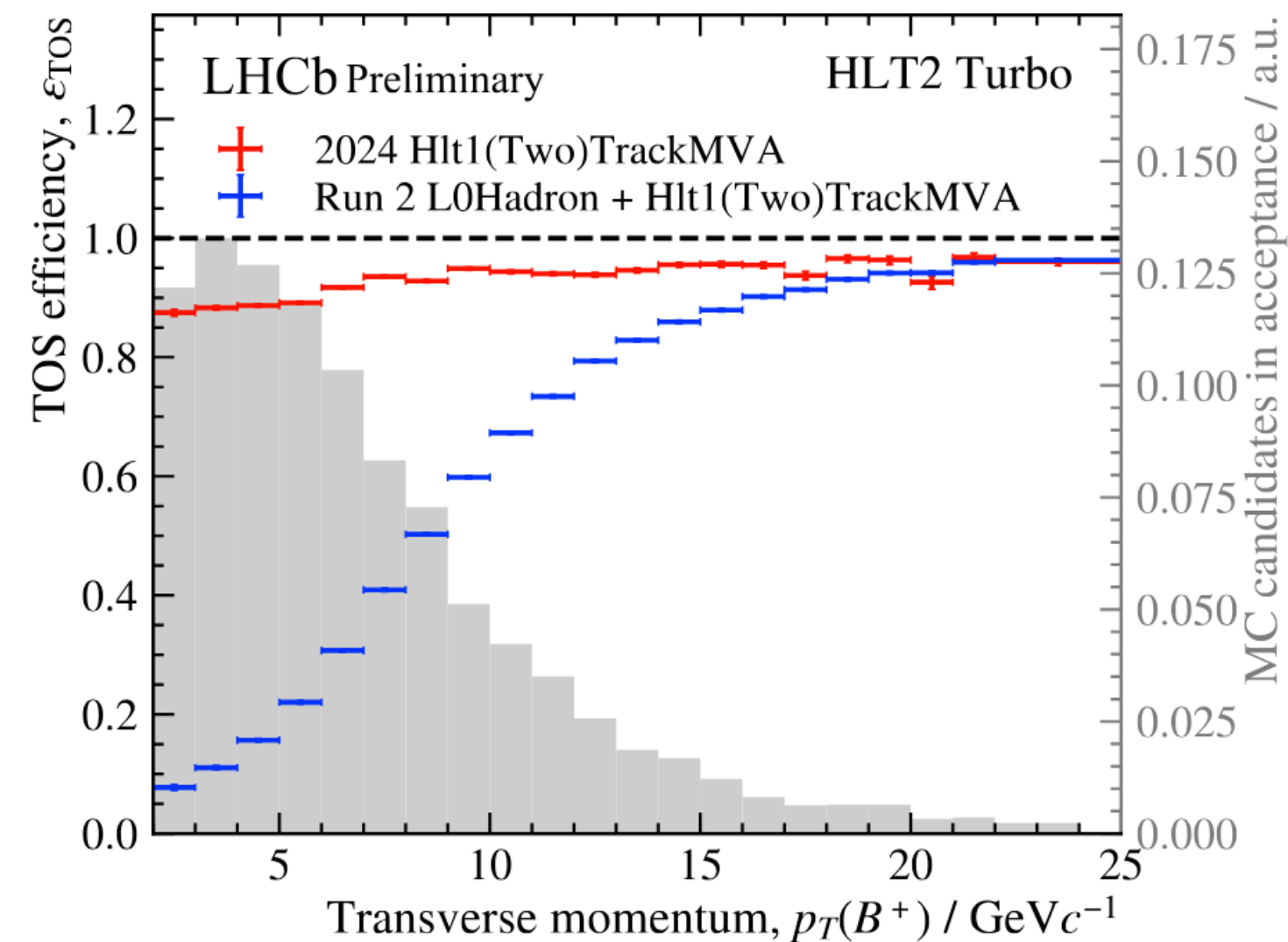


HLT1 in Allen: Run 3 performance

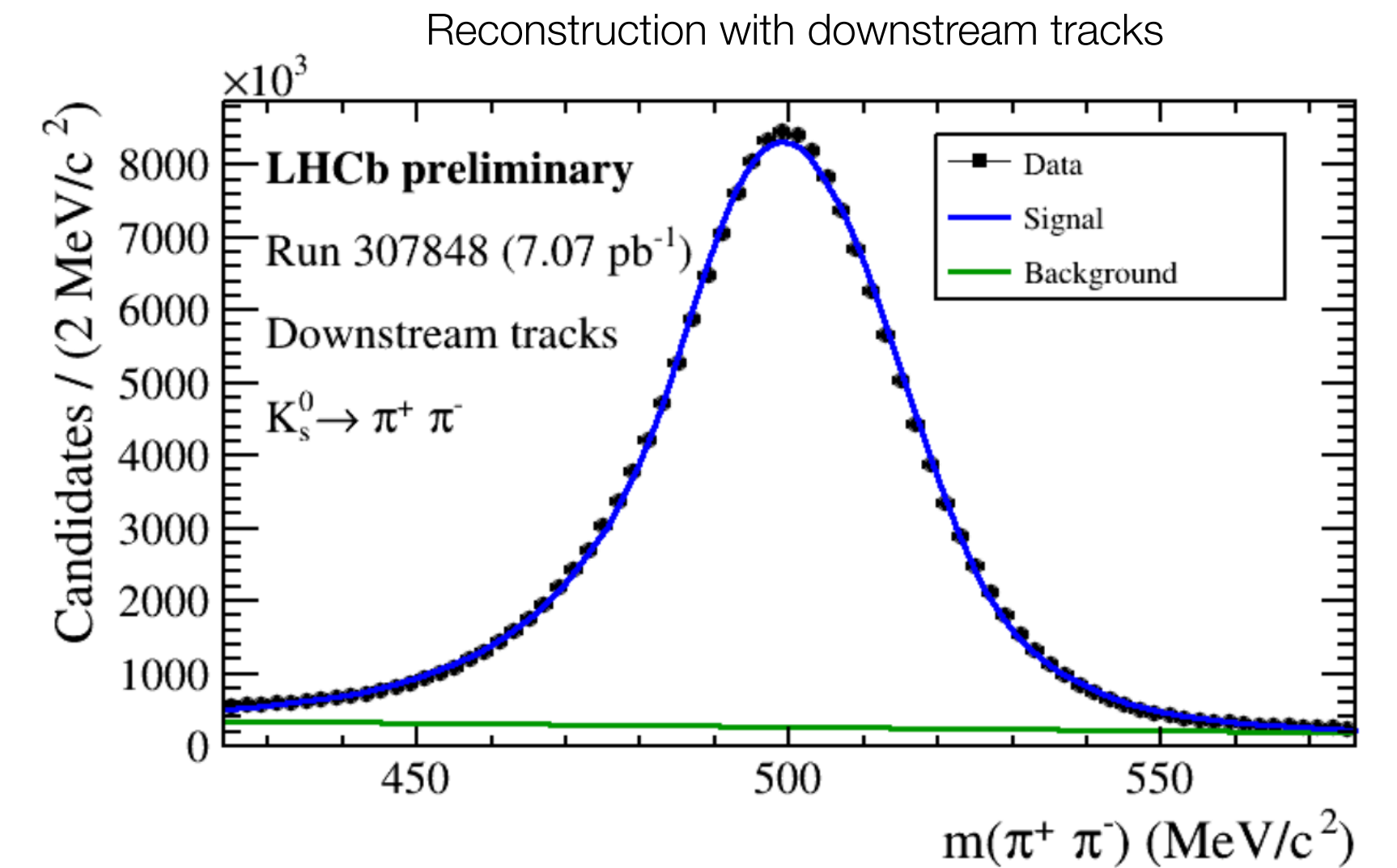


(a) TOS efficiencies in $B^+ \rightarrow J/\psi (e^+ e^-) K^+$.

LHCB-FIGURE-2024-030



(a) TOS efficiencies in $B^+ \rightarrow \bar{D}^0 (K^+ \pi^-) \pi^+$.

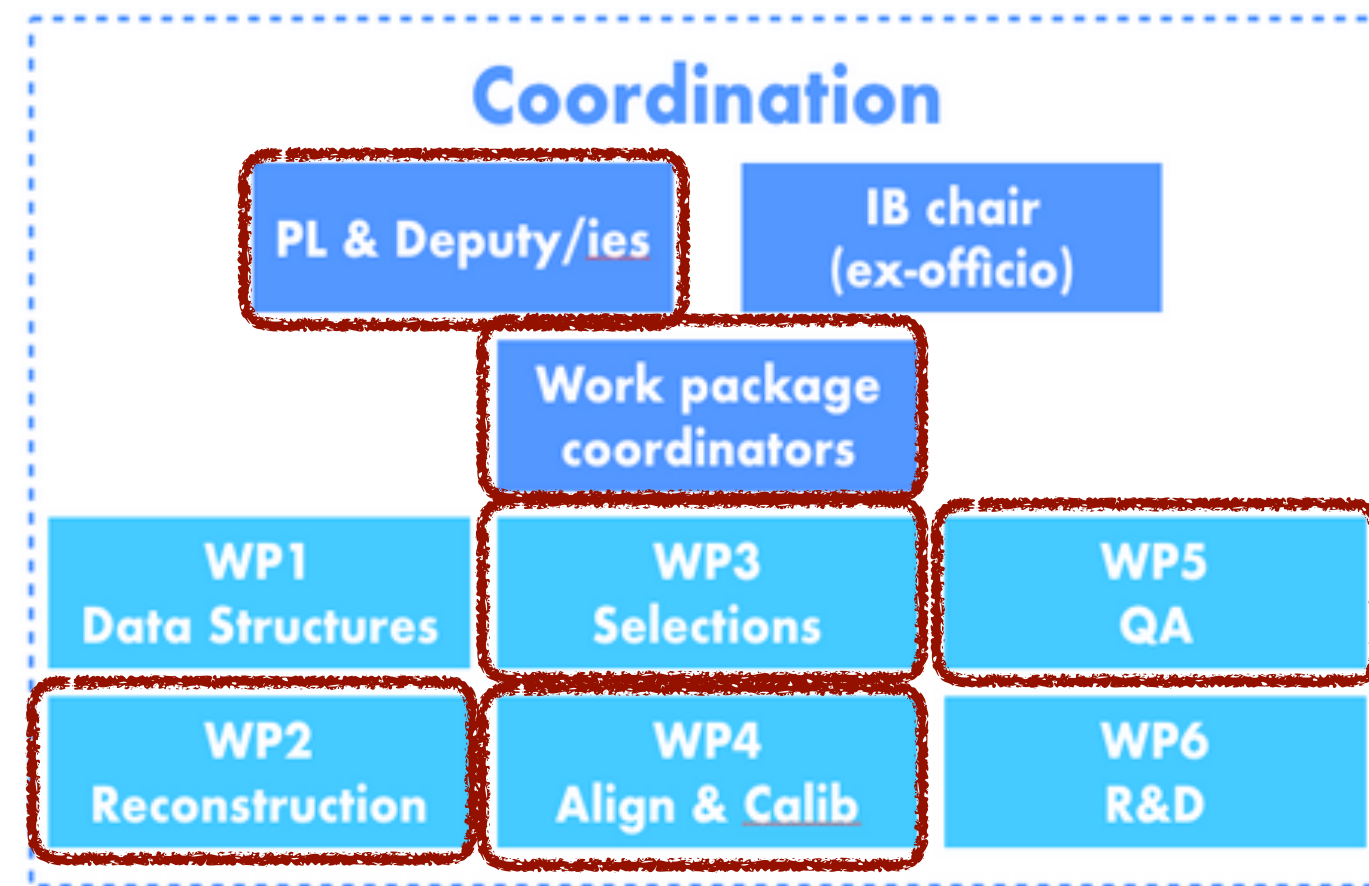


LHCB-FIGURE-2024-035

- Heterogeneous software trigger has been a full success
- HLT1 trigger performance greatly surpasses TDR goals set in 2014 and 2020

RTA for Run 3 - LHCb France involvement

Organisation of LHCb RTA project



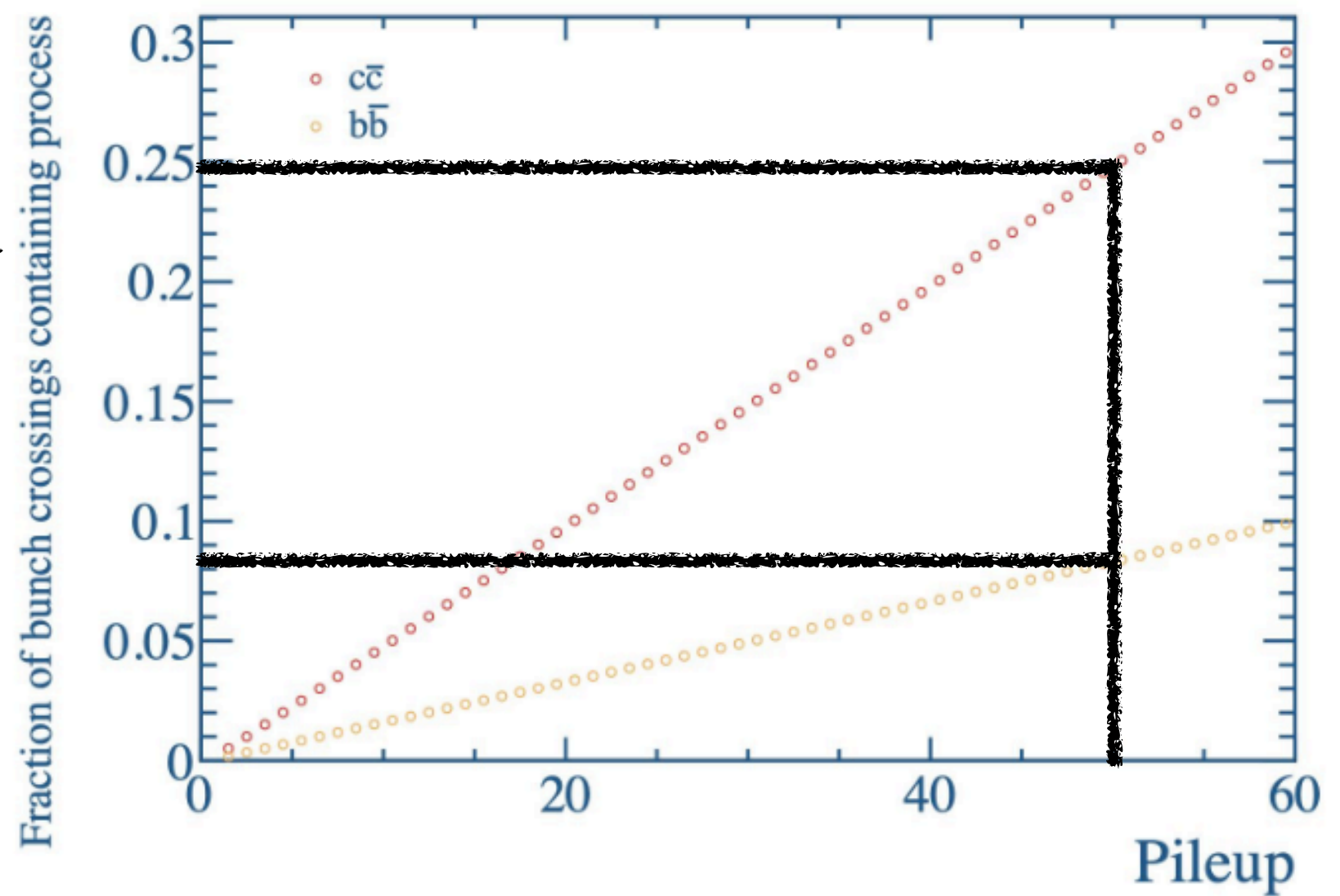
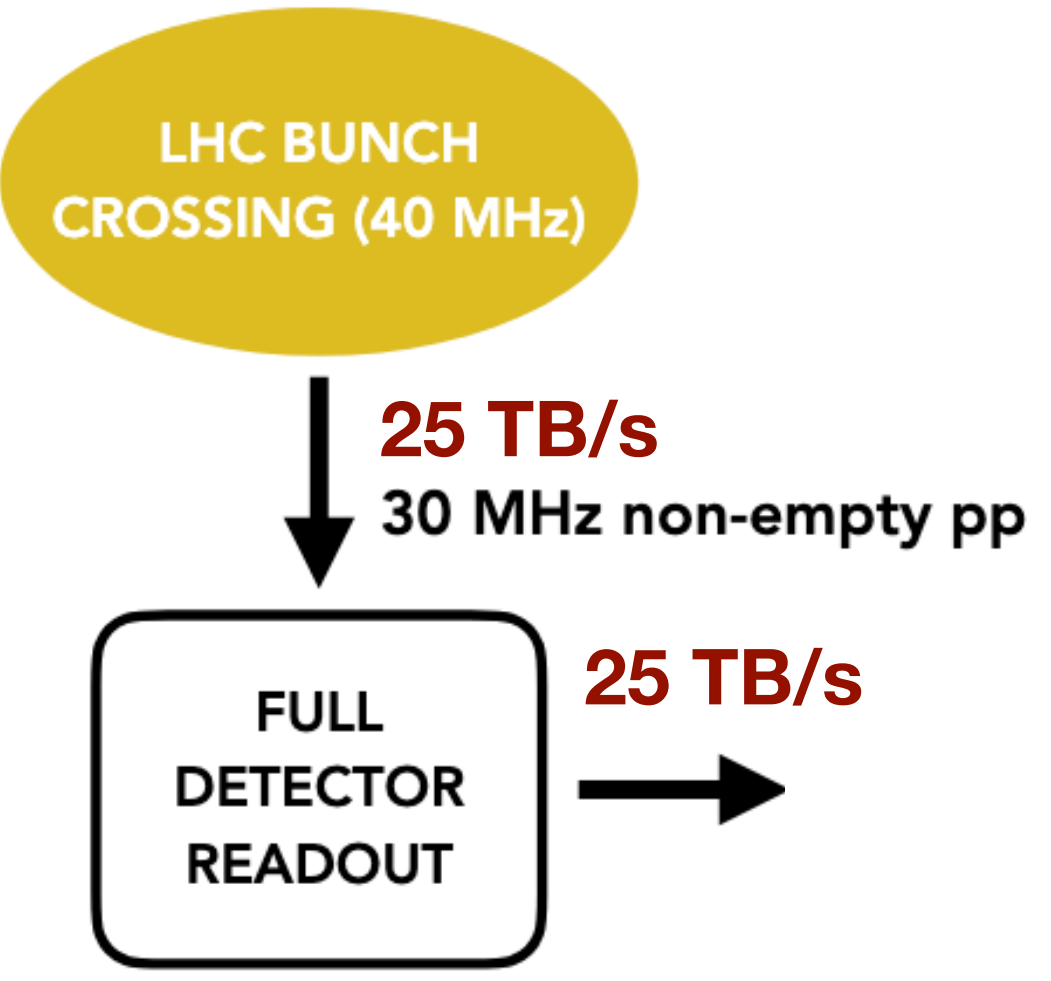
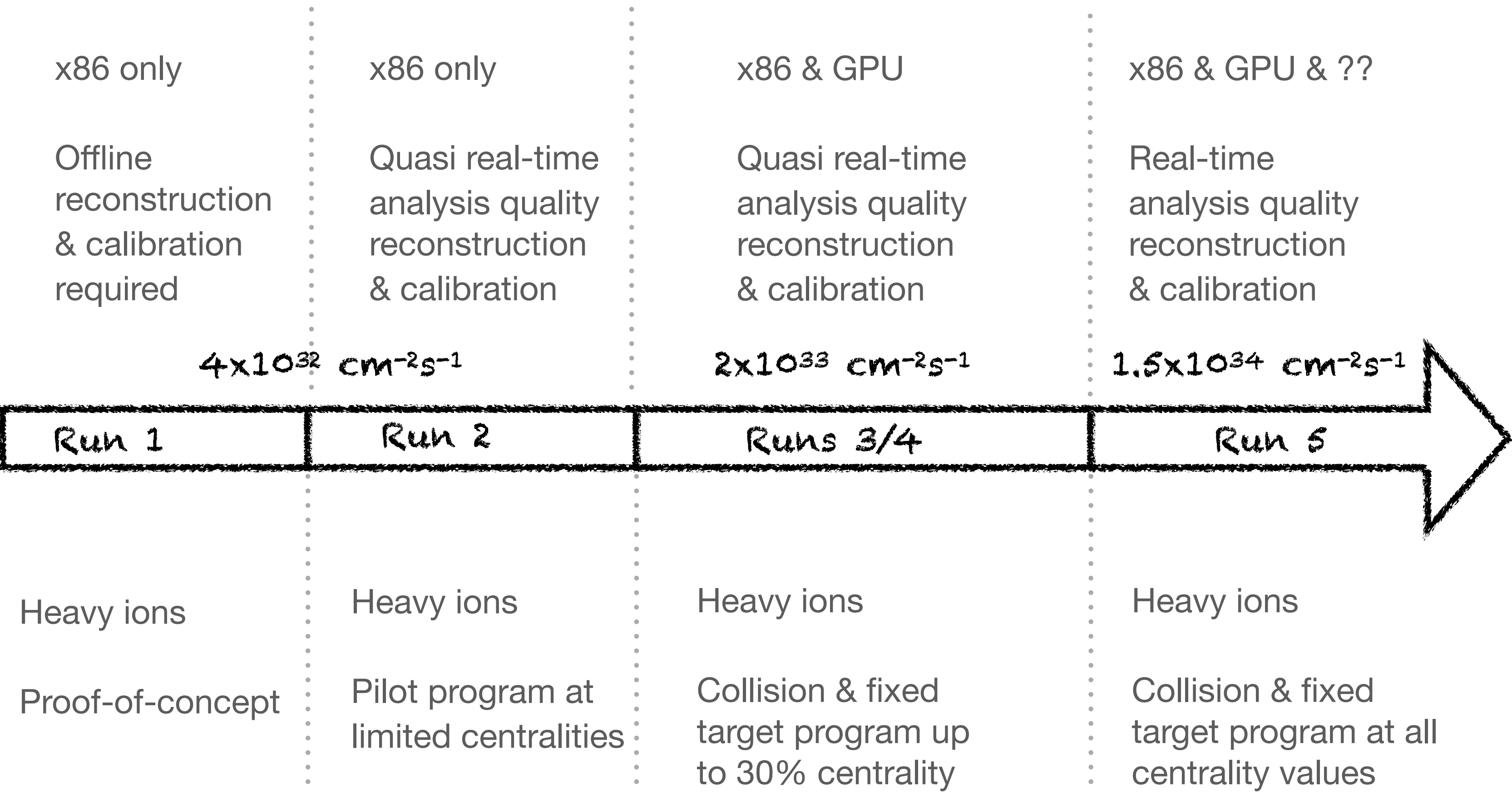
Leading roles of French scientists in the LHCb RTA project

- V. V. Gligorov (LPNHE): Initiator of LHCb RTA project and its project leader from 2019-2022
- D. vom Bruch (CPPM): Co-leader of Allen since 2018
- A. Poluektov (CPPM): « Selections » coordinator 2021-2024
- D. vom Bruch (CPPM): « Reconstruction » coordinator since 2023
- C. Agapopoulou (IJCLab): « Quality Assurance & operations » coordinator since 2023
- A. Poluektov (CPPM): « Alignment & Calibration » coordinator since 2024

French contributions to RTA in Run 3

- One out of three project leaders from France
- 3/10 working group conveners from France (current status)
- 2/4 ERC funded projects related to LHCb RTA hosted in French labs

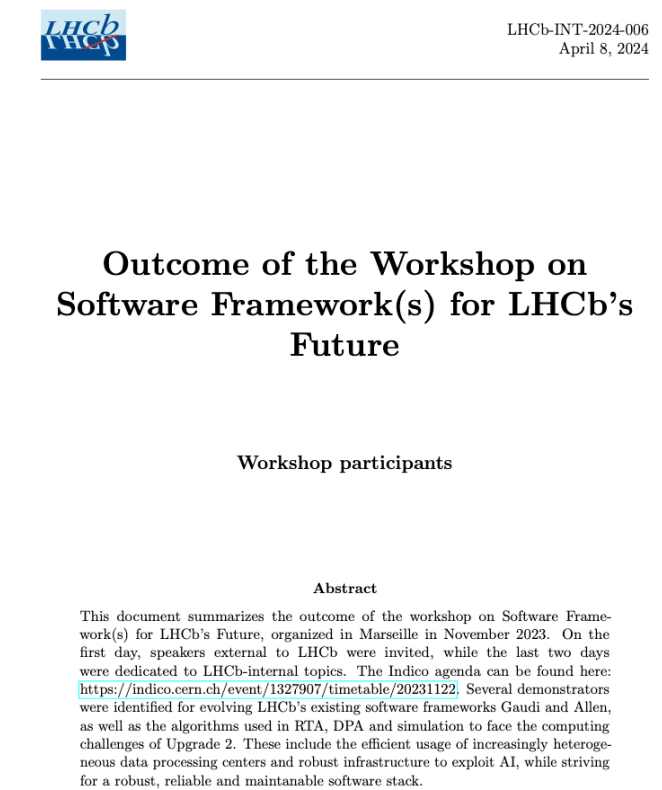
The LHCb trigger challenge in U2



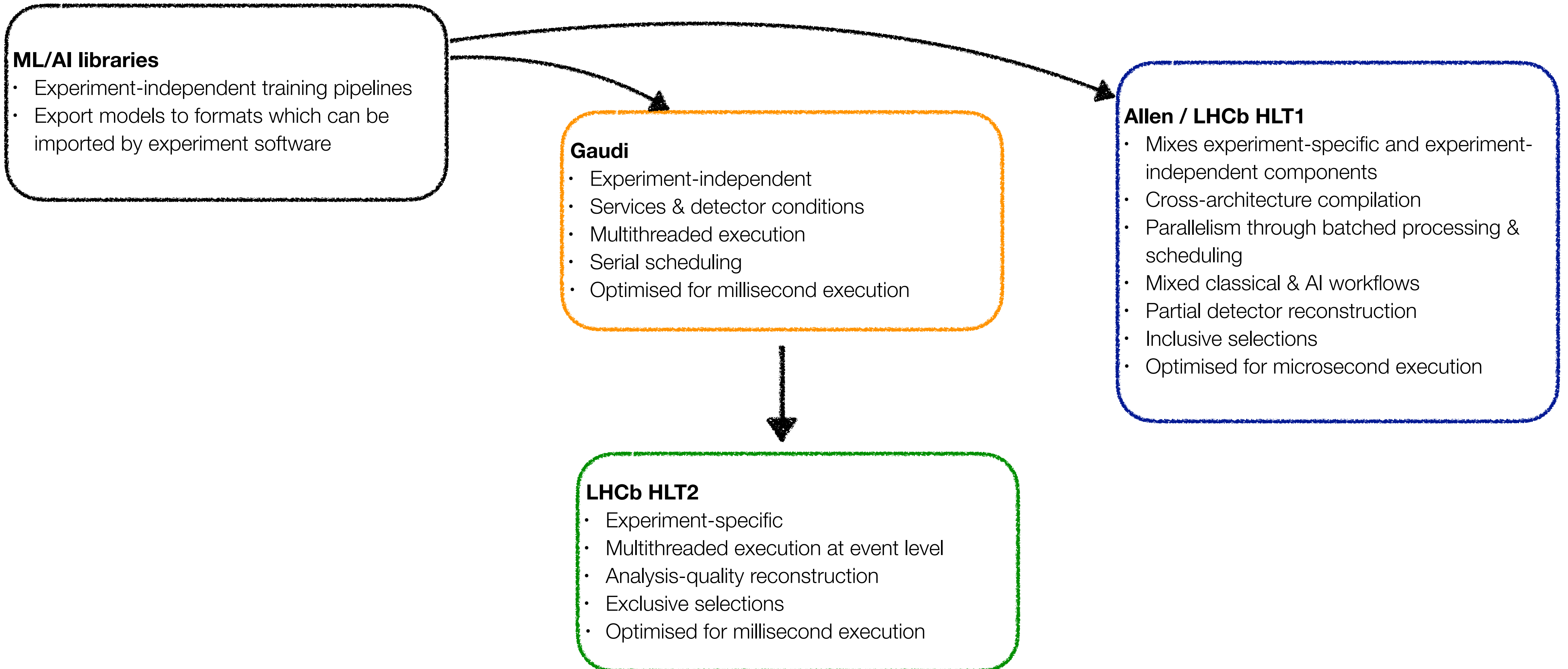
pp analysis-level performance at same level, despite higher pileup

RTA for LHCb Upgrade 2

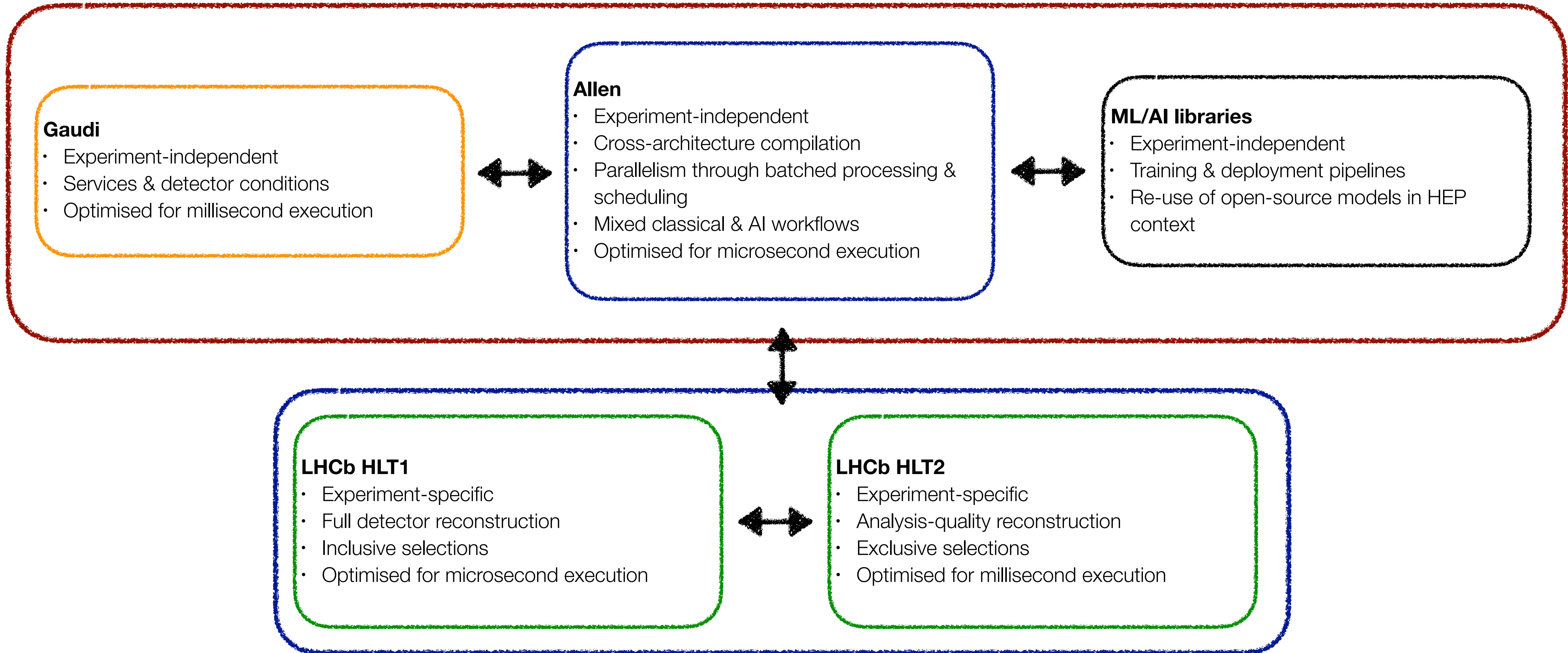
- Discussions on the design of the RTA system and the evolution of the software frameworks are ongoing within LHCb, with strong contribution from French physicists
 - Organised a workshop on future software frameworks for LHCb in Marseille in 11/2023
 - Co-organised Computing Workshop (01/2025 in A Coruna)
 - Next Computing Workshop in 2026 to be held in Orsay
- Only viable solution as of today is to process the full HLT2 reconstruction on GPUs
 - Software frameworks need to scale for this challenge
- Preparing an RTA document in 2026 to accompany sub-detector TDRs
- RTA TDR will be written in 2030



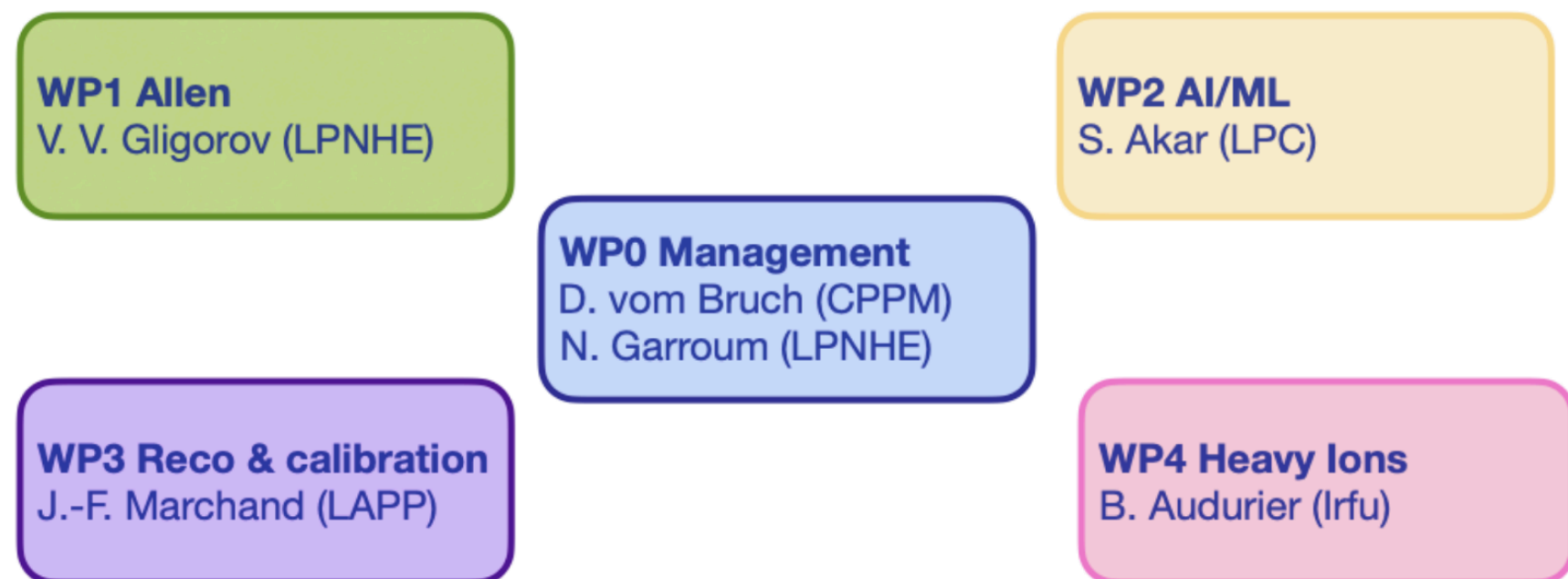
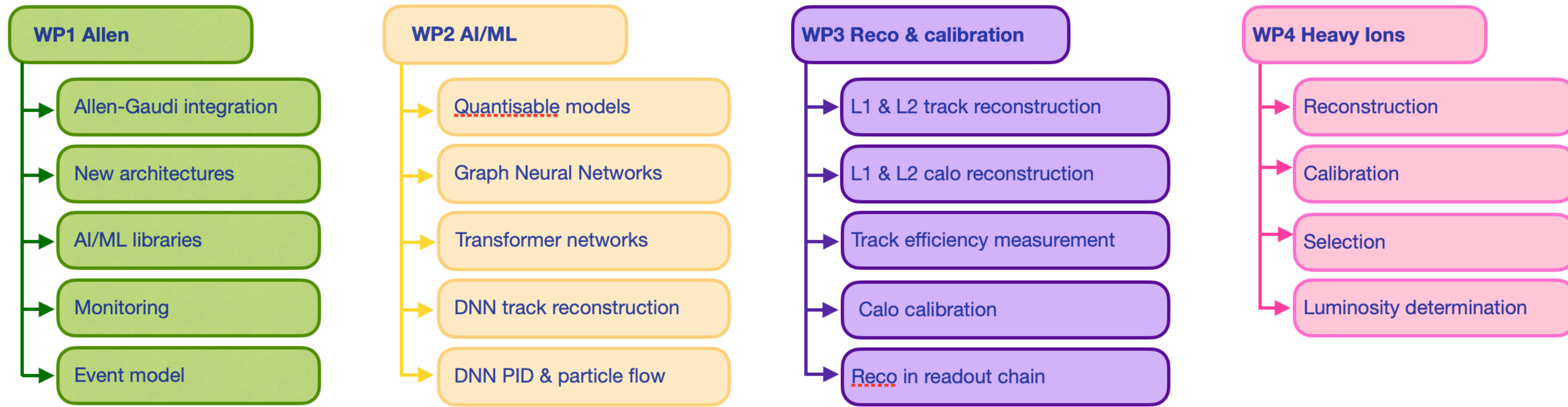
Current LHCb software



Vision of future LHCb software



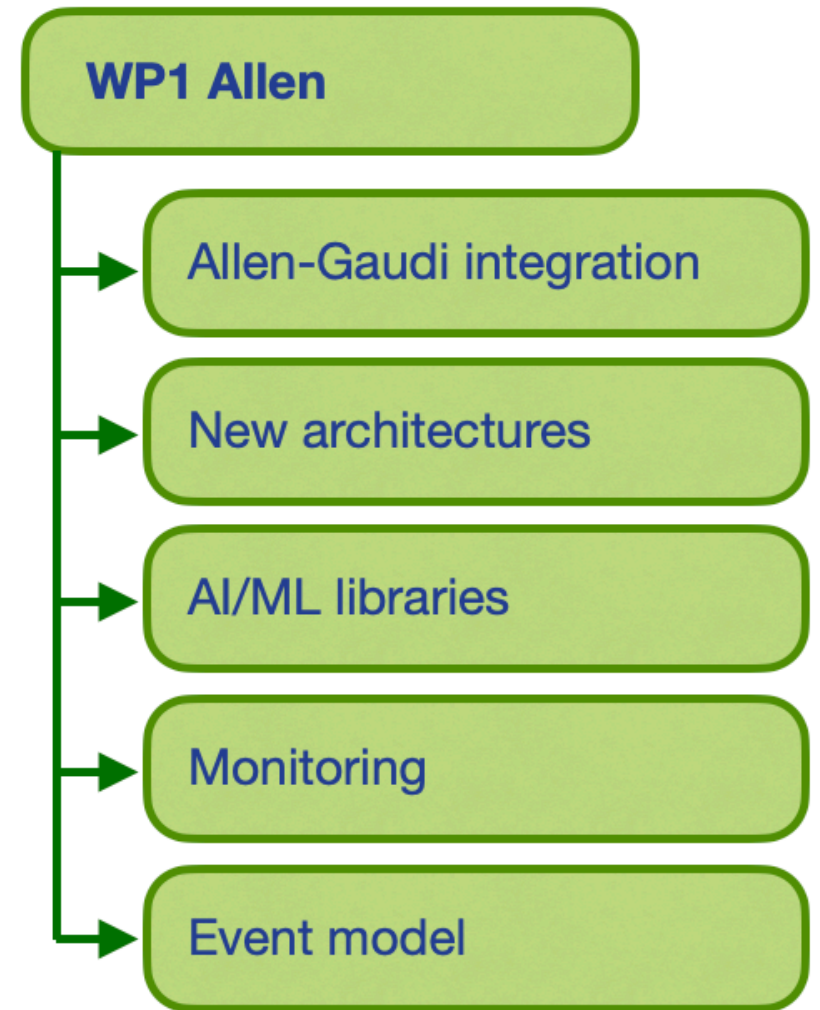
Proposed RTA-LHCb-U2-France project



Planned contributions within RTA-LHCb-U2-France project:

- Evolve Allen to cope with HLT2 reconstruction on GPUs
- Explore AI/ML methods for reconstruction scalability to U2 conditions
- Author reconstruction algorithms within Allen related to sub-detectors developed in France
- Adapt the reconstruction to heavy ion conditions

Planned outcome of WP1 « Allen »



Allen-Gaudi integration

- Enable Gaudi services in Allen
- Single memory manager for Gaudi & Allen algorithms:
 - Memory management from the GPU in Allen, remove transient event store in Gaudi
- R&D for graph-based scheduling (Cuda graphs, tbb)
- Single multi-event scheduler in Gaudi & Allen
- Harmonize algorithm syntax

AI/ML libraries

- Common ML interface for frontend: where models are deployed
 - For example using ONNX standard
 - Deployed on the architectures supported by Allen
- Fast inference for small modes is a unique challenge in LHCb

Monitoring

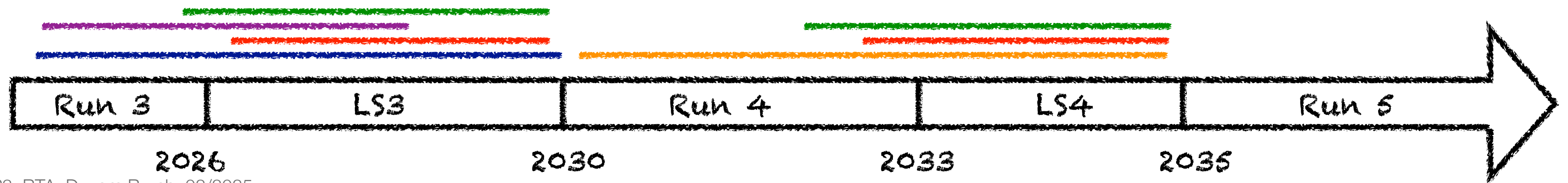
- Make Allen's monitoring more user-friendly
- Extend Gaudi's monitoring infrastructure to Allen

New architectures

- Enable ARM build of Allen
- Explore emerging architectures
 - FPGAs
 - RISC-V co-processors
 - ...

Event model

- Make Allen's event model more user-friendly
- Harmonize Gaudi and Allen transient event model
- Evolve event model for 4D reconstruction
- Enable high-level object persistency



Planned outcome of WP2 « AI/ML »

WP2 AI/ML

Quantisable models

Graph Neural Networks

Transformer networks

DNN track reconstruction

DNN PID & particle flow

Quantisable models

- Large models have to be quantised to fit into limited memory resources
- Develop generally usable and maintainable framework to quantise models

Graph Neural Networks (GNNs)

- Studied in detail for track finding and jet building in the last years
- Can have prohibitive memory usage
- R&D to understand how they can be used in limited resources at high throughput

DNN track reconstruction

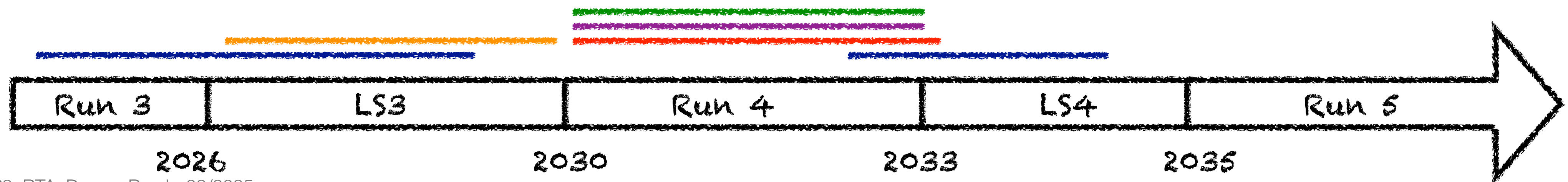
- Evolve towards fully NN based track finding and vertex reconstruction
- Include timing information

DNN PID & particle flow

- Update NN-based PID algorithms to include time
- R&D for DNN-based particle flow including calorimeter, tracking systems and hadron PID

Transformer networks

- Have gained prominence through Large Language Models (LLMs)
- R&D how to bridge low-level particle reconstruction and high-level physics analysis
- Could reduce maintenance and computational costs of thousands of exclusive selections



Planned outcome of WP3 « Reconstruction & Calibration »

WP3 Reco & calibration

L1 & L2 track reconstruction

L1 & L2 calo reconstruction

Track efficiency measurement

Calo calibration

Reco in readout chain

L1 & L2 track reconstruction

- Common denominator of historical French reconstruction developments is seeding
- Develop seeding techniques for parallel architectures
- Develop common methods to describe magnetic field map and material distributions
- Develop common methods to cope with misalignment and detector inefficiencies

L1 & L2 calo reconstruction

- Introduction of precise time measurement
- Development and benchmarking of AI/ML based clustering algorithms (interface with WP2)
- Integration of L2 reconstruction on GPUs
- Benchmark different implementations

Reconstruction in readout chain

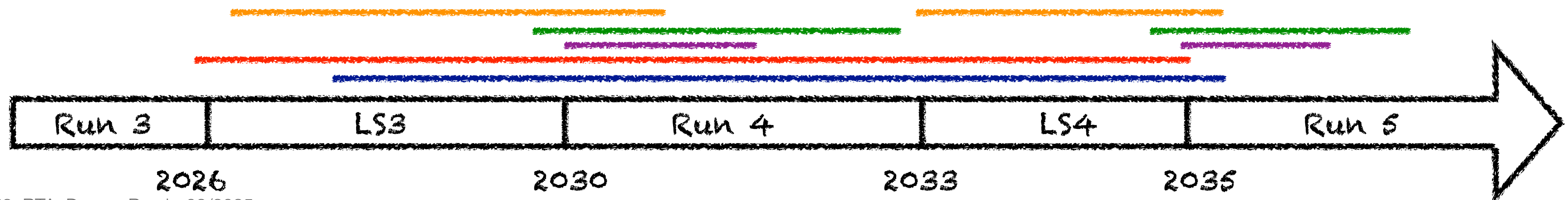
- R&D to evaluate which part of calo reconstruction can be carried out on the readout boards
- R&D for clustering or seeding of tracking detectors

Track efficiency measurement

- Historically, French groups are involved in the data-driven measurement of tracking efficiencies for electrons
- Evolve tools and perform measurements for Run 4 and Run 5

Calo calibration

- Historically, French groups implemented the automation of both absolute Ecal calibration using pi0 and relative Ecal to Hcal calibration using LEDs
- Update tools for Runs 4 & 5



Planned outcome of WP4 « Heavy Ions »

WP4 Heavy Ions

Reconstruction

Calibration

Selection

Luminosity determination

Reconstruction

- Co-development of NN-based track reconstruction and of seeding algorithms capable of reconstructing most central collisions
- Development of efficient reconstruction algorithms at the high occupancies of PbPb collisions
- Development of algorithms with low ghost rate at high occupancies

Calibration

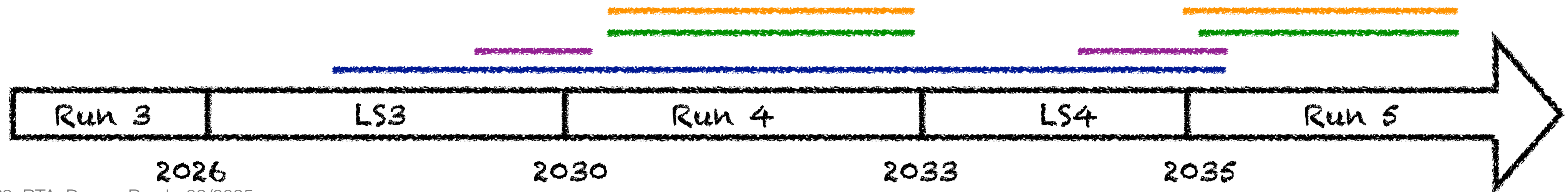
- Adjust references for ion and fixed-target operations
- Assess relationship between detector quantities and value of centrality in each collected data sample

Selection

- Adapt selections to improved reconstruction performance for central collisions

Luminosity determination

- Determine luminosity for unique heavy ion and fixed target conditions

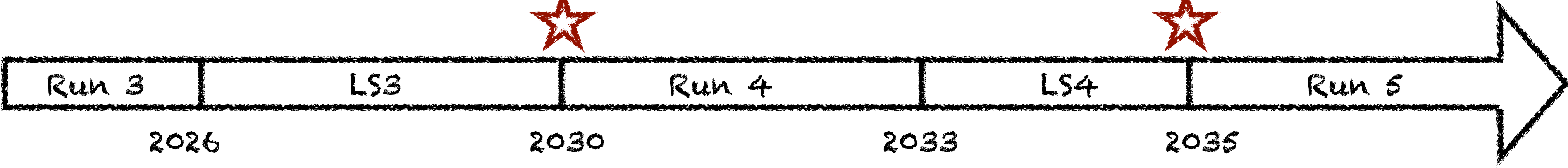


Timeline of RTA-LHCb-U2-France project

Generic developments

Allen integrated with Gaudi,
Allen core-functionality
for high-throughput heterogeneous computing
in standalone framework

Allen core-functionality
Extended for emerging architectures,
AI/ML support, energy-efficiency evaluation



LHCb-specific

LHCb reconstruction
For pp and heavy ions
ready for Run 4

Prototypes of ML/AI-based
reconstruction ready

LHCb reconstruction
For pp and heavy ions
ready for Run 5

Resource requirements from IN2P3

- RTA France community has been a key player in delivering the successful Run 3 trigger system
 - So far mostly financed with external funding (2 ERCs, 1 European infrastructure project)
- From 2028 onward, 0% of required engineers are covered
- Skills required:
 - Large software project maintenance & management
 - Development of AI/ML tools for reconstruction
 - Development of common AI/ML tools for inference
 - Development of core software tools for heterogeneous architectures
 - Programming for efficient high-throughput applications
 - Exploring new computing architectures

	Run 3	LS3					Run 4				LS4		Total
	2025	2026	2027	2028	2029	2030	2031	2032	2033	2034	2035		
Physicists													
WP1	0	0.5	0.5	0	0	0	0	0	0	0	0	1	
WP2	0	0	1	1	1	2	2	2	0.2	0.2	0.2	9.6	
WP3	0.3	0.55	0.75	1.25	1.75	2.05	2.05	1.3	0.3	0	0	10.3	
WP4	0	0	0	0.5	0.7	0.5	0.3	0.6	0.6	0.9	0.4	4.5	
Management	0.25	0.25	0.25	0.25	0.25	0.25	0.25	0.25	0.25	0.25	0.25	2.75	
Total	0.55	1.3	2.5	3	3.7	4.8	4.6	4.15	1.35	1.35	0.85	28.15	

	Run 3	LS3					Run 4				LS4		Total
	2025	2026	2027	2028	2029	2030	2031	2032	2033	2034	2035		
Engineers													
WP1	0.25	2.25	2	2	1.05	1	1	1.8	0.8	0.3	0.3	12.75	
WP2	0	0	0	0.3	0.3	0.2	0.2	0.2	0.2	0.2	0.2	1.8	
WP3	0	0	1	0.5	1.5	1.25	0.95	0.45	0.95	1.1	0.7	8.4	
WP4	0	0	0	0	0	0	0	0	0	0	0	0	
Management	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	1.1	
Total	0.35	2.35	3.1	2.9	2.95	2.55	2.25	2.55	2.05	1.7	1.3	24.05	

Table 4: Requested hardware support for the LHCb France RTA project.

Timescale	Request (k€)	Purpose
2026–2028	50	Two high-end CPU servers
2032–2034	50	Six co-processors (2 GPUs, 2 FPGAs, 2 cards with emerging architectures (such as RISC-V))

RTA impact on engineering skills of the IN2P3 community

- RTA allows the development of high level software engineering skills
- Connection between Allen/RTA and AI/ML is becoming a strategic development
- Synergies are developed through IN2P3, CERN and HEP community networks
 - Reprises IN2P3 engineering group for High Performance Computing
 - Gray Scott Challenge
- Synergies respond to critical need of identified engineering skills in IN2P3

Identification of RTA software engineering skills in the IN2P3 PECTIN database

Skill	Tension	Type de tension
324-Développement de systèmes temps réel	Importante	Importante
342-Logiciels ML-IA	Critique	Nombre
348-Ingénierie algorithmique	Critique	Nombre
324-Développement de systèmes temps réel	Critique	Nombre
342-Logiciels ML-IA	Critique	Nombre
320-Architecture logicielle	Importante	Technicité

Risks

- Slow-down of development cycle due to large legacy codebase
 - Modularise tools during Gaudi-Allen convergence, only use the required ones for U2
 - Preserve slimmed down Allen core test-bench for fast prototyping
- Development of commercial technologies
 - Design code with abstraction layers in modular way
 - Test new architectures in test bench as they emerge
- Slower progress on project than anticipated
 - Margins included in estimated FTE
 - Develop part of core functionality already during LS3, to test in production during Run 4
- No knowledge transfer due to non-appointment of permanent staff
 - Structured onboarding with documentation and training sessions
 - Long-term development is threatened if critical mass of permanent project members is not available

RTA-France U2 involvement - in the LHCb context

- Allen developments for U2 proposed from France only
- France is so far the only country where a country-specific RTA project for U2 is proposed
- Chance to keep leading role for LHCb RTA developments
- Chance to develop standalone Allen framework to be used outside of LHCb

Table 3: Distribution of RTA areas of interest across French LHCb teams.

Institution	Area of interest
WP1 Allen	CPPM, LPNHE, Subatech
WP2 ML/AI	LPCA, LPNHE, Subatech
WP3 Reconstruction & Calibration	CPPM, IJCLab, LAPP, LPNHE, LLR
WP4 Heavy Ions	Irfu, LLR

Interested IN2P3 permanent personnel

Physicists

C. Agapopoulou (IJCLab), S. Akar (LPCA), Y. Amhis (IJCLab), F. Fleuret (LLR), V. V. Gligorov (LPNHE), J.-F. Marchand (LAPP), Émilie Maurice (LLR), A. Poluektov (CPPM), D. vom Bruch (CPPM)

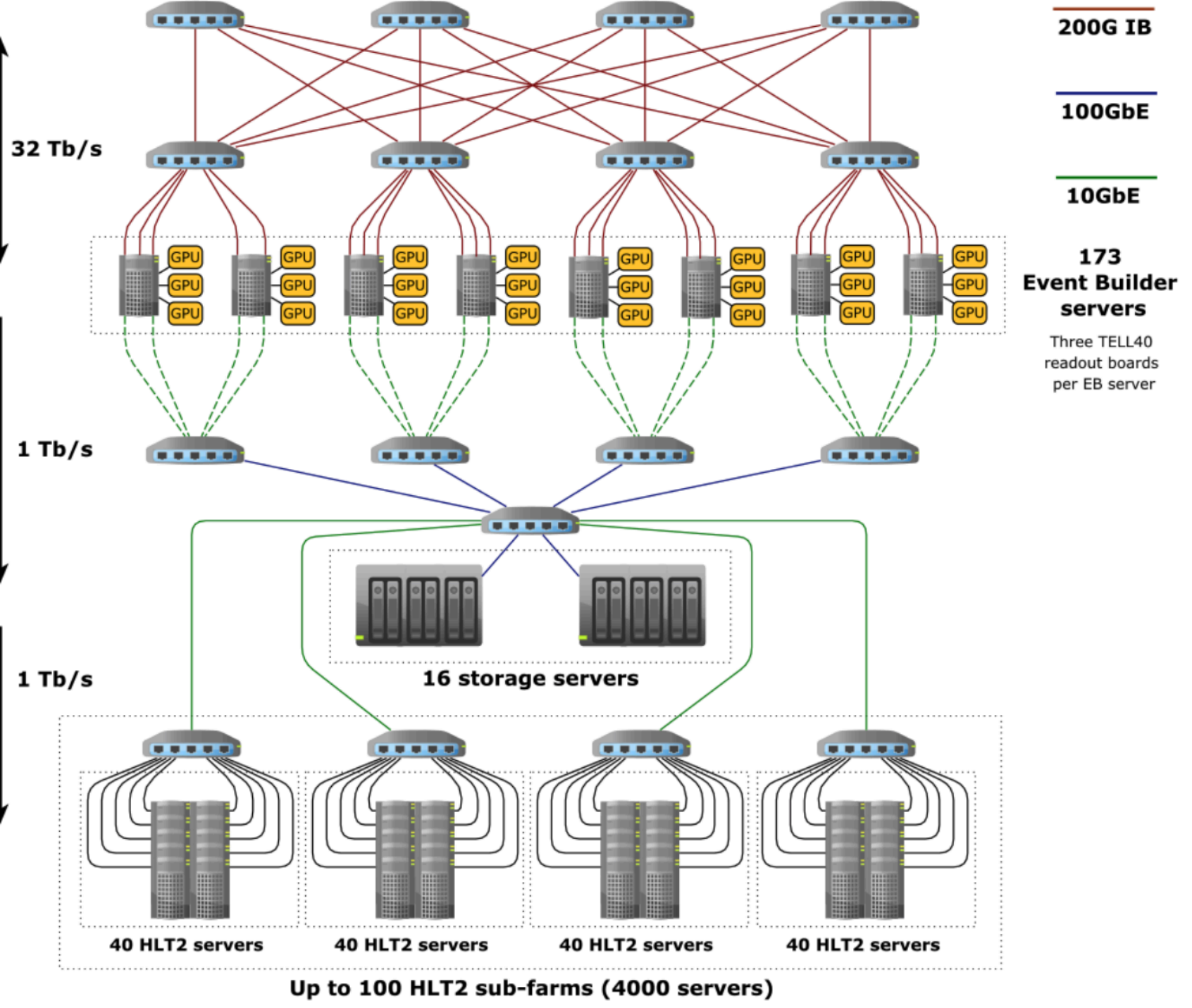
Engineers

N. Garroum (LPNHE), G. Grasseau (Subatech), D. Vintache (Subatech)

Synergies with detector projects in France

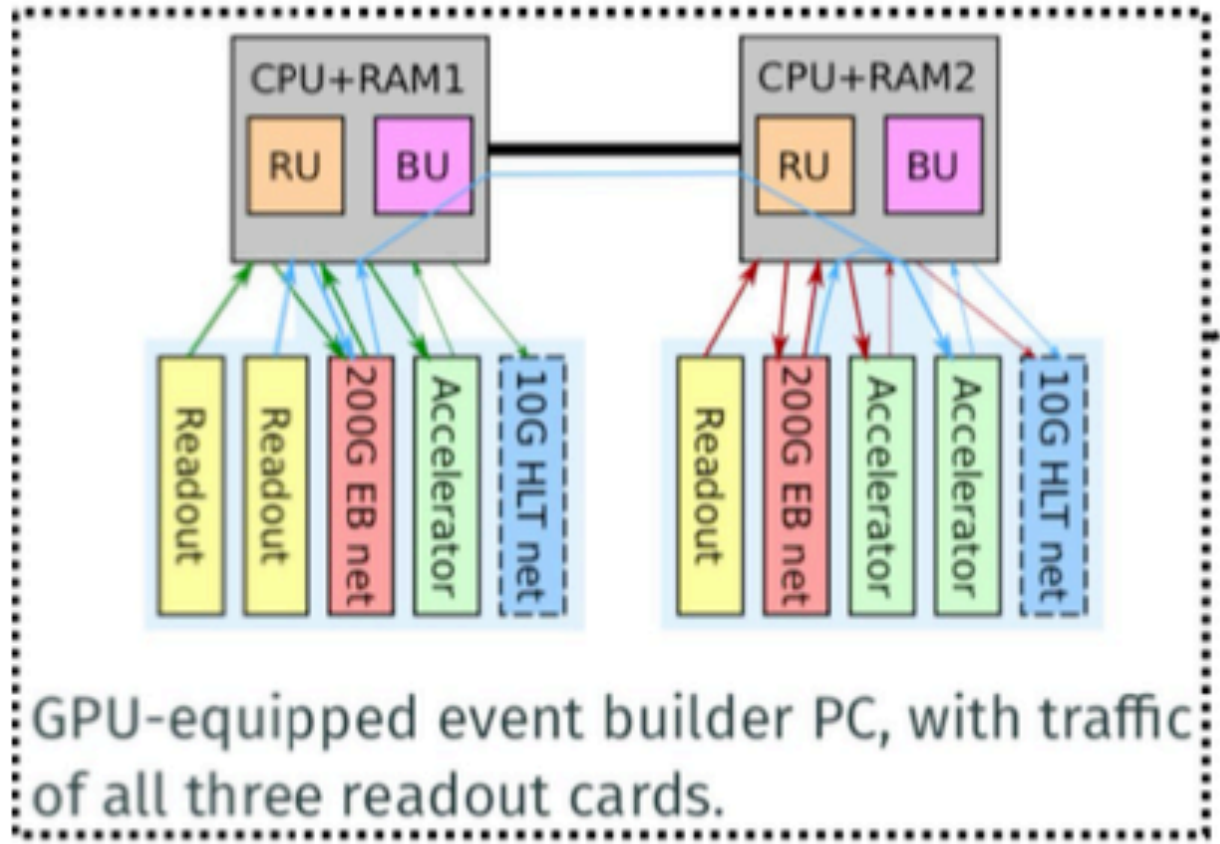
- Algorithm development closely linked to sub-detector projects
- RTA will focus on:
 - Designing reconstruction algorithms for parallel computing architectures within Allen
 - Designing global reconstruction using information from different subsystems
 - Performing alignment and calibration
- Within France, will focus algorithm developments on sub-detectors co-developed in France
 - Calorimeter and potentially the pixel tracker

Synergies with PCIe400 project in France



Comp. Soft for Big Science 6, 1 (2022)

- Real-time analysis should occur as early as possible in the readout chain
- Decoding of detector primitives are natural tasks to be carried out in the readout cards
 - Clustering of pixel detectors
 - Clustering and reconstruction of the calorimeter
- A strong collaboration between the RTA-France team and the PCIe400 group is foreseen for common developments



Synergies of the Allen framework with other experiments

Allen

- Experiment-independent
- Cross-architecture compilation, extensible to emerging architectures
- Parallelism through batched processing & scheduling
- Mixed classical & AI workflows
- Optimised for high-throughput parallel computing (microsecond execution)
- Energy efficiency included in performance metrics

*Suited for streaming
readout DAQ systems*

- Four IN2P3 labs are part of the European infrastructure project OSIDEE
 - Partners from academia and science, from SKA, LHCb and computer science
 - Focus on development of open AI techniques for real-time data processing
- Interest in Allen expressed from other experiments in the community
 - For example the ePIC experiment planned for EIC (data-taking starting in 2034)

Summary

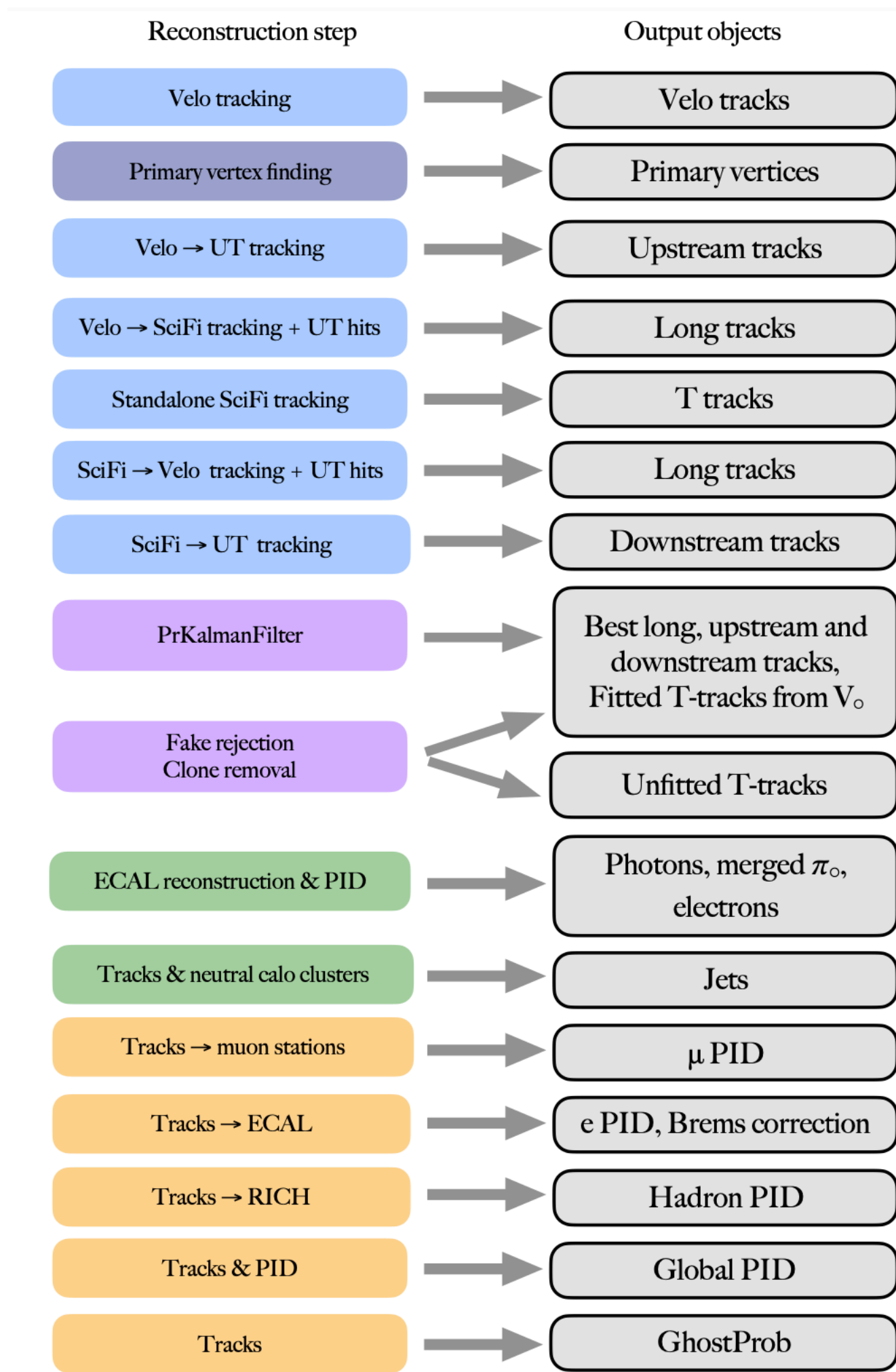
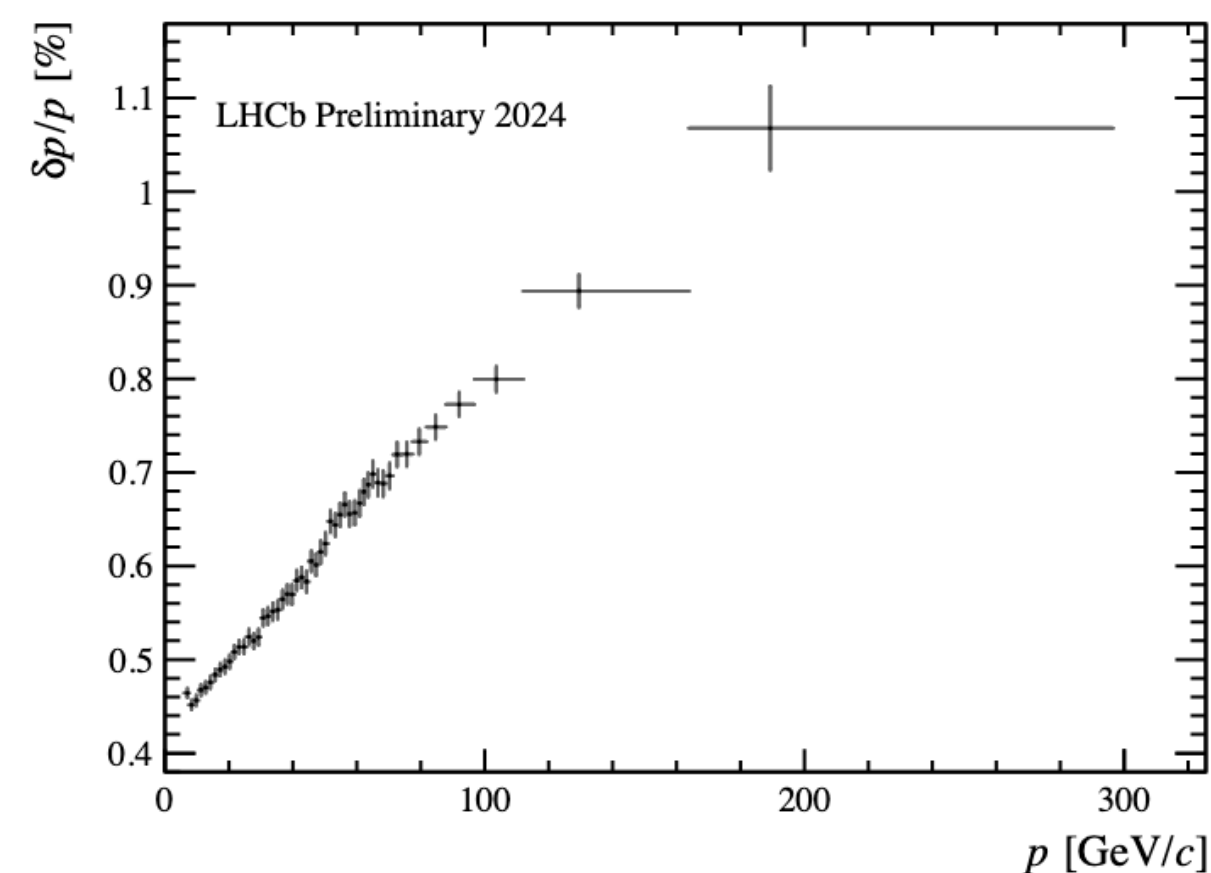
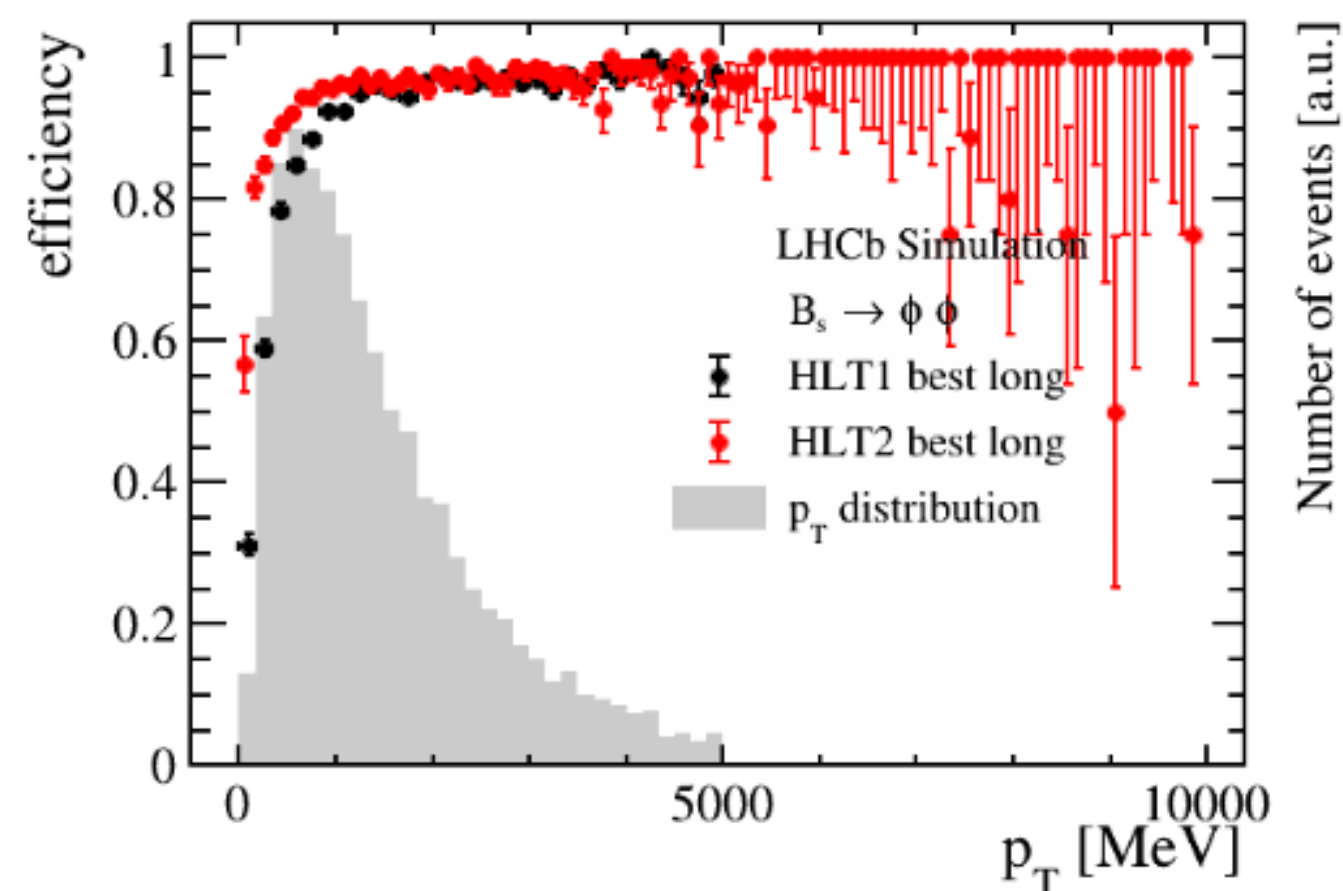
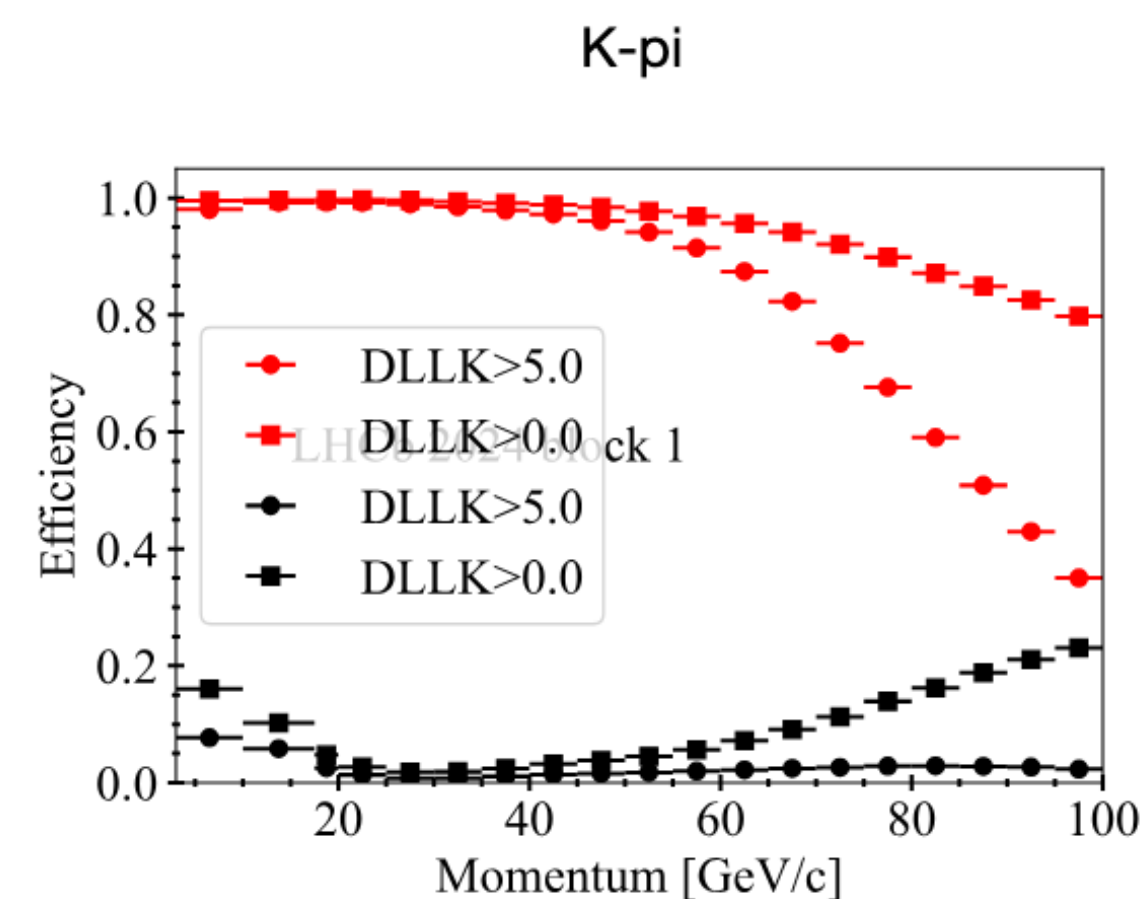
- LHCb has revolutionised its trigger system in Run 3 to use a fully software system based on GPUs and CPUs
 - LHCb France has played a key role in successfully delivering this system
- Allen framework has enabled HLT1 to go far beyond the original TDR design
 - Allen provides generic tools for high-performant processing on heterogeneous architectures
 - Future development will provide standalone framework to be used outside of LHCb
 - Core expertise and vision is located in France
- Proto-project for LHCb-RTA-U2-France has been created
 - Skills required for RTA correspond to critical software engineering skills at IN2P3
- LHCb France community proposes to keep leading role of RTA system for LHCb Upgrade II
 - Strong French implication in design of future system within LHCb
- For long-term support of the project, a coherent team of permanent physicists and engineers is crucial

Backup

HLT2 reconstruction

Offline quality reconstruction

- Track reconstruction algorithms tuned for efficiency
- Kalman filter with parameterised scattering errors
- Hadron and global PID



Allen: Memory manager & Multi-event scheduler

Multi-event scheduler

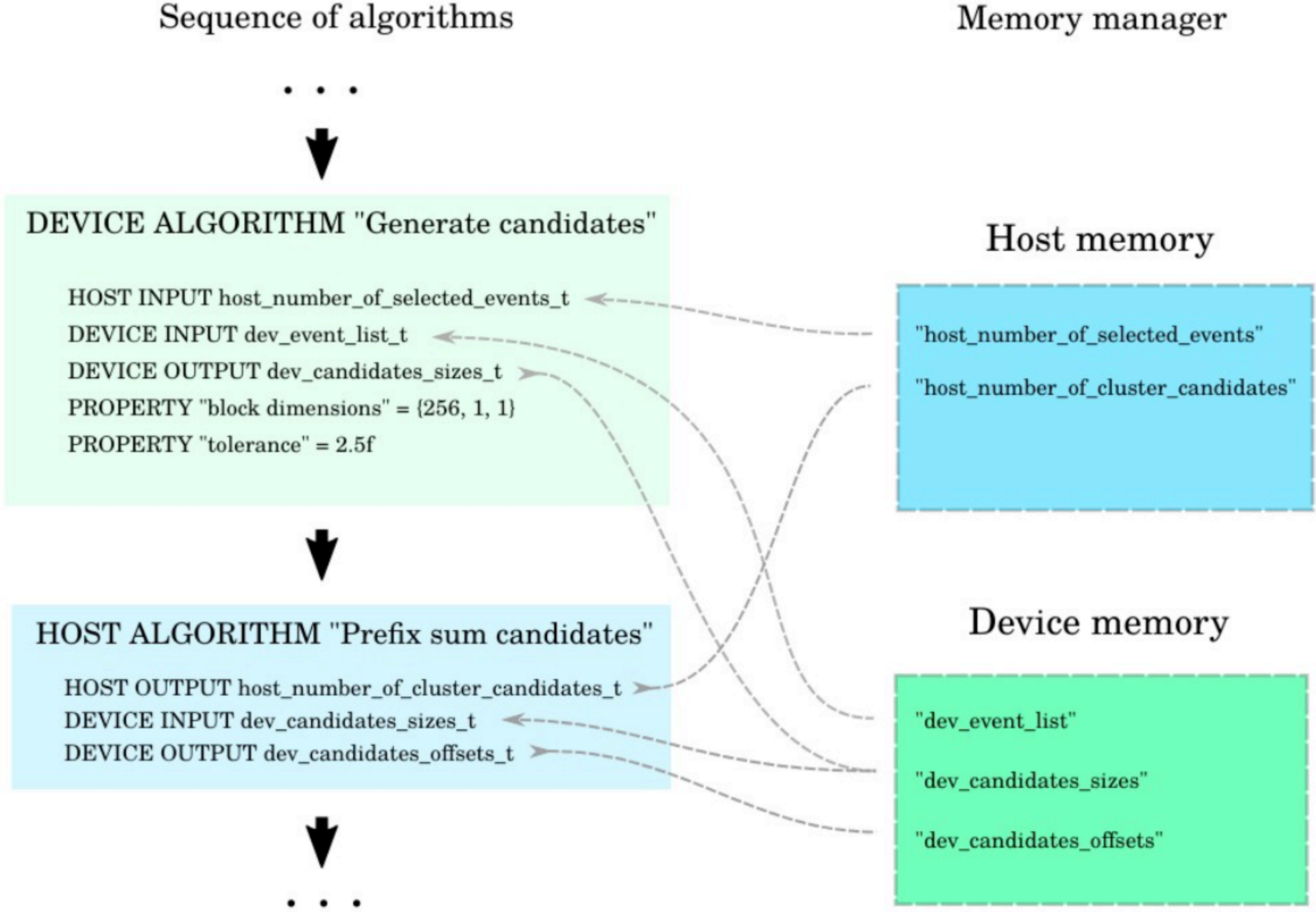
- For efficient GPU utilisation, every algorithm processes many events
- Scheduler generates static sequence of algorithms to be executed, considering all possible branching paths
- Event masks are generated storing outcome of algorithm for every event
- Masks are picked up by scheduler and required for the control flow

Memory manager

- Memory allocations on the GPU are slow
- Allocate chunk of memory at start of application
 - Bookkeeping of pointers during sequence processing
- Strong preference for « Count first, write later »
- Scheduler uses data dependencies to track lifetime of objects in memory
- Host and device memory managed analogously

Allen: Memory manager & Multi-event scheduler

- **Multi-event scheduler**
- For efficient GPU processes many events
- Scheduler generates sequence of candidates to be executed, compute paths
- Event masks are generated by algorithm for every event
- Masks are picked up by scheduler to control the control flow



GPU are slow
 start of application
 sequence processing
 read first, write later »
 mechanisms to track lifetime
 managed analogously

Allen: Python configuration

- Database of algorithms, inputs, outputs and properties built using code parsing with libclang
- Configure output (reconstructed) objects to be produced
- Configure algorithms with properties
- Multiple instances of an algorithm with separate inputs and outputs are possible
- Configuration in python using LHCb's [PyConf](#) package

```
seed_tracks = make_algorithm(  
    seed_confirmTracks_t,  
    name='seed_confirmTracks_{hash}',  
    host_number_of_events_t=number_of_events["host_number_of_events"],  
    dev_number_of_events_t=number_of_events["dev_number_of_events"],  
    dev_scifi_hits_t=decoded_scifi["dev_scifi_hits"],  
    dev_scifi_hit_count_t=decoded_scifi["dev_scifi_hit_offsets"],  
    dev_seeding_tracksXZ_t=xz_tracks["seed_xz_tracks"],  
    dev_seeding_number_of_tracksXZ_part0_t=xz_tracks[  
        "seed_xz_tracks_part0"],  
    dev_seeding_number_of_tracksXZ_part1_t=xz_tracks[  
        "seed_xz_tracks_part1"],  
    tuning_nhits=10,  
    tuning_tol_chi2=100,  
    tuning_tol=0.8,  
)
```

Hardware trigger saturation at LHCb for hadronic signatures

