



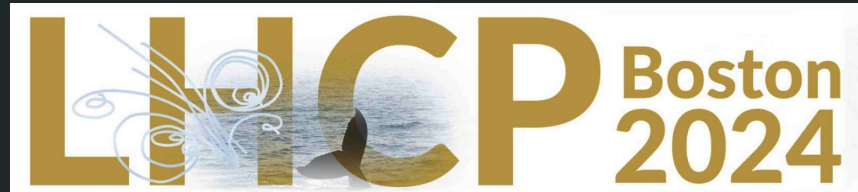
UiO • University of Oslo

# Software Upgrades for the High-Luminosity LHC

David Shope

On behalf of the ALICE, ATLAS, CMS, LHCb  
collaborations

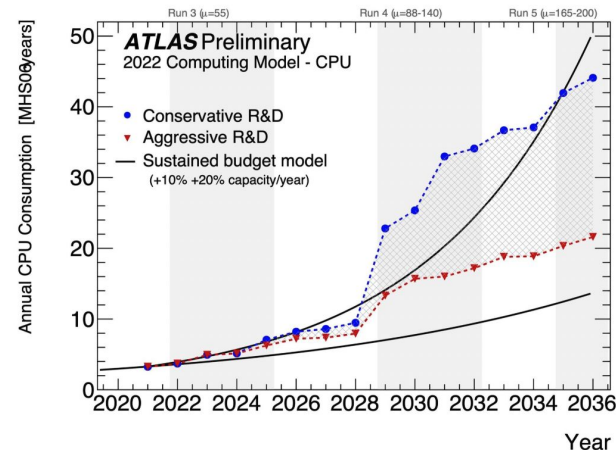
June 6, 2024



# Introduction

- HL-LHC data challenge:
  - Luminosity increase from  $2 \times 10^{34} \text{ s}^{-1} \text{ cm}^{-2}$  to  $7.5 \times 10^{34} \text{ s}^{-1} \text{ cm}^{-2}$ , with peak leveled pileup: 60 → 200 (ATLAS/CMS)
  - More frequent and larger events, with finer-grained detector readout
- Resources will fall short unless **significant** R&D occurs!
- Software upgrades for HL-LHC fall into several broad themes:

[ATLAS Software and Computing HL-LHC Roadmap (CERN-LHCC-2022-005)]



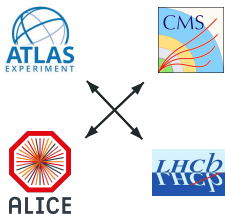
## Adapting to Heterogeneous Platforms



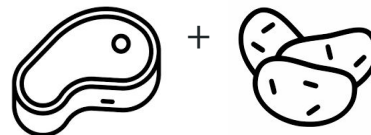
## New Approaches with Machine Learning



## Cross-experiment R&D Projects



## Meat-and-potatoes Developments



# A Brief Word On Porting To/Between GPU Devices



GPUs offer large parallel processing capabilities and excel at tasks such as the training of deep learning models - they've been used to great effect already in several key areas

However, the model for programming these devices **differs significantly** from classical x86 processors and supporting a heterogeneous set of accelerator technologies therefore also represents a **major** software development challenge



→ When porting to/between GPU devices, keep in mind:

- ◆ Not all applications are suitable for GPUs
- ◆ GPU languages are **evolving rapidly** with the hardware, difficult to predict which vendors will be most popular at the start of Run 4
  - ALICE: generic C++ code on GPUs in a vendor-agnostic approach (see also Gabriele Cimador's [parallel talk](#) tomorrow on performance in Run 3)
  - ATLAS: generic C++ code on GPUs, supporting multiple languages
  - CMS: ALPAKA for portability
  - LHCb: CUDA for NVIDIA GPUs in Run 3 trigger

# Software-Based Triggers



ALICE

- No trigger during Pb-Pb run since the start of Run 3
  - ◆ Continuous readout instead, with online compression of raw data in software
  - ◆ Designed to operate at 50 kHz for Pb-Pb run
  - ◆ Online computing farm consists of 350 servers, 8 GPUs each
    - 95% of online processing workload running on GPUs

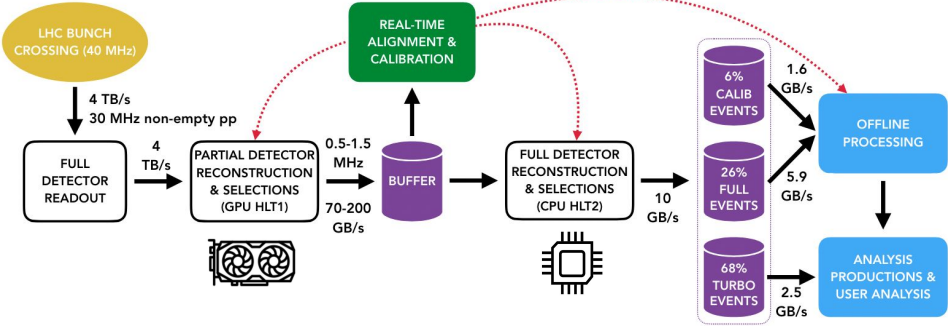


O<sup>2</sup> - common framework for **online/offline** data **reconstruction** and **analysis** [TDR]



- Fully software-based since the start of Run 3
  - ◆ Readout at full 30 MHz rate
  - ◆ Two-stage trigger:
    - HLT1, based on GPUs
      - align & calibrate in real time
      - partial reconstruction
    - HLT2, based on CPUs
      - full reconstruction
      - selection lines

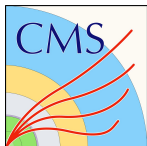
[The LHCb upgrade I]



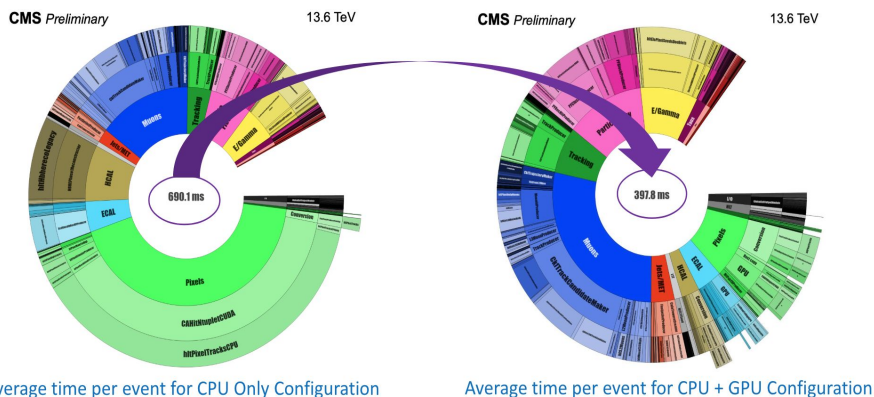
See also the parallel talk yesterday on [novel triggering strategies at the HL-LHC](#) by Marco Zanetti

# Software-Based Triggers

Run 4 trigger rate is  $\sim 3x$  that of Run 3 trigger, data size / event at least  $\sim 3x$  as well  
→ Both CMS and ATLAS TDAQ systems require **significant upgrades during LS3**



→ Using GPU enabled reconstruction in HLT since start of Run 3. Execution time / event **reduced by  $\sim 40\%$**



[G. Parida - Run-3 Commissioning of CMS Online HLT reconstruction using GPUs]



→ HLT runs on CPUs in Run 3

→ Overhaul to software level trigger for Phase II upgrade

- ◆ Event Filter (EF), data reduction: **1 MHz  $\rightarrow$  10 kHz** (100 kHz  $\rightarrow$  1 kHz for Run 3 HLT)
- ◆ Multiple types of computational units
  - Commodity CPU-servers
  - Possibly **accelerators**: GPU, FPGA
    - Decision for exact composition expected in 2025
- ◆ **EF tracking**: Many ongoing FPGA & GPU studies for track seeding from ITk inputs
- ◆ **EFCalo**: Studies with GPUs for topological clustering
- ◆ **EFMuon**: Muon reconstruction using ML algorithms

More info in parallel talk tomorrow on [Trigger performance \(including data scouting and GPU\) at CMS and ATLAS](#) by Silvio Donato

# Software-Based Triggers

Run 4 trigger rate is  $\sim 3x$  that of Run 3 trigger, data size / event at least  $\sim 3x$  as well  
→ Both CMS and ATLAS TDAQ systems require **significant upgrades during LS3**

## NextGen Triggers Project

A project created to support the **investigation of innovative computing technologies** (both hardware and software) in the design of **future data acquisition strategies** for HEP experiments based on the experience of ATLAS and CMS

- Funded by a not-for-profit foundation (*Eric and Wendy Schmidt Fund for the Strategic Innovation*)
- Approved by the CERN council in October 2023

### Work Packages:

- WP1 - Common infrastructure
- WP2 - ATLAS
- WP3 - CMS
- WP4 - Education and outreach

→ Supplements the existing Phase-II TDAQ upgrades with additional funding (48M USD over 5 years), enabling lines of research that are otherwise not feasible within existing financial and technological resources limits

Average time per event for CPU Only Configuration

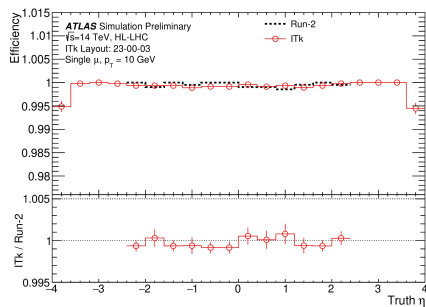
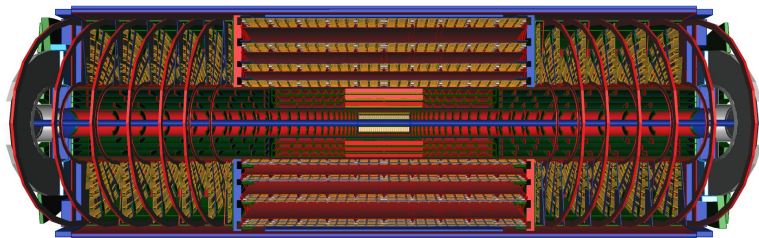
Average time per event for CPU + GPU Configuration

[G. Parida - Run-3 Commissioning of CMS Online HLT reconstruction using GPUs]

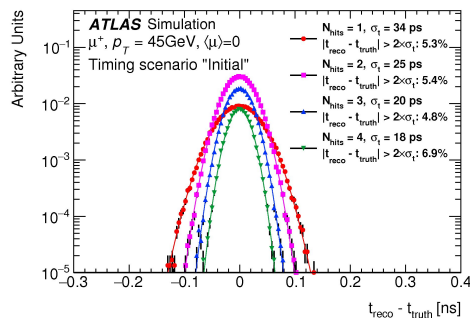
- ◆ **EFCalo**: Studies with GPUs for topological clustering
- ◆ **EFMuon**: Muon reconstruction using ML algorithms

More info in parallel talk tomorrow on [Trigger performance \(including data scouting and GPU\) at CMS and ATLAS](#) by Silvio Donato

# Meat-and-potatoes Development (ATLAS ITk+HGTD Case Study)



[\[Expected tracking and related performance with the updated ATLAS Inner Tracker layout at the High-Luminosity LHC \(ATL-PHYS-PUB-2021-024\)\]](#)



[\[ATLAS HGTD TDR \(CERN-LHCC-2020-007\)\]](#)

→ New detectors entail new software, much of which written from scratch - facilitated most often by experts which were around the last time a new detector was introduced (not a regular activity!)

- ◆ Detector description
- ◆ Digitization (electronics simulation, including extensive modeling of radiation damage)
- ◆ Reconstruction algorithms

→ ITk extends tracking to the forward region of ATLAS ( $|\eta| < 4.0$ ), where HGTD will also provide timing information ( $2.4 < |\eta| < 4.0$ )

- ◆ Requires updates to downstream software in other domains to properly utilize these forward tracks as well as the time

→ Critical that baseline software workflows are in place for **studying performance up to physics object level** to support decisions about upgrade projects (many of which are irreversible)

More on exploiting the time dimension in parallel session talks:

- [Detectors with timing capabilities](#) by Tim Evans
- [Pileup suppression with timing detectors](#) by Simone Pagan Griso

# Event Generators

- MC generators are projected to use **10-20% of the HL-LHC computing resources**
- ◆ Run 4 will see the need for both high-statistics inclusive samples as well as the efficient population of exclusive phase-spaces, all while maintaining the best available accuracy
  - ◆ Many ways to speed them up, see [review paper](#) from the HEP Software Foundation (HSF) Generator WG

**MadGraph4GPU:** speed up of matrix element calculation in MG5aMC on GPUs and vector CPUs

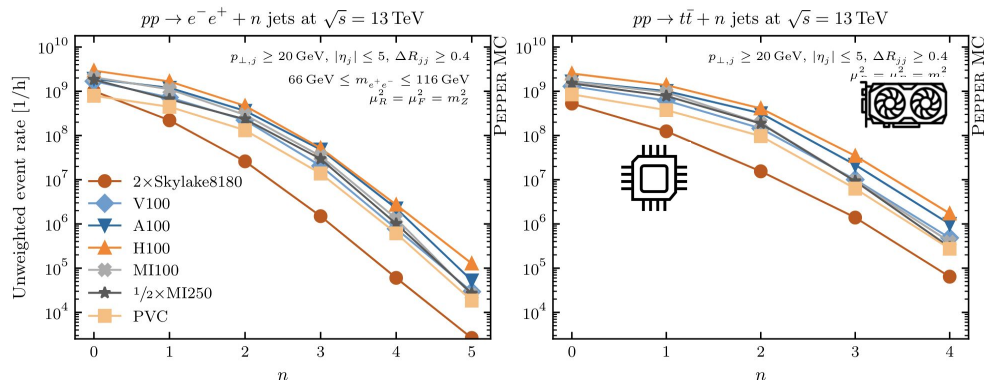
FORTRAN parts limit total achievable speed up (Amdahl's law)

Process	Madevent 262144 events			Standalone CUDA
	Total	Momenta+unweight	Matrix elm	ME Throughput
$e^+e^- \rightarrow \mu^+\mu^-$	17.9 s	10.2 s	7.8 s	$1.9 \times 10^6 \text{s}^{-1}$
	+CUDA Tesla A100	10.0 s	10.0 s	$633.8 \times 10^6 \text{s}^{-1}$
	1.8 x	1.0 x	390 x	334 x
$gg \rightarrow t\bar{t}gg$	209.3 s	7.8 s	201.5 s	$2.8 \times 10^3 \text{s}^{-1}$
	+CUDA Tesla A100	8.4 s	7.8 s	$758.9 \times 10^3 \text{s}^{-1}$
	24.9 x	1.0 x	336 x	271 x
$gg \rightarrow t\bar{t}ggg$	2507.6 s	12.2 s	2495.3 s	$1.1 \times 10^2 \text{s}^{-1}$
	+CUDA Tesla A100	30.6 s	14.1 s	$170.7 \times 10^2 \text{s}^{-1}$
	82.0 x	0.9 x	151 x	155 x

[S. Hageboeck - Madgraph5\_aMC@NLO on GPUs and vector CPUs: experience with the first alpha release]

**Pepper: Portable Engine for the Production of Parton-level Event Records**

Emphasis on portability, using [Kokkos](#)



[E. Bothmann - Pepper - A Portable Parton-Level Event Generator for the High-Luminosity LHC]



# R&D for Fast Simulation

Detector simulation is (today) the largest CPU consumer on the GRID, with time overwhelmingly spent in calorimeters

- Producing physics-accurate simulations in a fraction of the current time will be critical for HL-LHC programs
  - Traditionally, fast simulation methods have relied on parameterizations of the detector response using principal component analysis (PCA)
  - In recent years, generative ML techniques have shown significant promise as a replacement:

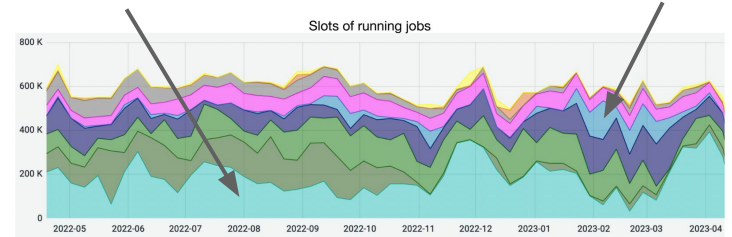
ALICE - **100x speed-up** of ZDC simulation wrt Geant4 using variational auto-encoders (VAEs) and generative adversarial networks (GANs)



[Machine Learning methods for simulating particle response in the Zero Degree Calorimeter at the ALICE experiment. CERN]

Full simulation

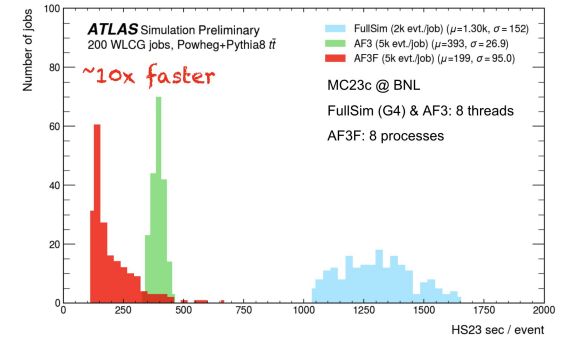
Fast simulation



[D. South - R&D in ATLAS Distributed Computing towards HL-LHC]

ATLAS - investigating fast tracking techniques (FATRAS) in addition to [FastCaloGAN](#) used in AtlFast3 (AF3):

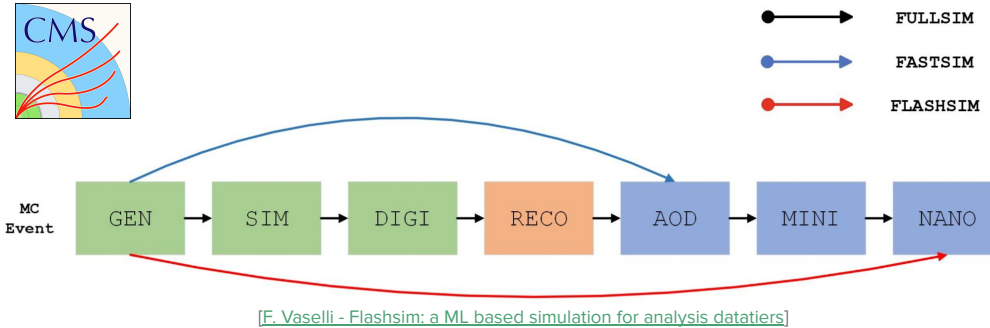
FATRAS+AF3 in ATLAS



[R. Wang - FATRAS integration for ATLAS fast simulation at HL-LHC]

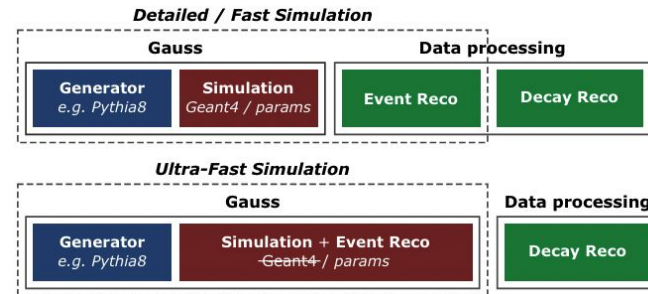
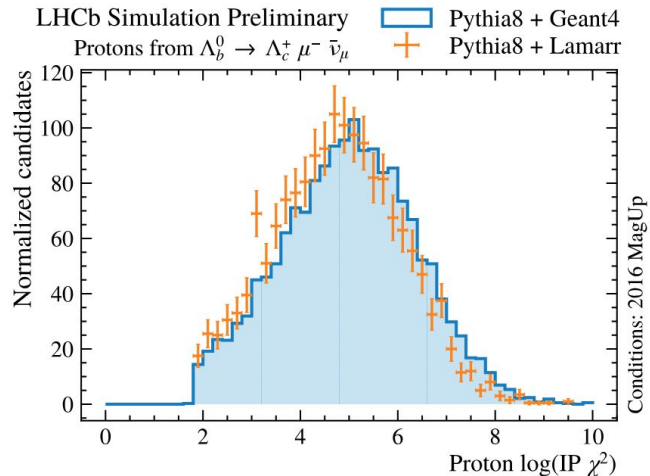
# R&D for Fast Simulation

See parallel talk tomorrow by Liza Mijovic on [Fast simulation with generative models at the LHC](#) for more details



Parameterization of the high level response of the detector can also enable a fast reconstruction

- CMS exploring **normalizing flows** (generative models for PDFs) to jump from generated events straight to NanoAOD
- LHCb developed **LAMARR**, a Gaudi-based framework for ultra-fast simulation
  - Currently GBDT for efficiency modeling, GAN for high-level physics distributions



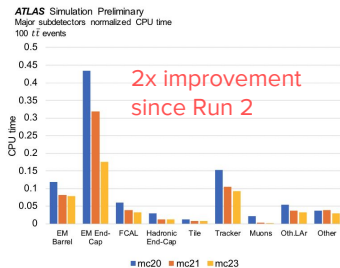
[Lamarr: LHCb ultra-fast simulation based on machine learning models deployed within Gauss]



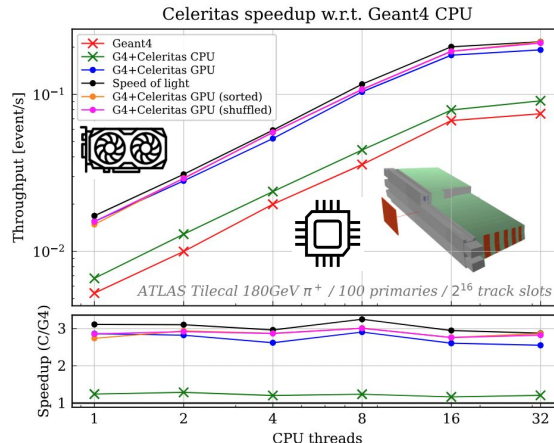
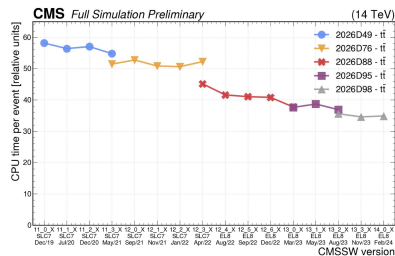
# R&D for Full Simulation

Two main avenues to reach HL-LHC demands:

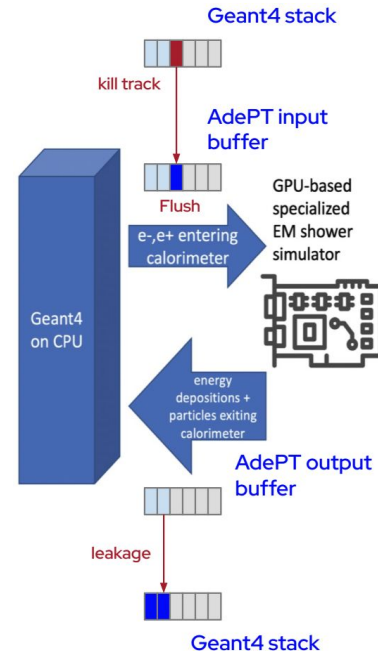
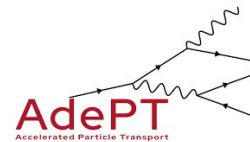
## A. Geant4 Optimization



[M. Schmidt - Optimizing the ATLAS Geant4 detector simulation]



[J. Esseiva - Celeritas: evaluating performance of HEP detector simulation on GPUs]



[A. Gheata - Accelerated demonstrator of electromagnetic Particle Transport (AdePT) status and plans]

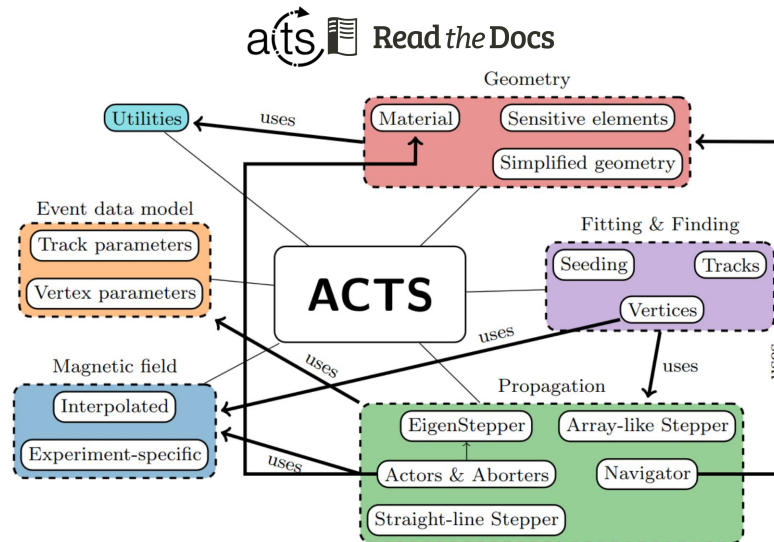
## B. GPU Utilization

- Two projects addressing the problem of general HEP particle transport on GPUs:
  - [AdePT](#) (CERN) - Accelerated demonstrator of electromagnetic Particle Transport
  - [Celeritas](#) (DOE)
- Use of GPUs to offload simulation of optical photons in LHCb: [Opticks](#) or [Mitsuba3](#)

# ACTS - A Cross-Experiment Tracking R&D Project

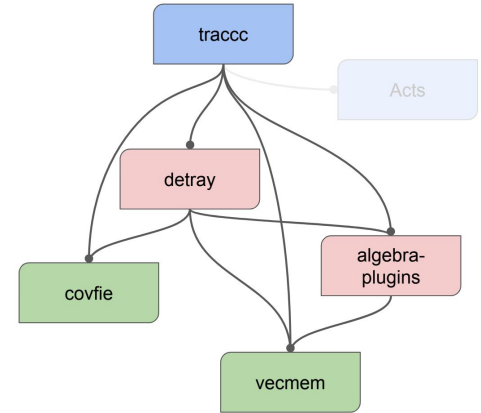
→ A Common Tracking Software (ACTS):

- ◆ Experiment-independent toolkit for the reconstruction of charged particle trajectories
- ◆ Core implementation based on modern (C++) and thread-safe code with a focus on maintainability
- ◆ Component library structure relying on minimal dependencies (CMake, Eigen, BOOST + optional plugins) that can be integrated into an experiment software
- ◆ Experiments with ACTS in use/data taking:  
ATLAS, FASER, sPHENIX (many other experiments in design/feasibility stage)



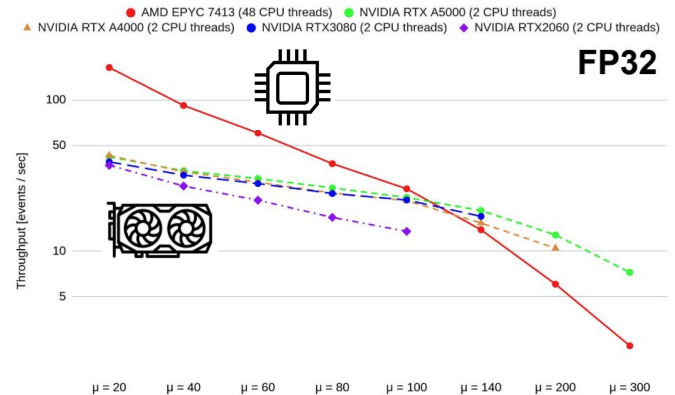
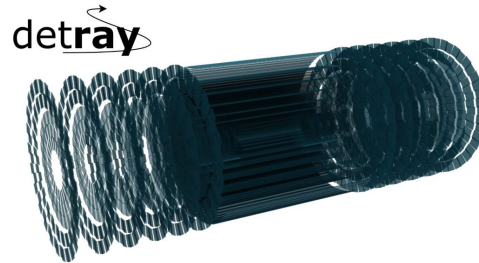
# ACTS on GPU - tracc

- Goal of tracc:
  - mapping of procedural (non-ML) tracking algorithms such as the Combinatorial Kalman Filter onto GPU **without algorithmic compromises**
- Full tracking chain (with exception of ambiguity resolution) now running on device with CUDA!
- Initial tests utilizing the OpenDataDetector (ODD) geometry show promising results, even with preliminary un-optimized algorithms



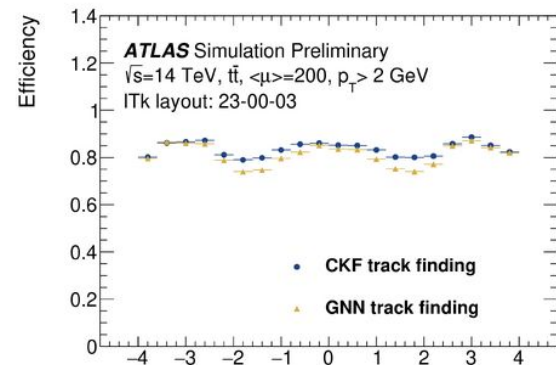
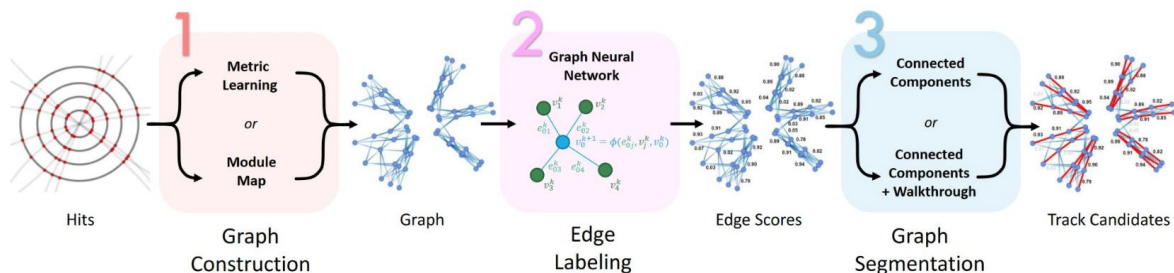
Category	Algorithms	CPU	CUDA	SYCL	Alpaka	Kokkos	Futhark
Clusterization	CCL / FastSv / etc.	✓	✓	✓	●	●	✓
	Measurement creation	✓	✓	✓	●	●	✓
Seeding	Spacepoint formation	✓	✓	✓	●	●	●
	Spacepoint binning	✓	✓	✓	✓	✓	●
	Seed finding	✓	✓	✓	✓	●	●
	Track param estimation	✓	✓	✓	✓	●	●
Track finding	Combinatorial KF	✓	✓	●	●	●	●
Track fitting	KF	✓	✓	✓	●	●	●
Ambiguity resolution	Greedy resolver	✓	●	●	●	●	●

✓: exists, ●: work started, ●: work not started yet

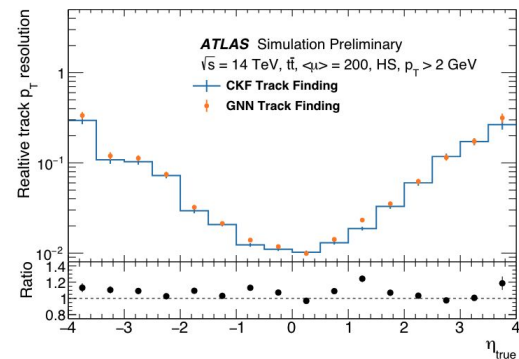


# Growth of ML Solutions - GNN4ITk

- Application of a graph neural network to the problem of track finding in ATLAS
  - Treat each hit as a node and each edge as a hypothesis for the two connected nodes to be consecutive hits on a particle track
- First look at the physics performance of the GNN-based pipeline compared with the Combinatorial Kalman Filter (CKF):
  - GNN provides competitive tracking efficiency (even in challenging dense environments), as well as a high quality track parameter resolution
- Work ongoing to assess computational performance (integration with ACTS and GPU), but preliminary results are promising



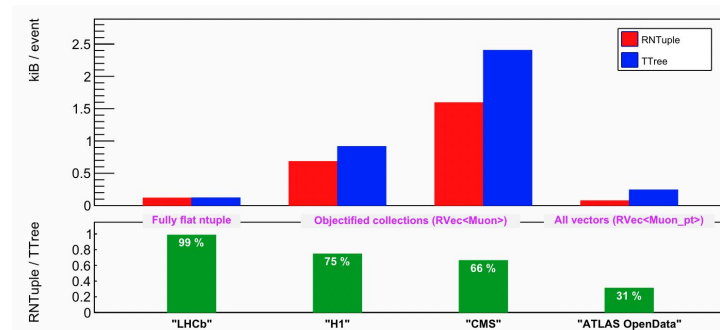
[[H. Torres - Physics Performance of the ATLAS GNN4ITk Track Reconstruction Chain](#)]



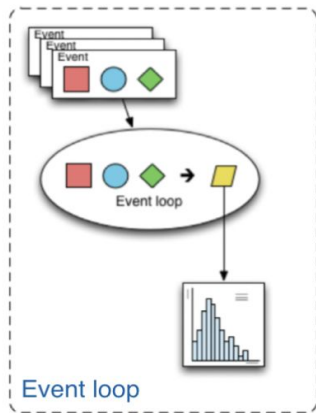
# Evolution of Data Analysis

Analysing the vast data of HL-LHC is itself a challenge as **current workflows do not scale**

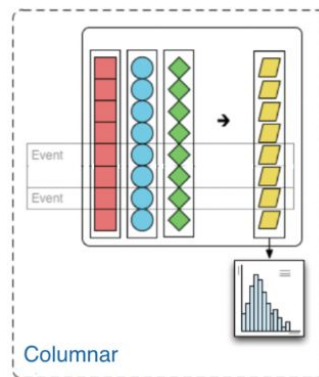
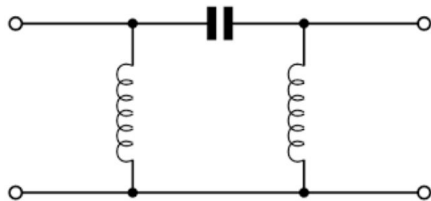
- Disk/Tape are limited/expensive and analysis data are projected to occupy ~30% (ATLAS model). Adaptation is needed.
  - ◆ Evolved encodings: ROOT's [RNTuple](#)
  - ◆ Reduce data copies and intermediate formats
  - ◆ More compact data formats (e.g. NanoAOD and [PHYSLITE](#))
- Target is < 10 kB per event for super-fast lightweight analysis



[J. Blomer - ROOT's RNTuple I/O Subsystem: The Path to Production]



Impedance mismatch between current analysis paradigm in HEP and tools now available in the wider software ecosystem



Clear trend towards **columnar analysis**, operating on arrays for a batch of events



Coffea

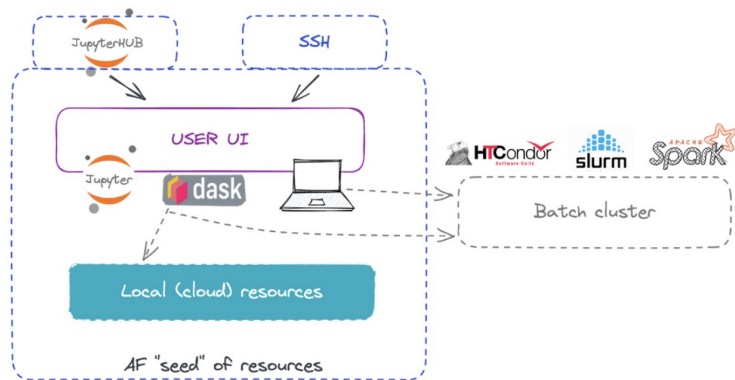


[L. Gray - Fine-Grained HEP Analysis Task Graph Optimization with Coffea and Dask]

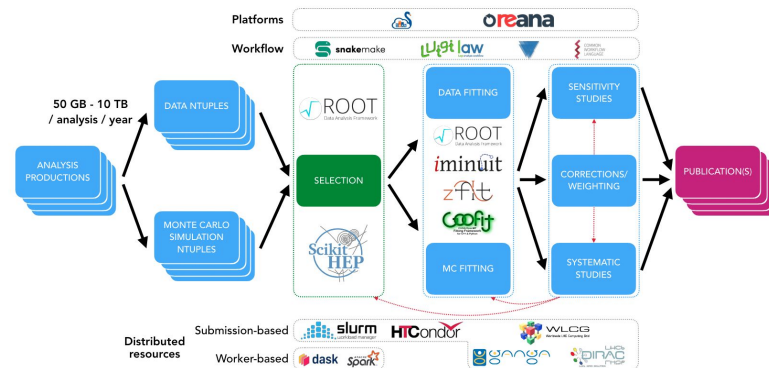
# Analysis Facilities

In the future, it may become possible to make analysis a highly interactive experience, minimizing the time to iterate between ideas and plots (see also Rob Gardner's talk on [Computing at the HL-LHC and beyond](#))

- **Key features** of an analysis facility:
  - ◆ Data consolidated at the site
  - ◆ Accessible to people from outside the site
  - ◆ Set up for both a distributed and interactive analysis style, with all tools readily available (including ML resources)
- Clearly a desire to make the Grid an easier tool to use for analysis, although **significant questions remain**:
  - ◆ What are the exact use cases, analysis model differences Run 3 to Run 4, organization, benchmarks, dedicated hardware needed



[[HSF Analysis Facilities white paper](#)]

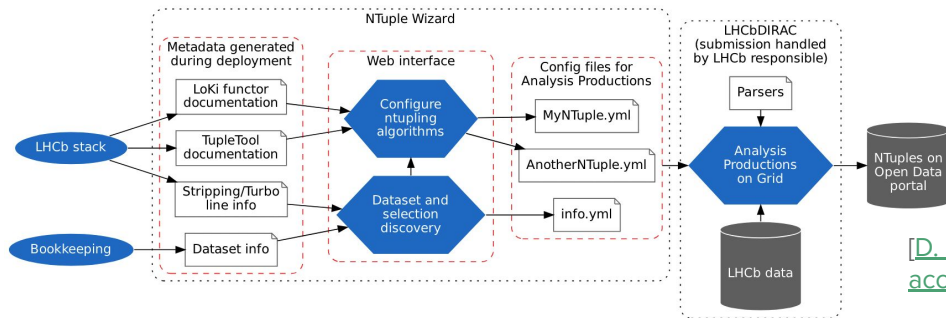


[[Dataflow diagrams for LHCb user analysis in Run 3](#)]



# Support for Legacy Data

- Complying with CERN Open Science policy [\[link\]](#)
  - ◆ Providing open data to the public [\[link\]](#)
    - Not only published results and outreach materials (usually in the form of highly processed ntuples), but reconstructed data as well (preprocessed to derive physics objects)
    - Takes real work for this to be highly useful to the community!
  - ◆ Analysis preservation through reproduction of workflows (e.g. [REANA](#))



LHCb Ntuple Wizard

[\[D. Fitzgerald - An Ntuple production service for accessing LHCb Open Data: the Ntuple Wizard\]](#)

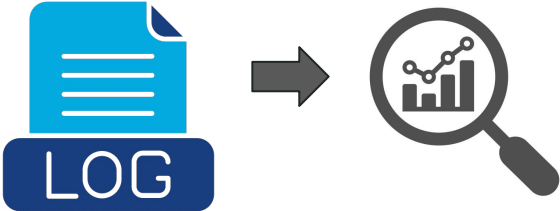
- Consideration must be given in addition to the feasibility to run old data/MC (from simulation through reconstruction) in the latest version of software (or rather, the implications of not being able to do so)
  - ◆ Real challenge posed if older data gets cut off
  - ◆ Alternatively, major investment likely needed to keep older data modern
    - One example (out of many): support for TRT in ACTS (ATLAS)
    - May be difficult to find developers willing to maintain detectors that are no longer taking data

# Leveraging Foundational Models for HEP


- LLMs fine-tuned for an experiment can summarize collaboration documentation that's disperse and difficult to navigate
  - TWikis, CDS, Indico, JIRA, GitLab, e-groups, etc.



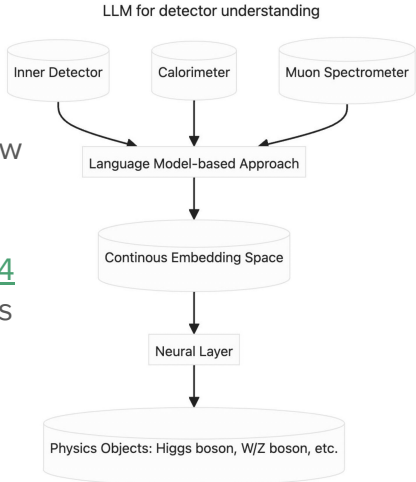
- Auto-parsing log files to provide useful data, e.g. summary of errors in human-readable form:



what is a foundational model in machine learning?

 A foundational model in machine learning refers to a **large, pre-trained model** that serves as a base for various downstream tasks. These models are typically trained on vast amounts of data across diverse domains and **can be fine-tuned for specific applications**. Foundational models include large language models (LLMs) like GPT-3 and GPT-4, as well as vision models like those based on convolutional neural networks (CNNs) or transformers.

- Foundational models have revolutionized the machine understanding of natural language, how about **scientific data**?
  - Various creative (potential) applications discussed during a [session in ACAT 2024](#)
  - Common observation: pre-training results in better performance of a downstream task than training that task from scratch



# Summary by ChatGPT-4o



Software Upgrades for HL-LHC (2)...

PDF

Please provide a summary of this presentation in just a few sentences



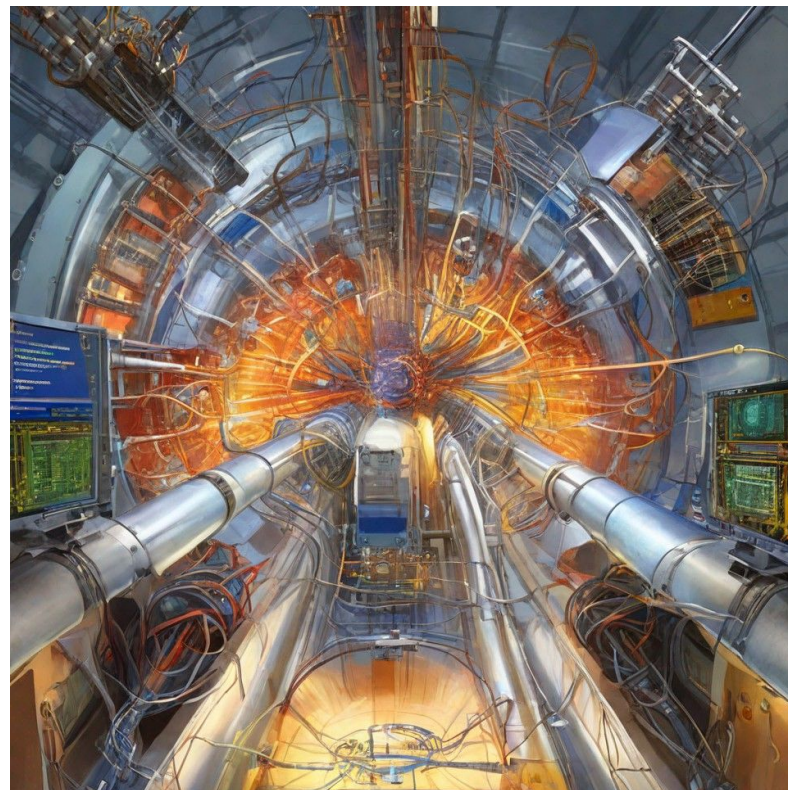
The presentation "Software Upgrades for the High-Luminosity LHC" discusses the various software development initiatives necessary to handle the increased data and complexity expected with the High-Luminosity Large Hadron Collider (HL-LHC). Key themes include adapting to heterogeneous computing platforms, incorporating machine learning, and extensive R&D to avoid resource shortfalls. The presentation highlights efforts in porting applications to GPU devices, developing software-based triggers, enhancing detector simulations, and optimizing tracking algorithms with tools like ACTS and tracc. Additionally, it emphasizes the need for evolving data analysis methods, supporting legacy data, and exploring foundational models in high-energy physics.



# Summary by a Human

- Meeting the challenges brought on by the HL-LHC requires not only hard work and innovation in hardware, but in software as well
  - ◆ Cross-experiment projects can help ease this burden!
- Critical that baseline upgrade-related software workflows (including meat-and-potatoes developments) are already in place today
  - ◆ Inform studies impacting the new detectors as (or before) they are built
- Many promising ideas for leveraging new technologies are actively being pursued:
  - ◆ Heterogeneous platforms
  - ◆ Machine learning applications

Thanks for your attention!



“Software Upgrades for the HL-LHC” - DALL-E 3

BACKUP SLIDES