# Deep-learning-based tracking demonstrator for the HL-LHC

Jan Stark

Laboratoire des 2 Infinis – Toulouse (L2IT)

CAF users meeting, November 21<sup>st</sup> 2023

### Approach of this talk

Will not give a comprehensive overview of techniques that have developed, nor of the striking the results that have been obtained. These are described in detail in other talks, e.g.:

- Charline Rougier at *Connecting the Dots 2022*, Princeton (clickable link)
- Xiangyang Yu at CHEP 2023, Norfolk (clickable link)
- Sylvain Caillou at CHEP 2023, Norfolk (clickable link)
- Heberth Torres at Connecting the Dots 2023, Toulouse (clickable link)

In instead, will try to describe our work in the bigger context.

### ATLAS HL-LHC S&C roadmap



# Predictions are difficult, especially when they concern the future (George Bernard Shaw, Winston Churchill or Niels Bohr)

#### What Future Processors Will Look Like

#### 2k f 72 X 286 in 1.1k <

AMD CTO Mark Papermaster talks about why heterogeneous architectures will be needed to achieve improvements in PPA.

JULY 13TH, 2022 - BY: ED SPERLING

Mark Papermaster, CTO at AMD, sat down with Semiconductor Engineering to talk about architectural changes that are required as the benefits of scaling decrease, including chiplets, new standards for heterogeneous integration, and different types of memory. What follows are excerpts of that conversation.

#### SE: What does a processor look like in five years? Is it a bunch of chips in a package? Is it a CPU and FPGA and GPU?

**Papermaster**: There is no question that the future of processing is heterogeneous. It's multiple compute engines working in tandem, because massive data and



graphics processing is needed everywhere. It's needed in data centers and in PCs, and the explosion of data from the Internet of Things requires analysis and visualization across that whole food chain. There clearly is a requirement of domain-specific architectures. The CPU is fantastic for general processing, and there are a gazillion applications

**~**1

that run on x86 and on Arm. For more specialized graphics and vector processing, FPGAs or ASICs can provide very specialized computing. We saw this future at AMD well over a decade ago, and we pointed our R&D efforts toward this future. We've been in mass production for years with an APU that combines the CPU and GPU for our client and embedded markets, and now we're bringing that APU into the CTO = Chief Technical Officer

#### PPA = Power, performance, area

I don't have a crystal ball either to predict the future. But I think that, as a field, we must be able to run our software on the GPU-heavy heterogeneous architectures that may well be the future.

INTEL RODUCTS SUPPORT SOLUTIONS DEVELOPERS PARTNERS FOUNDRY

#### Why Data Center GPUs Are Essential to Innovation

Data center graphics processing units (GPUs) are discrete accelerators that enable and enhance emerging technologies such as artificial intelligence (AI), rendering, analytics, and simulation/modeling.

#### GPUs vs. CPUs

Why Data Center GPUs? Use Cases Intel® Solutions

#### ₽ 🖂

Q Search Intel.com

To support AI, analytics, 3D rendering, and other advanced workloads, GPUs must play an expanded role in your data center environment. By augmenting CPUs with powerful parallel processing capabilities, data center GPUs help speed outcomes and accelerate innovation.

#### Key Takeaways

 Data center GPUs are used alongside CPUs to meet the elevated computational demands



#### CPUs are no longer the centre of the data centre



PUs) are shifting the balance of power in the data centre. ABI Research expects this arket to grow significantly, driven by the emergence of highly specialised workloads

| Contrast univer stray filts see VLCA limiting und<br>a specy choir<br>The changing face of retailing<br>to colopacity and the colopacity with<br>famic to i<br>Colud to bein prolite amotter Air powered digit<br>thanks for churchel's monter ways for<br>modernising under development.<br>Categories<br>ADRICULTURE<br>A IND MACHINE LEAREND<br>APPLICATION AM DEMOLEYARE   | R       | ecent Posts                                   |
|--|---------|---|
| Un support count of intelling<br>Trodgomotal resolutionises that delivery with<br>franct to T<br>Cloud to their prolline monther Air-powered digit<br>tractics for future-to Timoter ways for<br>moderning under delivery for<br>moderning under delivery of the<br>delivery of the second second second<br>additional second second second second second second<br>additional second second second second second second second<br>additional second second second second second second second second<br>additional second se | 6       | artner survey says 91% see VUCA limiting valu |
| Categories  | т       | he changing face of retailing                 |
| Could be help realits smarter Al-powered digit<br>twins<br>moderning urban development<br>Categories<br>Addicut TURE<br>A AND ALCHONE (LATERNO<br>APPLICATION AND MEDICIVARE   | Fi<br>S | oodpanda revolutionises food delivery with    |
| PodChots for FutureIoT: Smother ways for<br>modernising urban development<br>Cotegories<br>AGRECULTURE<br>AL AND MACHINE LEARNING<br>APPLICATION AND MIDDLEWARE  | C ty    | loud to help realise smarter Al-powered digit |
| moderniting urban development<br>Categories<br>AGRICULTURE<br>AI AND MACLINE LEARNING<br>APPLICATION<br>MERICATION   | P       | odChats for FutureIoT: Smarter ways for       |
| Cotegories<br>AGRICULTURE<br>AI AND MACHINE LEARNING<br>APPLICATION<br>APPLICATION AND MIDDLEWARE  | n       | nodernising urban development                 |
| AGRICULTURE<br>AI AND MACHINE LEARNING<br>APPLICATION<br>APPLICATION AND MIDDLEWARE  | c       | ategories                                     |
| AI AND MACHINE LEARNING<br>APPLICATION<br>APPLICATION AND MIDDLEWARE   | A       | GRICULTURE                                    |
| APPLICATION<br>APPLICATION AND MIDDLEWARE  | A       | I AND MACHINE LEARNING                        |
| APPLICATION AND MIDDLEWARE   | A       | PPLICATION                                    |
|  | A       | PPLICATION AND MIDDLEWARE                     |

### ATLAS HL-LHC S&C roadmap



### A few words on SIMD and GPUs



SIMD = single instruction, multiple data

CUDA C programming guide.

Figure taken from the

### A few words on SIMD and GPUs



### Consequence of SIMD: warp divergence in GPUs

if (threadIdx.x < 4) {
 A;
 B;
} else {
 X;
 Y;
}
Z;</pre>

X; Y; A; B; A; B; A; B;

Time

## Consequence of SIMD: warp divergence in GPUs



And this is just one "if". Imagine nested if statements ... the GPU quickly becomes idle.

Need to learn to design algorithms (almost) without *if* statements ... otherwise GPUs are useless.

Taken from the <u>CUDA C</u> programming guide.

### 5.1. Overall Performance Optimization Strategies

Performance optimization revolves around three basic strategies:

- Maximize parallel execution to achieve maximum utilization;
- Optimize memory usage to achieve maximum memory throughput;
- > Optimize instruction usage to achieve maximum instruction throughput.

Which strategies will yield the best performance gain for a particular portion of an application depends on the performance limiters for that portion; optimizing instruction usage of a kernel that is mostly limited by memory accesses will not yield any significant performance gain, for example. Optimization efforts should therefore be constantly directed by measuring and monitoring the performance limiters, for example using the CUDA profiler. Also, comparing the floating-point operation throughput or memory throughput - whichever makes more sense - of a particular kernel to the corresponding peak theoretical throughput of the device indicates how much room for improvement there is for the kernel.

### 5.2. Maximize Utilization

To maximize utilization the application should be structured in a way that it exposes as much parallelism as possible and efficiently maps this parallelism to the various components of the system to keep them busy most of the time.

Taken from the <u>CUDA C</u><u>programming guide</u>.

### 5.1. Overall Performance Optimization Strategies

Performance optimization revolves around three basic strategies:

- Maximize parallel execution to achieve maximum utilization;
- Optimize memory usage to achieve maximum memory throughput;
- > Optimize instruction usage to achieve maximum instruction throughput.

Which strategies will yield the best performance gain for a particular portion of an application depends on the performance limiters for that portion; optimizing instruction usage of a kernel that is mostly limited by memory accesses will not yield any significant performance gain, for example. Optimization efforts should therefore be constantly directed by measuring and monitoring the performance limiters, for example using the CUDA profiler. Also, comparing the floating-point operation throughput or memory throughput - whichever makes more sense - of a particular kernel to the corresponding peak theoretical throughput of the device indicates how much room for improvement there is for the kernel.

### 5.2. Maximize Utilization

To maximize utilization the application should be structured in a way that it <u>exposes as much</u> parallelism as possible and efficiently maps this parallelism to the various components of the system to keep them busy most of the time.

An obsolete/cheap gaming GPU like the RTX 2070 (an expensive A100 datacentre GPU) can run ~1k threads (~30k threads) in parallel. In addition, they have powerful mechanisms to switch from one thread to another to avoid waiting for things (e.g. memory access). ⇒ Need to break down our algorithms into tens to hundreds of thousands of simple sub-tasks that can be executed using the SIMD paradigm. VERY HARD in experimental particle physics.

Taken from the <u>CUDA C</u> programming guide.

### 5.1. Overall Performance Optimization Strategies

Performance optimization revolves around three basic strategies:

- Maximize parallel execution to achieve maximum utilization;
- Optimize memory usage to achieve maximum memory throughput;
- > Optimize instruction usage to achieve maximum instruction throughput.

Which strategies will yield the best performance gain for a particular portion of an application depends on the performance limiters for that portion; optimizing instruction usage of a kernel that is mostly limited by memory accesses will not yield any significant performance gain, for example. Optimization efforts should therefore be constantly directed by measuring and monitoring the performance limiters, for example using the CUDA profiler. Also, comparing the floating-point operation throughput or memory throughput - whichever makes more sense - of a particular kernel to the corresponding peak theoretical throughput of the device indicates how much room for improvement there is for the kernel.

### 5.2. Maximize Utilization

To maximize utilization the application should be structured in a way that it <u>exposes as much</u> parallelism as possible and efficiently maps this parallelism to the various components of the system to keep them busy most of the time.

Whenever we talk about GPU usage, we should quote numbers for GPU utilisation, memory bus efficiency, etc. This is the *nerf de la guerre* in GPU computing.

If people (like the traccc demonstrator) don't give the numbers, then ask for them.

An obsolete/cheap gaming GPU like the RTX 2070 (an expensive A100 datacentre GPU) can run ~1k threads (~30k threads) in parallel. In addition, they have powerful mechanisms to switch from one thread to another to avoid waiting for things (e.g. memory access). ⇒ Need to break down our algorithms into tens to hundreds of thousands of simple sub-tasks that can be executed using the SIMD paradigm. VERY HARD in experimental particle physics.

Example of an algorithm that runs efficiently on GPUs: matrix operations on large matrices (e.g. A\*B = C)



Neural networks can be expressed in terms of matrix operations.

Example of an algorithm that does not run efficiently on GPUs: Combinatorial Kalman filtering in a detector with a complex geometry and a non-zero magnetic field.



Example of an algorithm that runs efficiently on GPUs: matrix operations on large matrices (e.g.  $A^*B = C$ )



Neural networks can be expressed in terms of matrix operations.

Example of an algorithm that does not run efficiently on GPUs: Combinatorial Kalman filtering in a detector with a complex geometry and a non-zero magnetic field.



Simply knowing on which detector module to look for the next hit Is a quagmire of "if" statements. Then do this for the next ~15 layers in ITk ....

## Simple detectors, simple algorithms

Computing and Software for Big Science (2020) 4:7 https://doi.org/10.1007/s41781-020-00039-7

**ORIGINAL ARTICLE** 

Check for

#### Allen: A High-Level Trigger on GPUs for LHCb

R. Aaij<sup>1</sup> · J. Albrecht<sup>2</sup> · M. Belous<sup>3,4</sup> · P. Billoir<sup>5</sup> · T. Boettcher<sup>6</sup> · A. Brea Rodríguez<sup>7</sup> · D. vom Bruch<sup>5</sup> · D. D. H. Cámpora Pérez<sup>1,8</sup> · A. Casais Vidal<sup>7</sup> · D. C. Craik<sup>6</sup> · P. Fernandez Declara<sup>9,10</sup> · L. Funke<sup>2</sup> · V. V. Gligorov<sup>5</sup> · B. Jashal<sup>11</sup> · N. Kazeev<sup>3,4</sup> · D. Martínez Santos<sup>7</sup> · F. Pisani<sup>9,12,13</sup> · D. Pliushchenko<sup>4,14</sup> · S. Popov<sup>3,4,15</sup> · R. Quagliani<sup>5</sup> · M. Rangel<sup>16</sup> · F. Reiss<sup>5</sup> · C. Sánchez Mayordomo<sup>11</sup> · R. Schwemmer<sup>9</sup> · M. Sokoloff<sup>17</sup> · H. Stevens<sup>2</sup> · A. Ustyuzhanin<sup>3,4,15</sup> · X. Vilasis Cardona<sup>18</sup> · M. Williams<sup>6</sup>

Received: 18 December 2019 / Accepted: 3 April 2020 / Published online: 30 April 2020 © The Author(s) 2020

#### Abstract

We describe a fully GPU-based implementation of the first level trigger for the upgrade of the LHCb detector, due to start data taking in 2021. We demonstrate that our implementation, named Allen, can process the 40 Tbit/s data rate of the upgraded LHCb detector and perform a wide variety of pattern recognition tasks. These include finding the trajectories of charged particles, finding proton–proton collision points, identifying particles as hadrons or muons, and finding the displaced decay vertices of long-lived particles. We further demonstrate that Allen can be implemented in around 500 scientific or consumer GPU cards, that it is not I/O bound, and can be operated at the full LHC collision rate of 30 MHz. Allen is the first complete high-throughput GPU trigger proposed for a HEP experiment.

Keywords GPU · Real-time data selection · Trigger · LHCb



Fig. 4 Upgraded LHCb detector. The y-component of the magnetic field  $B_y$  is overlaid to visualize in which parts of the detector trajectories are bent. The maximum  $B_y$  value is 1.05 T

- LHCb raw events have an average size of 100 kB. When copying raw data to the GPU, the PCIe connection between the CPU and the GPU poses no limitation to the system, even when several thousand events are processed in parallel.

#### 🖂 R. Aaii

R. Aaıj raaij@cern.ch

D, vom Bruch

- dovombru@cern.ch D. H. Cámpora Pérez
- dcampora@cern.ch
- <sup>1</sup> Nikhef National Institute for Subatomic Physics, Amsterdam, The Netherlands
- <sup>2</sup> Fakultät Physik, Technische Universität Dortmund, Dortmund, Germany
- <sup>3</sup> National Research University Higher School of Economics, Moscow, Russia
- <sup>4</sup> Yandex School of Data Analysis, Moscow, Russia
- <sup>5</sup> LPNHE, Sorbonne Université, Paris Diderot Sorbonne Paris Cité, CNRS/IN2P3, Paris, France
- <sup>6</sup> Massachusetts Institute of Technology, Cambridge, USA
- <sup>7</sup> Instituto Galego de Física de Altas Enerxías (IGFAE), Universidade de Santiago de Compostela, Santiago de Compostela, Spain

- Faculty of Science and Engineering, Maastricht University, Maastricht, The Netherlands
- <sup>9</sup> European Organization for Nuclear Research (CERN), Geneva, Switzerland
- <sup>10</sup> Department of Computer Science and Engineering, University Carlos III of Madrid, Madrid, Spain
- <sup>11</sup> Instituto de Física Corpuscular, Centro Mixto Universidad de Valencia, CSIC, Valencia, Spain
- <sup>12</sup> INFN Sezione di Bologna, Bologna, Italy
- <sup>13</sup> Università di Bologna, Bologna, Italy
- <sup>14</sup> National Research University Higher School of Economics, Saint Petersburg, Russia
- <sup>15</sup> National University of Science and Technology MISIS, Moscow, Russia
- <sup>16</sup> Instituto de Física, Universidade Federal do Rio de Janeiro (UFRJ), Rio de Janeiro, Brazil
- <sup>17</sup> University of Cincinnati, Cincinnati, OH, USA
- <sup>18</sup> DS4DS, la Salle, Universitat Ramon Llull, Barcelona, Spain

🙆 Springer

### Simple detectors, simple algorithms



#### **Velo Detector**

The Velo detector consists of 26 planes of silicon pixel sensors placed around the interaction region. Its main purpose lies in reconstructing the pp collisions (primary vertices or PVs) and in creating seed tracks to be further propagated through the other LHCb detectors. The Velo track reconstruction is fully described in an earlier publication [17] and is recapped here for convenience.

The reconstruction begins by grouping measurements caused by the passage of a particle within each silicon plane into clusters, an example of a more general process known as connected component labeling. Allen uses a clustering algorithm employing bit masks, which searches for clusters locally in small regions. Every region can be treated independently, allowing for parallel processing.

Straight-line tracks are reconstructed by first forming seeds of three hits from consecutive layers ("triplets"), and then extending these to the other layers in parallel. We exploit the fact that prompt particles produced in ppcollisions traverse the detector in lines of constant  $\phi$  angle (within a cylindrical coordinate system where the cylinder axis coincides with the LHC beamline) and sort hits on every layer by  $\phi$  for fast look-up when combining hits to tracks.

# On fancy detectors



In this region, it is relatively clear in which direction to look for the next hit.

# Graph Neural Networks (GNN)

Data from HEP detectors in general, and tracking detectors in particular, are sparse.

> ATLAS ITk tracker at HL-LHC: 9 billion channels "only" 300k hits in one given event

The detectors are inhomogeneous (combine different technologies) and have complex geometry.

Such data are hard to represent as images.

Graphs are a natural tools to represent such data. GNNs are neural networks that operate on graphs of any topology and complexity.



#### ARTIFICIAL INTELLIGENCE



## Tracking based on GNNs



Represent the data using a graph



One node of the graph = one hit in the detector

Connect two nodes using an edge if "it seems possible" that the two hits are two (consecutive) hits on a track Goal: classify the edges of the graph



High classification score

=> high probability that the edge is part of a track

Low classification score => low probability that the edge is part of a track

### Tracking based on GNNs



### **GNN4ITk** Edge labeling with GNN



### GNN config:

- 2 layers per MLP
- 128D latent space
- 8 message-passing

#### New w.r.t. CTD 2022:

- Non-recurrent interaction network
- Doing batch norm
- Heterogeneous data





### Efficiency and purity vs (r, z)

**CTD2022** 





sylvain.caillou@l2it.in2p3.fr | CHEP2023 | Novel fully-heterogeneous GNN designs for track reconstruction at the HL-LHC | 9/13



- Problem addressed by passing info of the two individual strip clusters to the GNN; node features: •
  - Strip barrel:  $r_{\rm hit}, ..., r_{\rm cl1}, ..., r_{\rm cl2}, ...$ \_
    - Pixel:  $r_{\rm hit}, ..., r_{\rm hit}, ..., r_{\rm hit}, ...$

(Heterogeneous data format)

Other alternatives under study: hand-engineered edge features based on hit pair info, & heterogeneous GNN model •

\_

Heberth



## **Tracking efficiency**

• Competitive "physics" efficiency (excluding electrons)





### **Tracking efficiency** Tracking inside jets



• Competitive "physics" efficiency even in dense environment (excluding electrons)



Slide from Heberth at CTD2023

### Impact parameter resolution

• Given the good pixel hit content, good impact parameter resolution





#### 8<sup>th</sup> International CTD workshop **Université Paul Sabatier, Toulouse, France**

https://indico.cern.ch/e/CTD2023 ctd2023-loc@l2it.in2p3.fr

satellite event on Real time Tracking: triggering events with tracks (October 13th)



Local Organizing Committee

International Advisory Committee

Alberto Annovi (INFN Pisa) Paolo Calafiura (LBNL) Giuseppe Cerati (FNAL) Michel De Cian (EPFL) Matthias Danninger (SFU) Markus Elsing (CERN)

cnr

de physique nucléaire

et de physique des particules

Frank Gaede (DESY) Jose E. Garcia (IFIC Valencia) Maurice Garcia-Sciveres (LBNL) Vladimir Gligorov (LPNHE) Heather Gray (UC Berkeley/LBNL) Phil Harris (MIT)

David Lange (Princeton) Salvador Marti (IFIC Valencia) Fabrizio Palla (INFN Pisa) David Rousseau (IJCLab) Andi Salzburger (CERN) Louise Skinnari (Northeastern U.)



T Université

"Ç,

de Toulouse

## Talk at CTD2023 from LHCb/LPNHE

#### **GNN-based pipeline for track finding in the Velo at LHCb**

| Allen: current HLT1,        |
|-----------------------------|
| classical algorithms on GPL |

#### Ext4velo:

ML algorithm (GNN-based) on GPU

| 3. Track-Finding Performance <sup>3</sup>                                     |                |        |   |   |   |  |  |  |
|---|----------------|--------|---|---|---|--|--|--|
| Category  | Metric         | Allen  | $s_{\text{triplet}} > 0.32$<br>Etx4velo<br>$d_{\text{max}}^2 = 0.010$ | $s_{\text{triplet}} > 0.36$<br>Etx4velo<br>$d_{\text{max}}^2 = 0.020$ | Evaluation with 5,000 events                                    |  |  |  |
| Long, no electrons  | Efficiency     | 99.26% | 99.28%  | 99.51%  | • Track matched to a  |  |  |  |
| <ul> <li>In acceptance</li> <li>Reconstructible in the velo</li> </ul>        | Clone rate     | 2.54%  | 0.96%   | 0.89%   | particle if at least 70% of its<br>bits belong to this particle |  |  |  |
| <ul> <li>✓ Reconstructible in the SciFi</li> <li>✓ Not an electron</li> </ul> | Hit efficiency | 96.46% | 98.73%  | 98.90%  | This belong to this particle                                    |  |  |  |
|   | Hit Purity     | 99.78% | 99.94%  | 99.94%  | Allen algorithm described in                                    |  |  |  |
| Long electrons  | Efficiency     | 97.11% | 98.80%  | 99.22%  | <u>arxiv.2207.03930v2</u>                                       |  |  |  |
| <ul> <li>In acceptance</li> <li>Reconstructible in the velo</li> </ul>        | Clone rate     | 4,25%  | 7.42%   | 7.31%   | • 2 different GNN trainings for                                 |  |  |  |
| <ul> <li>✓ Reconstructible in the SciFi</li> <li>✓ Electron</li> </ul>        | Hit efficiency | 95.24% | 96.54%  | 96.79%  | $a_{\rm max} = 0.010$ and $a_{\rm max} = 0.020$                 |  |  |  |
|   | Hit purity     | 97.11% | 98.46%  | 98.46%  |   |  |  |  |
| Long, from strange  | Efficiency     | 97.69% | 97.50%  | 98.06%  | Long categories   |  |  |  |
| <ul> <li>In acceptance</li> <li>Reconstructible in the velo</li> </ul>        | Clone rate     | 2.50%  | 0.92%   | 0.81%   |   |  |  |  |
| <ul> <li>✓ Decays from a strange</li> <li>Good proxy for displaced</li> </ul> | Hit efficiency | 97.69% | 98.22%  | 98.77%  |   |  |  |  |
| tracks  | Hit purity     | 99.34% | 99.68%  | 99.68%  | Worse Better  |  |  |  |
| X   | Ghost rate     | 2.18%  | 0.76%   | 0.81%   |   |  |  |  |

39

(link to talk)

# **Tracking based on GNNs**

GNN4ITK demo plans ☆ & & ↔ Fichier Édition Affichage Insertion Format Données Outils Extensions Aide ⊞

Q Menus 5 순 등 및 100% ▼ \$ % 0 0 123 Pardé... ▼ - 10 + B I 중 A 🌺 ⊞ 용 ♥ 트 ★ ♥ ♥ ★ ♥ ♥ A ♥ G 🗉

Recently detailed our schedule toward the TDR in Q3 of 2024.

We have the physics performance that we need for a demonstrator.

Moving a lot of focus on implementation aspects.

- full chain (from clusters to fitted tracks) on GPU, without any intermediate transfers to/from the host
- use "GPU EDM" and track fit from ACTS/traccc

This domain (implementation) is new for many in our field (in ATLAS and elsewhere). It is crucial that we coordinate our activities with the LHCb colleagues at LPNHE, our ACTS/ATLAS colleagues at IJCLab and the Reprises project.

| A1 | -   | fx ID                                |  |                        |                     |                   |
|----|-----|--------------------------------------|--|------------------------|---------------------|-------------------|
|    | Α   | В                                    | С  | D                      | E                   | F                 |
| 1  | ID  | Goal                                 | Task   | Effort (person*months) | Start               | End               |
| 2  | 1   | Physics Performance Paper            |  | 2                      | 1 October 24, 2023  | March 24, 2024    |
| 3  | 1.1 |                                      | Request new (ITk 3.0.0) samples                    | :                      | 2 October 24, 2023  | December 24, 202  |
| 4  | 1.2 |                                      | Validate new samples                               | :                      | 2 December 24, 2023 | January 24, 2024  |
| 5  | 1.3 |                                      | Single-particle performance                        | 2                      | 2 October 24, 2023  | December 24, 2023 |
| 6  | 1.4 |                                      | Track-by-track CKF Comparison                      | :                      | October 24, 2023    | January 24, 2024  |
| 7  | 1.5 |                                      | Resolution loss from singlet clusters              | :                      | 2 January 24, 2024  | March 24, 2024    |
| 8  | 1.6 |                                      | Robustness studies                                 | 4                      | November 15, 2023   | March 15, 2024    |
| 9  | 1.7 |                                      | Dense environment tracking                         | :                      | B December 15, 2023 | March 15, 2024    |
| 10 | 1.8 |                                      | Large radius tracking                              | :                      | B December 15, 2023 | March 15, 2024    |
| 11 | 1.9 |                                      | Inference on other physics processes               | :                      | 3 January 15, 2024  | April 14, 2024    |
| 12 |     |                                      |  |                        |                     |                   |
| 13 | 2   | Compute Optimization & Compute Paper |  | 20                     | October 24, 2023    | May 15, 2024      |
| 14 | 2.1 |                                      | Perform initial timing study                       |                        | 2 October 24, 2023  | December 15, 2023 |
| 15 | 2.2 |                                      | Perform initial memory study                       |                        | 2 October 24, 2023  | December 15, 2023 |
| 16 | 2.3 |                                      | Regional tracking study                            | :                      | 3 October 24, 2023  | January 24, 2024  |
| 17 | 2.4 |                                      | Quantization, Pruning, Distillation study          | 4                      | October 24, 2023    | February 24, 2024 |
| 18 | 2.5 |                                      | Module map GPU-ification                           | :                      | October 24, 2023    | January 24, 2024  |
| 19 | 2.6 |                                      | NN (GNN, metric learning, etc) optimization        | :                      | B February 15, 2024 | May 15, 2024      |
| 20 | 2.7 |                                      | Custom/fused kernels                               | :                      | B February 15, 2024 | May 15, 2024      |
| 21 |     |                                      |  |                        |                     |                   |
| 22 |     |                                      |  |                        |                     |                   |
| 23 | 3   | Integrations & Infrastructure        |  | 10                     | October 24, 2023    | February 13, 2024 |
| 24 | 3.1 |                                      | Release CTD version of acorn, with trained models  | · · · · · ·            | October 24, 2023    | November 14, 2023 |
| 25 | 3.2 |                                      | Update acorn to latest Pytorch, Lightning, etc.    |                        | November 14, 2023   | December 14, 2023 |
| 26 | 3.3 |                                      | Add GPU ability to Docker Gitlab runner            | 0.9                    | 5 October 24, 2023  | November 7, 2023  |
| 27 | 3.2 |                                      | Recipe/walkthrough on dumping objects              | 0.25                   | 5 October 31, 2023  | October 31, 2023  |
| 28 | 3.2 |                                      | Recipe/walkthrough on training pipeline            | 0.25                   | 5 November 7, 2023  | November 7, 2023  |
| 29 | 3.3 |                                      | Recipe/walkthrough on configuring GNN              | 0.25                   | 5 November 14, 2023 | November 14, 2023 |
| 30 | 3.3 |                                      | Recipe/walkthrough on running Athena + IDPVM       | 0.25                   | 5 November 21, 2023 | November 21, 2023 |
| 31 | 3.5 |                                      | Simple inference script for physics analysis       | 0.5                    | 5 October 24, 2023  | November 7, 2023  |
| 32 | 3.6 |                                      | Move to athena 24, ITk layout 3.0.0                | 2                      | 2 November 12, 2023 | December 12, 2023 |
| 33 | 3.7 |                                      | Create the Dump module in main branch> ROOT file   |                        | 2 November 12, 2023 | December 12, 2023 |
| 34 | 3.8 |                                      | Full chain test in athena                          |                        | 2 January 1, 2024   | February 13, 2024 |
| 35 | 3.9 |                                      | Onnx conversion of all models                      |                        | 2 November 14, 2023 | January 14, 2024  |
| 36 |     |                                      |  |                        |                     |                   |
| 37 | 4   | ML R&D                               |  | 1:                     | 2 October 24, 2023  | May 15, 2024      |
| 38 | 4.1 |                                      | Single-particle generalization                     | 2                      | 2 October 24, 2023  | December 24, 2023 |
| 39 | 4.2 |                                      | Singlet cluster model & training                   | 2                      | 2 October 24, 2023  | December 24, 2023 |
| 40 | 4.3 |                                      | Electron-targeted model & training                 | 2                      | 2 January 15, 2024  | March 15, 2024    |
| 41 | 4.4 |                                      | Low-pt model & training                            | 2                      | 2 March 15, 2024    | May 15, 2024      |
| 42 | 4.5 |                                      | Improved metric learning graph construction purity | 2                      | 2 October 24, 2023  | December 24, 2023 |
| 43 | 4.6 |                                      | Fine-tuned models for different channels           |                        | 2 December 15, 2023 | February 15, 2024 |

# Additional material

### AISSAI conference on heterogeneous data in Toulouse – spring 2024



All nodes are pink, regardless of their position in the detector (tracker, calorimeter, muon detectors) ! Not only in this illustration, this reflects what is typically done in practice.

Ideally, nodes in different subdetectors would represent different types of measurements in different subdetectors (3D hits in the tracking pixel detectors, 2D measurements in the tracking strip detectors, energy deposits in the calorimeters, ...) and they would be shown in different colours in illustrative figures as the one above.

Developments of models and techniques, initially driven by applications in particle physics, could accelerate developments in this domain.

#### Subtopics:

- Heterogeneous GNN architectures.
- The next big thing in geometric deep learning ? Modelling complex systems requires going beyond graphs.
- Green AI is an integral part of this. Heterogeneous architectures will likely be run on low-level reconstruction tasks like particle flow and track reconstruction, i.e. run on essentially every event. This is where the big potential gains in energy savings are.

#### The definition of the contours of the conference is curently being finalised

- Spring 2024
- ~80 participants
- At the same beautiful venue right in the city centre as CTD 2023 (Le Village by CA and Flashback café, this has a start-up flair to it).
- With contributions from ANITI chairs
- Fix date before end of Nov. 2023