



A Comparison of CPU and GPU implementations for the LHCb Experiment Run 3 Trigger

Authors are listed on the following pages

Abstract

The LHCb experiment at CERN is undergoing an upgrade in preparation for the Run 3 data taking period of the LHC. As part of this upgrade the trigger is moving to a fully software implementation operating at the LHC bunch crossing rate. We present an evaluation of a CPU-based and a GPU-based implementation of the first stage of the High Level Trigger. After a detailed comparison both options are found to be viable. This document summarizes the performance and implementation details of these options, the outcome of which has led to the choice of the GPU-based implementation as the baseline.

Published in Computing Software for Big Science 6, Article number: 1 (2022)

R. Aaij²¹, M. Adinolfi³⁴, S. Aiola¹⁷, S. Akar⁴³, J. Albrecht¹⁰, M. Alexander³⁹, S. Amato¹, Y. Amhis⁸, F. Archilli¹², M. Bala²⁵, G. Bassi^{20,e}, L. Bian⁴⁵, M.P. Blago³¹, T. Boettcher⁴², A. Boldyrev⁴⁸, S. Borghi⁴⁰, A. Brea Rodriguez²⁹, L. Calefice^{10,9}, M. Calvo Gomez⁴⁹, D.H. Cámpora Pérez^{47,31}, A. Cardini¹⁹, M. Cattaneo³¹, V. Chobanova²⁹, G. Ciezarek³¹, X. Cid Vidal²⁹, J.L. Cobble Dick⁴⁰, J.A.B. Coelho⁸, T. Colombo³¹, A. Contu¹⁹, B. Couturier³¹, D.C. Craik⁴², R. Currie³⁸, P. d'Argent³¹, M. De Cian³², D. Derkach⁴⁸, F. Dordei¹⁹, M. Dorigo^{20,f}, L. Dufour³¹, P. Durante³¹, A. Dziurda²⁴, A. Dzyuba²⁶, S. Easo³⁷, S. Esen⁹, P. Fernandez Declara³¹, S. Filippov²⁷, C. Fitzpatrick⁴⁰, M. Frank³¹, P. Gandini¹⁷, V.V. Gligorov⁹, E. Golobardes⁴⁹, G. Graziani¹⁴, L. Grillo⁴⁰, P. A. Günther¹², S. Hansmann-Menzemer¹², A.M. Hennequin³¹, L. Henry^{17,30}, D. Hill³², S.E. Hollitt¹⁰, J. Hu¹², W. Hulsbergen²¹, R.J. Hunter³⁶, M. Hushchyn⁴⁸, B.K. Jashal³⁰, C.R. Jones³⁵, S. Klaver²¹, K. Klimaszewski²⁵, R. Kopecna¹², W. Krzemien²⁵, M. Kucharczyk²⁴, R. Lane³⁴, F. Lazzari^{20,e}, R. Le Gac⁷, P. Li¹², J.H. Lopes¹, M. Lucio Martinez²¹, A. Lupato⁴⁰, O. Lupton³⁶, X. Lyu⁴, F. Machefert⁸, O. Madejczyk²³, S. Malde⁴¹, J.F. Marchand⁶, S. Mariani^{14,a,31}, C. Marin Benito³¹, D. Martinez Santos²⁹, F. Martinez Vidal³⁰, R. Matev³¹, M. Mazurek³¹, B. Mitreska⁴⁰, D.S. Mitzel³¹, M.J. Morello^{20,c}, H. Mu², P. Muzzetto^{19,31}, P. Naik³⁴, M. Needham³⁸, N. Neri^{17,b}, N. Neufeld³¹, N.S. Nolte^{10,31}, D. O'Hanlon³⁴, A. Oyanguren³⁰, M. Pepe Altarelli³¹, S. Petrucci³⁸, M. Petruzzo¹⁷, L. Pica^{20,e}, F. Pisani³¹, A. Piucci¹², F. Polci⁹, A. Poluektov⁷, E. Polycarpo¹, C. Prouve²⁹, G. Punzi^{20,d}, R. Quagliani⁹, R.I. Rabadan Trejo⁷, M. Ramos Pernas³⁶, M.S. Rangel¹, F. Ratnikov^{41,48}, G. Raven²², F. Reiss⁹, V. Renaudin⁴¹, P. Robbe⁸, A. Ryzhikov⁴⁸, M. Santimaria¹⁵, M. Saur¹⁰, M. Schiller³⁹, R. Schwemmer³¹, B. Sciascia¹⁵, A. Solomin^{34,50}, F. Suljik⁴¹, N. Skidmore⁴⁰, M.D. Sokoloff⁴³, P. Spradlin³⁹, M. Stahl⁴³, S. Stahl³¹, H. Stevens¹⁰, L. Sun⁴⁵, A. Szabelski²⁵, T. Szumlak²³, M. Szymanski³¹, D.Y. Tou⁹, G. Tuci^{20,d}, A. Usachov²¹, N. Valls Canudas²⁸, R. Vazquez Gomez²⁹, S. Vecchi¹³, M. Vesterinen³⁶, X. Vilasis-Cardona⁴⁹, D. Vom Bruch⁷, Z. Wang³³, T. Wojton²⁴, M. Whitehead³⁴, M. Williams^{42,44}, M. Witek²⁴, Y. Xie⁵, A. Xu³, H. Yin⁵, M. Zdybal²⁴, O. Zenaiev³¹, D. Zhang⁵, L. Zhang², X. Zhu²

¹Universidade Federal do Rio de Janeiro (UFRJ), Rio de Janeiro, Brazil

²Center for High Energy Physics, Tsinghua University, Beijing, China

³School of Physics State Key Laboratory of Nuclear Physics and Technology, Peking University, Beijing, China

⁴University of Chinese Academy of Sciences, Beijing, China

⁵Institute of Particle Physics, Central China Normal University, Wuhan, Hubei, China

⁶Univ. Grenoble Alpes, Univ. Savoie Mont Blanc, CNRS, IN2P3-LAPP, Annecy, France

⁷Aix Marseille Univ, CNRS/IN2P3, CPPM, Marseille, France

⁸Université Paris-Saclay, CNRS/IN2P3, IJCLab, Orsay, France

⁹LPNHE, Sorbonne Université, Paris Diderot Sorbonne Paris Cité, CNRS/IN2P3, Paris, France

¹⁰Fakultät Physik, Technische Universität Dortmund, Dortmund, Germany

¹¹Max-Planck-Institut für Kernphysik (MPIK), Heidelberg, Germany

¹²Physikalisches Institut, Ruprecht-Karls-Universität Heidelberg, Heidelberg, Germany

¹³INFN Sezione di Ferrara, Ferrara, Italy

¹⁴INFN Sezione di Firenze, Firenze, Italy

¹⁵INFN Laboratori Nazionali di Frascati, Frascati, Italy

¹⁶INFN Sezione di Genova, Genova, Italy

¹⁷INFN Sezione di Milano, Milano, Italy

¹⁸INFN Sezione di Milano-Bicocca, Milano, Italy

¹⁹INFN Sezione di Cagliari, Monserrato, Italy

²⁰INFN Sezione di Pisa, Pisa, Italy

²¹Nikhef National Institute for Subatomic Physics, Amsterdam, Netherlands

²²Nikhef National Institute for Subatomic Physics and VU University Amsterdam, Amsterdam, Netherlands

²³AGH - University of Science and Technology, Faculty of Physics and Applied Computer Science,

Kraków, Poland

²⁴ Henryk Niewodniczanski Institute of Nuclear Physics Polish Academy of Sciences, Kraków, Poland

²⁵ National Center for Nuclear Research (NCBJ), Warsaw, Poland

²⁶ Petersburg Nuclear Physics Institute NRC Kurchatov Institute (PNPI NRC KI), Gatchina, Russia

²⁷ Institute for Nuclear Research of the Russian Academy of Sciences (INR RAS), Moscow, Russia

²⁸ ICCUB, Universitat de Barcelona, Barcelona, Spain

²⁹ Instituto Galego de Física de Altas Enerxías (IGFAE), Universidade de Santiago de Compostela, Santiago de Compostela, Spain

³⁰ Instituto de Física Corpuscular, Centro Mixto Universidad de Valencia - CSIC, Valencia, Spain

³¹ European Organization for Nuclear Research (CERN), Geneva, Switzerland

³² Institute of Physics, Ecole Polytechnique Fédérale de Lausanne (EPFL), Lausanne, Switzerland

³³ Physik-Institut, Universität Zürich, Zürich, Switzerland

³⁴ H.H. Wills Physics Laboratory, University of Bristol, Bristol, United Kingdom

³⁵ Cavendish Laboratory, University of Cambridge, Cambridge, United Kingdom

³⁶ Department of Physics, University of Warwick, Coventry, United Kingdom

³⁷ STFC Rutherford Appleton Laboratory, Didcot, United Kingdom

³⁸ School of Physics and Astronomy, University of Edinburgh, Edinburgh, United Kingdom

³⁹ School of Physics and Astronomy, University of Glasgow, Glasgow, United Kingdom

⁴⁰ Department of Physics and Astronomy, University of Manchester, Manchester, United Kingdom

⁴¹ Department of Physics, University of Oxford, Oxford, United Kingdom

⁴² Massachusetts Institute of Technology, Cambridge, MA, United States

⁴³ University of Cincinnati, Cincinnati, OH, United States

⁴⁴ School of Physics and Astronomy, Monash University, Melbourne, Australia, associated to ³⁶

⁴⁵ School of Physics and Technology, Wuhan University, Wuhan, China, associated to ²

⁴⁶ Universität Bonn - Helmholtz-Institut für Strahlen und Kernphysik, Bonn, Germany, associated to ¹²

⁴⁷ Universiteit Maastricht, Maastricht, Netherlands, associated to ²¹

⁴⁸ National Research University Higher School of Economics, Moscow, Russia, associated to Yandex School of Data Analysis, Moscow, Russia

⁴⁹ DS4DS, La Salle, Universitat Ramon Llull, Barcelona, Spain, associated to ²⁸

⁵⁰ Institute of Nuclear Physics, Moscow State University (SINP MSU), Moscow, Russia

^a Università di Firenze, Firenze, Italy

^b Università degli Studi di Milano, Milano, Italy

^c Scuola Normale Superiore, Pisa, Italy

^d Università di Pisa, Pisa, Italy

^e Università di Siena, Siena, Italy

^f INFN Sezione di Trieste, Trieste, Italy

Contents

1	Introduction	1
1.1	DAQ and HLT architectures	2
1.1.1	CPU-only architecture	3
1.1.2	Hybrid architecture	3
1.2	Real Time Analysis	4
2	Assumptions and boundary conditions	5
2.1	Use of storage and computing resources during and outside datataking	6
2.2	Existing and pledged HLT2 computing resources	6
2.3	Event building and data-flow	6
2.4	Disk buffer boundary conditions	7
3	HLT1 sequence	7
4	Throughput	9
5	Physics performance	9
5.1	Tracking efficiencies and ghost rates	10
5.2	Track parameter resolution	11
5.3	Primary vertex reconstruction efficiency and resolution	11
5.4	Muon ID efficiency	12
5.5	HLT1 Efficiency	14
5.6	HLT1 rates	16
5.7	Global Event Cut efficiency	18
5.8	Summary of the physics performance	18
6	Cost benefit analysis	19
6.1	Nominal cost-benefit	19
6.2	Disk buffer simulation	20
6.2.1	Quantification of HLT1 and HLT2 throughputs	21
6.2.2	Additional free computing resources for HLT2	21
6.2.3	Results	21
6.2.4	Results assuming a factor two performance increase in HLT2	22
6.3	Deferred purchasing	23
7	Conclusion	23
	References	25

1 Introduction

The LHCb experiment is a general purpose spectrometer instrumented in the forward direction based at the LHC. [1] Although optimized for the study of hadrons containing beauty and charm quarks, LHCb’s physics programme gradually expanded over the course of its first datataking period¹, taking in Kaon physics, prompt spectroscopy, Electroweak physics, and searches for putative heavy and long-lived exotic particles beyond the Standard Model.

The LHC provides a non-empty bunch crossing rate of up to 30 MHz, with a variable number of inelastic proton-proton interactions per bunch crossing (μ) which can be adjusted to suit the physics programme of each experiment. During Runs 1 and 2 LHCb took data with a μ of between 1.1 and 2.5, corresponding to an instantaneous luminosity of around $4 \cdot 10^{32} \text{ cm}^{-2}\text{s}^{-1}$. A fixed-latency hardware trigger (L0), based on calorimeter and muon system information, reduced the 30 MHz LHCb bunch crossing rate to ~ 1 MHz at which the detector readout operated. These events were then passed to a two stage asynchronous High Level Trigger (HLT) system entirely implemented in software. In the first stage, HLT1, a partial event reconstruction selected events based on inclusive signatures, reducing the event rate by an order of magnitude. Accepted events were stored using a 11 PB disk buffer in order to align and calibrate the detector. After this procedure, events were passed to the second stage, HLT2, which had sufficient computing resources to run the full offline-quality detector reconstruction. A multitude of dedicated selections deployed in HLT2 reduced the data to an output rate of 12.5 kHz using a combination of full and reduced event formats [2].

The ambitious goal of the upgraded LHCb experiment in Run 3 (i.e. starting from 2022) is to remove the hardware trigger and directly process the full 30 MHz of data at an increased luminosity of $2 \cdot 10^{33} \text{ cm}^{-2}\text{s}^{-1}$ in the HLT. At this luminosity, corresponding to a μ of around 6, it is no longer possible to efficiently identify bunch crossings of interest based purely on calorimeter and muon system information, as there is too much QCD background generated by the pileup pp collisions. [3] It is necessary to instead fully read the detector out for every bunch crossing and fully perform the real-time processing in the HLT. This allows the much more sophisticated selections, in particular selections based on charged particle trajectories reconstructed in the whole of LHCb’s tracker system, to be deployed already in HLT1. Such a fully software trigger will not only allow LHCb to maintain its Run 1 and 2 efficiencies for muonic signatures, but will lead to a factor two improvement in efficiency for hadronic signatures compared to the calorimeter-based L0 trigger despite the harsher Run 3 environment.

LHCb’s Run 3 datataking conditions imply a data volume of around 32 Terabits (Tb) per second, comparable to what the ATLAS [4] and CMS [5] software triggers will be required to process during High-Luminosity LHC running from 2027 onwards. The design and delivery of LHCb’s high-level trigger is therefore also one of the biggest computing challenges that the field of High Energy Physics is facing today. The closest current parallel to LHCb’s system is that of the ALICE experiment [6], which will also operate a triggerless readout in Run 3 with an objective to reduce an input data rate of

¹The first datataking period was broken into two “runs”, with Run 1 taking place from 2009 to 2013 and Run 2 taking place from 2015 to 2018.

roughly 10 Tb/s to a manageable amount by performing a full detector reconstruction in real-time.

The concept of a pure CPU-based solution for this approach was reviewed during the preparation of LHCb’s Trigger and Online TDR [7] in 2014 followed by a systematic rewrite of the LHCb trigger software infrastructure, which enabled data taking in these conditions. In parallel, R&D efforts have explored a possible usage of GPUs for HLT1 [8], referred to as the hybrid approach in the following. An intensive effort was launched to demonstrate if such a hybrid system could be delivered in time for Run 3 and concluded in a positive review of its TDR [9] in early 2020. After a detailed comparison of both options the collaboration selected the hybrid approach as the new Run 3 baseline. This decision parallels that of ALICE, which pioneered the use of GPUs among LHC experiments during the last decade [10] and whose Run 3 triggerless readout and full real-time reconstruction mentioned earlier will be mainly implemented on GPUs.

This document compares both options and summarizes the salient points which led to the decision. It reflects the status of both implementations at the time the decision was made, in April 2020. Further significant improvements [11] in both throughput and physics performance, which will not be discussed in this document, have been achieved since then.

This document is structured as follows. In Section 1, the DAQ and HLT architecture of both systems is summarized. Section 2 describes the boundary conditions within which the implemented HLT1 triggers must operate, including available financial resources and operational constraints. Section 3 describes the HLT1 trigger sequence and algorithms, which are the basis of the performance comparison. Section 4 summarizes the performance of the two architectures in terms of throughput, while Section 5 presents their physics performance. Section 6 combines the performance assessment into a cost-benefit analysis.

1.1 DAQ and HLT architectures

LHCb’s DAQ and event building (EB) infrastructure is described in the Trigger and Online TDR [7]. The full detector is read out for all LHC bunch-crossings, and information from subdetectors is received by around 500 custom backend FPGA “TELL40” boards hosted in a farm of EB servers, with three TELL40 boards per server. These subdetector fragments are then assembled into “events”, with one event corresponding to one LHC bunch crossing, and sent to the HLT for processing. Both the event building and HLT are fully asynchronous and no latency requirements exist in the system.

The following two HLT processing architectures are under consideration.

1. **CPU-only**, which implements both HLT1 and HLT2 using Event Filter Farm (EFF) CPU servers.
2. **Hybrid**, which implements HLT1 using GPU cards installed in the EB servers with the CPU-based HLT2 running as before on the EFF.

The HLT2 software and processing architecture are identical in both cases.

1.1.1 CPU-only architecture

The CPU-only architecture is illustrated in Fig. 1. Briefly, it consists of

- A set of custom FPGA cards, called TELL40, which together receive on average 32 Tb/s of data from LHCb’s subdetectors;
- A set of EB servers which host the TELL40 cards and implement a network protocol to bring the subdetector data fragments produced in a single LHC bunch crossing (“event” in LHCb nomenclature) together and then group $O(1000)$ of these events into multi-event packets (MEPs) to minimize I/O overheads further down the line. The EB servers will be equipped with 32-core AMD EPYC (7502) CPUs;
- An EFF illustrated on the bottom left of Fig. 1 which receives MEPs from the EB and executes HLT1. The memory available in the EB servers and HLT1 EFF nodes allows data to be buffered for $O(20)$ seconds in case of temporary network or processing issues. The HLT1 EFF servers are assumed to be equipped with the same 32-core AMD EPYC (7502) CPUs as the EB servers. When LHCb is not taking data this HLT1 EFF can also receive events from the disk servers (described below) and run the HLT2 process on them;
- A high performance network, including high-speed network interface cards (NICs) located in the EB servers and EFF nodes as well as a large switch, connects the EB and HLT1 EFF and allows to transmit at a rate of 32 Tb/s.
- The HLT1 process reduces the event rate by a factor of between 30-60, and an array of disk servers buffers this HLT1 output data while the detector alignment and calibration are performed in quasi-real-time. This disk buffer, whose size is a tunable parameter of the system as discussed later in Section 6, allows HLT1 output data to be buffered for $O(100)$ hours in case of problems with the alignment and calibration which require specialist intervention;
- A second EFF receives events from the disk servers once the alignment and calibration constants are available and runs the HLT2 process on them. Because of limitations in network bandwidth this second EFF cannot be used to process HLT1.

Advances in server technology have permitted a substantial reduction of the number of servers required by the EB, from 500 in the TDR [7] to around 173. This allows to host three TELL40 cards per EB server rather than the one card foreseen in the TDR. This improvement means that the EB will be much more compact and as a consequence easier to upgrade in future.

1.1.2 Hybrid architecture

The hybrid architecture is illustrated in Fig. 2. It follows the same processing logic as the CPU-only solution: the full detector data is received by TELL40 boards and assembled into MEP packets by the EB servers, those MEP packets are then sent for processing by HLT1 and the events selected by HLT1 are recorded to the disk buffer for later processing by HLT2. Compared to the CPU-only solution it replaces the HLT1 EFF by GPU cards running HLT1 which are installed in the spare PCI express slots available in the EB

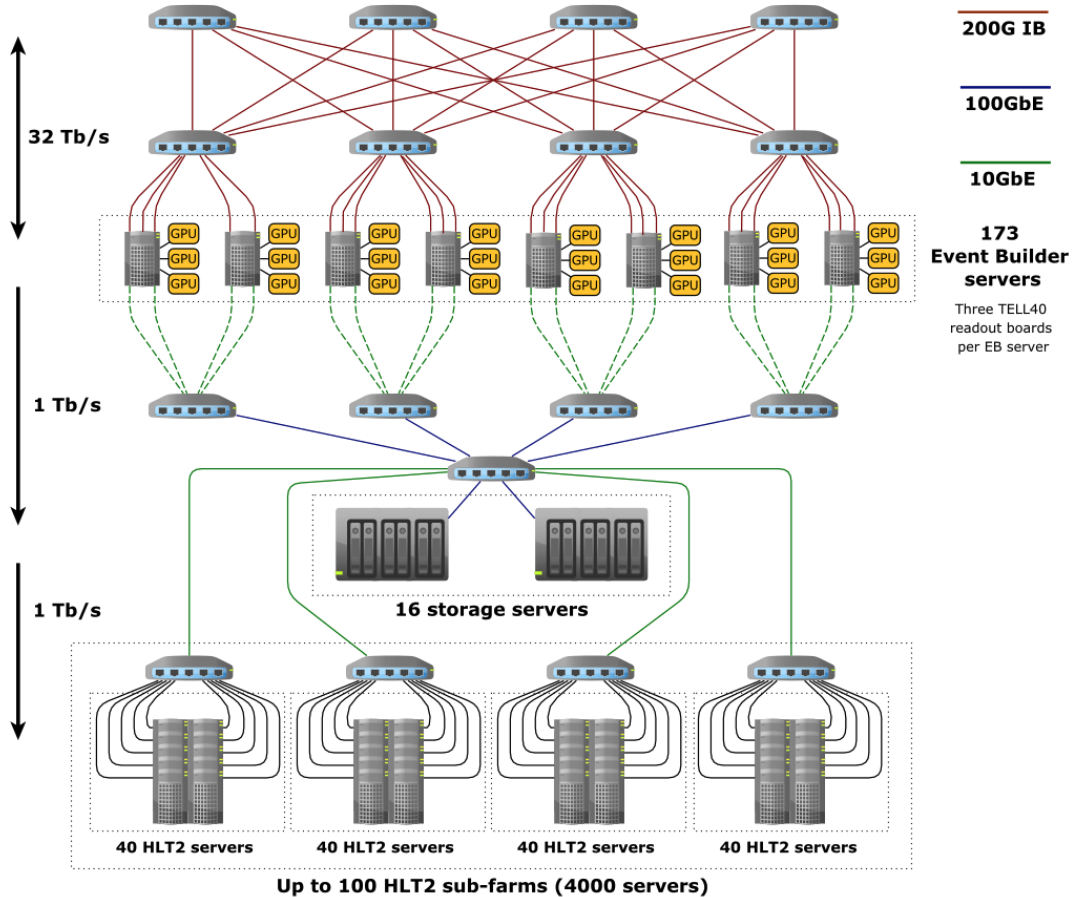


Figure 2: Run 3 DAQ architecture in the case of the hybrid solution, with GPUs placed in the EB servers to reduce the data rate. Labels are the same as in Figure 1.

physics programme is written to the TURBO stream [14], a reduced format, which on one hand allows flexible event information to be added in a selective manner and on the other hand can discard a user-specified fraction of both the raw and reconstructed detector data.

- Consequently, HLT2 must be able to run the complete offline-quality reconstruction. Therefore HLT2 must use the highest quality alignment and calibration at all times. A substantial disk buffer must therefore be purchased to allow events selected by HLT1 to be temporarily saved while the full detector is aligned and calibrated in real-time. This buffer must be big enough not only to cover the steady-state data taking conditions but also to permit recovery from unforeseen operational issues in a reasonable timescale without loss of data.

2 Assumptions and boundary conditions

Having described the overall design of LHCb's Run 3 DAQ and HLT, as well as the processing technologies under consideration, we will now describe the boundary conditions

which these technologies have to respect, as well as common assumptions relevant to the cost-benefit comparison in Section 6.

2.1 Use of storage and computing resources during and outside datataking

Throughout this document it is assumed that the LHC is in “data-taking” for 50% of the year, and in either the winter shutdown or longer technical stops for the other 50%. When in data-taking, it is assumed that the LHC is in stable beams 50% of the time. During data-taking, it is assumed that all CPU resources are used to process HLT1 and/or HLT2. Outside data-taking, it is assumed that all HLT CPU resources are used to produce simulation for LHCb analyses. GPU resources can only be used to process HLT1 in-fill, and cannot be used opportunistically during data-taking. They cannot yet be used for producing simulation, and there is no realistic prospect of this changing on a short timescale. However in principle GPU resources could be used out of data-taking if use-cases can be found, as discussed in Ref. [9]. When LHCb is not taking data, the EB nodes will be used to produce simulation. However they will not be used for any task other than event building while taking data, as all the available memory bandwidth in these nodes will be required for the event building process and transferring data to the HLT1 application.

2.2 Existing and pledged HLT2 computing resources

We quantify the computing resources available for HLT2 in terms of a reference QuantaPlex (“Quanta”) server consisting of two Intel E5-2630v4 10-core processors, which was the workhorse of our Run 2 HLT. This reference node corresponds to approximately 380 HEPSPEC.² We currently have roughly 1450 equivalent such servers available for Run 3 processing, with a further 1200 equivalent servers pledged, corresponding to a total capacity of slightly more than one million HEPSPEC. These servers can only be used to process HLT2 as it would not be cost-effective to equip so many old servers with the high-speed NICs required to process HLT1. So far no economical way has been found to reuse the predominantly very small disk-drives in the old Run 2 servers, so there are no free storage resources available.

2.3 Event building and data-flow

In the CPU-only scenario no data reduction happens before the EFF, so all the data collected in the building network has to be distributed to the HLT1 CPU nodes at 32 Tb/s. Using newly available AMD 32-core CPUs a cost-effective implementation is to use a dual-socket server with a total of 64 physical cores and two network interfaces of 100 Gb/s. Each EB node requires two high-speed network connections for sending the event fragments between EB nodes while they are built. In addition, the distribution

²The HEPSPEC benchmark is defined at <http://w3.hepix.org/benchmarking>. In reality our HLT2 farm will consist of a mixture of servers of different generations and with different numbers of physical cores, but because of the asynchronous nature of HLT2 processing load-balancing between these servers is an implementation detail and it is more convenient to quantify the available resources in units of the reference node as if the system were fully homogeneous.

network needs to connect all HLT1 CPU nodes to the EB nodes as well as to at least 30 storage servers. These connections need to be optical.

In the hybrid scenario the GPUs hosted within the EB nodes execute the HLT1 process and reduce the data rate, so that only 0.5 – 1.0 Tb/s has to be sent from the EB to the disk buffer servers. The EB servers can therefore use their on-board 10 Gigabit interfaces to send the data and the distribution network needs significantly fewer optical connections.

2.4 Disk buffer boundary conditions

The disk buffer needs to be able to handle at least 1 MHz of events coming from HLT1, with a potential upgrade to be able to handle 2 MHz as Run 3 progresses. A typical minimum bias event in LHCb Run 3 conditions is expected to be around 100 kB, however events selected by HLT1 are bigger than average minimum bias events since they typically contain a hard QCD process leading to the production of a heavy flavour hadron. Therefore assuming an event size of 120 kB to account for this effect, this implies 120 GB/s both for writing events coming from HLT1 and reading events out into HLT2. The nominal rate of a 12 TB disk is 100 MB/s, thus 50 MB/s for writing and 50 MB/s for reading. However the read and write speed of hard disks decreases as they fill up. As the system must be costed so that there is no dead-time even when the buffer is almost full, an effective sustainable write and read rate of a single 12 TB disk is assumed to be 35-40 MB/s.

Since this part of the system is hardware limited and must be able to handle burst data rates, a minimum of 2880 physical disks is required assuming minimal redundancy, and 3120 physical disks with adequate redundancy. A final point to note on the disk buffer is that the usable disk sizes are only around 80% of the nominal disk size, so the 12 TB disks actually only provide around 9.6 TB of usable storage each. In practice, the cost of such a minimal disk buffer is so large compared to the overall budget discussed earlier that spending money on buying additional disks or bigger disks is not really an interesting option.

3 HLT1 sequence

The CPU-only and hybrid approaches under study in this document implement an HLT1 configuration which broadly corresponds to the one used in Run 2 [2] and whose physics objectives have been described in reference [9]. The reconstruction consists of the following components, whose performance in terms of efficiency, purity and resolution is discussed in Section 5.

- Vertex Locator (Velo) detector decoding, clustering, and reconstruction. Conceptually very similar algorithms are used here in the CPU and the GPU implementation. Only minor differences in physics performance are expected due to a limited number of architecture-specific optimizations.
- The primary vertex (PV) finding with tracks reconstructed in the Velo. Again only minor differences are expected in the physics performance of the CPU- and GPU-based implementations.

- Decoding of the UT and SciFi tracking detector raw banks. The raw bank formats are converted to the input of the pattern recognition algorithms. These are specific to the data layout and event model used by each architecture.
- Reconstruction of Velo-UT track segments. Different algorithms are used here for the CPU and the GPU implementations, which lead to slight differences in the physics performance.
- Reconstruction of long tracks³ starting from reconstructed Velo-UT track segments. Both the CPU and GPU tracking algorithms use a parameterization of particle trajectories in the LHCb magnetic field and the initial Velo-UT momentum estimate⁴ to speed up their reconstruction. One major difference in strategy is that the CPU algorithm matches Velo-UT track segments to doublets in the SciFi x-layers, while the GPU algorithm matches them to triplets. In addition, the CPU algorithm applies a 490 MeV transverse momentum threshold when defining its search windows. These choices, and the differences in the Velo-UT algorithms, lead to somewhat different performance.
- Decoding of the muon raw banks and calculation of crossing points in the muon system, as well as implementation of the muon identification algorithm. Both architectures use very similar algorithms here.
- A simplified Velo-only Kalman filter which uses the momentum estimate from the forward tracking to calculate a covariance matrix and estimate an uncertainty on the track impact parameter. Again the underlying algorithm is the same for both implementations.

The HLT1 sequence described covers most use cases required by bottom and charm physics. There are however a few dedicated algorithms missing, such as a reconstruction for high momentum muons for electroweak physics, reconstruction of beam gas events or a reconstruction of low momentum tracks down to a transverse momentum of 80 MeV, which is motivated by the strange physics programme. While preliminary versions of these algorithms were ready in time for this comparison, they were not yet fully optimized in the same way as the other described algorithms. We did use these preliminary versions to test whether their inclusion impacted on the relative throughput between the CPU and GPU implementations, and found that the relative slowdown from including these algorithms was similar across the two architectures. It is therefore expected that these missing components will not change the conclusions of this document.

In addition to the reconstruction algorithms listed above, both the CPU and GPU HLT1 implement a representative subset of HLT selections including the finding of displaced vertices, monitoring, as well as writing of decision and selection reports. While the specific components are not exactly the same in both cases, these parts of the system consume comparatively little throughput and hence the differences are not relevant for

³Long tracks are tracks that traverse the entire tracking system from Velo to SciFi. They deliver the best parameter estimate in terms of position and momentum and thus are the most valuable tracks for physics analysis.

⁴The momentum resolution of Velo-UT tracks is about 15% with significant non-Gaussian tails due to the small and inhomogeneous B field between Velo and UT.

this comparison.

A global event cut that removes 7% of the busiest minimum bias events is applied in all selections before the reconstruction is performed for both the CPU and GPU HLT1 implementations. The criterion used is that the total number of hits in the UT and SciFi tracking detectors be below a certain value. This criterion is chosen for historical reasons and because the UT and SciFi reconstructions are the most sensitive to occupancy, especially when trying to reconstruct low p_T signatures. However in practice any subdetector occupancy could be used, and LHCb will likely use the occupancy of whichever subdetector shows the best data-simulation agreement once the new detector is commissioned in Run 3. The impact of this criterion on the physics is given in Section 5.7.

4 Throughput

We define throughput as the number of events which can be processed by a given architecture per second in steady-state conditions, that is to say neglecting the time it takes to initialize the HLT1 application at the start of each datataking period. They can be converted into GB/s by multiplying by the average expected Run 3 event size of 100 kB. An event is processed when it is read into HLT1, the HLT1 reconstruction and selection sequences used to decide whether or not to keep this event, and the event is finally written out (or not). The throughput of the CPU and GPU HLT1 implementations are measured using the same minimum bias samples. The GPU throughput is measured directly on the candidate production card. The CPU throughput is measured on both the Quanta reference nodes used for HLT2, and the dual-socket AMD 7502 EPYC nodes used for HLT1. While the HLT1 throughput measurements include both reconstruction, selection, and saving of trigger candidates, the HLT2 throughput measurement only includes the reconstruction and not the selection and saving of trigger candidates. Additional costs associated with these missing components of HLT2 are expected but we expect that they can be counterbalanced with future performance improvements in the reconstruction.

The measured throughputs used in the rest of this document are

- CPU HLT1 : 171 kHz;
- GPU HLT1 : 92 kHz;
- CPU HLT2 on an AMD EPYC node : 471 Hz;
- CPU HLT2 on a Quanta node : 134 Hz.

Although we will discuss the cost-benefit of the two architectures later in Section 6, we can already conclude that both the CPU and GPU HLT1 architectures can be implemented using rather compact systems of $O(170)$ CPU servers or $O(330)$ GPU cards.

5 Physics performance

This section presents key figures which are evaluated in a like-for-like comparison of the GPU and CPU performance. This includes performance numbers for track reconstruction

and fitting, for PV reconstruction, for Muon-ID and the overall HLT performance for some representative trigger selections. Identical algorithms are used to fill the histograms and produce the plots based on the output of HLT1. The source code is compiled to operate on both the GPU and the CPU. The output is also translated to the same format. This ensures that the same definitions of physics performance parameters such as efficiencies and resolutions are used when doing the comparison. The CPU compiled version of the GPU code has been checked to give results which agree to within $10^{-4} - 10^{-3}$ with results obtained with the GPU version. More details of this comparison can be found in reference [9].

An overview of the samples used for these studies is given in Tab. 1. The specific samples used for the individual studies are listed in the corresponding subsections.

Sample	# Events
$J/\psi \rightarrow \mu^+ \mu^-$	10k
$B_s^0 \rightarrow \phi \phi$	10k
$B^0 \rightarrow K^{*0} e^+ e^-$	10k
$B^0 \rightarrow K^{*0} \mu^+ \mu^-$	10k
$D_s \rightarrow K \pi \pi$	10k
$Z \rightarrow \mu^+ \mu^-$	10k
Minimum Bias	10k

Table 1: Simulated samples used to evaluate the physics performance. All samples are simulated using the same pileup assumed in the Run 3 minimum bias samples.

5.1 Tracking efficiencies and ghost rates

All plots shown in this section are based on the B and D signal samples in Tab. 1. Tracking efficiencies are defined as the number of reconstructed tracks out of the number of reconstructible tracks in a given sub-sample. In this section only particles which are reconstructible as long tracks are studied, which essentially means that they have left signals on at least 3 pixel sensors in the Velo and one x and one stereo cluster in each of the 3 fibre tracker (FT) stations. Furthermore, the studies are restricted to B and D daughter particles which are in the range $2 < \eta < 5$, which have true transverse momenta $p_T \geq 0.5$ GeV/c and true momenta $p \geq 3$ GeV/c. Electrons are explicitly excluded because their performance has not yet been optimised to the same extent as that of hadrons and muons, however we have checked the electron performance for both the CPU and GPU implementations of HLT1 and it does not change our conclusions.

Efficiencies to reconstruct these particles in the Velo, in both the Velo and UT and in the Velo, UT, and FT are shown as functions of the true p and p_T of the particles in Fig. 3. The performances of both technologies are very similar, the only difference is seen in the momentum dependence of the long track reconstruction efficiency. This difference is minor in the context of the overall trigger efficiency since the majority of HLT1 triggers use only a subset of tracks coming from a given signal decay to select an event, and are therefore inherently robust to differences in single-track reconstruction

efficiencies at the level of a few percent. While the CPU implementation is more efficient at low momenta, the GPU implementation is better at high momenta and the overall single-track efficiency integrated over the kinematic range of interest agrees to better than 1% between the two algorithms. This difference is expected to be entirely related to the tuning of the algorithms and not to the underlying technology.

To check for a possible decay time bias in the Velo reconstruction, the Velo efficiency for long reconstructible tracks is studied as a function of the distance of closest approach to the beamline, $docaz$ and as a function of the z position of the primary vertex in the event (Fig. 4). Again the performance of both implementations is very similar. The loss in efficiency at large $docaz$ is caused by the use of search windows which favour tracks coming from the beamline, and can be recovered at a moderate throughput cost.

Finally, Fig. 5 shows the fraction of ghost tracks among all long tracks. Ghost tracks are tracks which can not be assigned to a true particle, thus are fake combinations of signals in the detector. Again, the performance of both implementations is very similar.

5.2 Track parameter resolution

For the comparison of the impact parameter resolution, σ_{IP_x} , and the momentum resolution, σ_p/p , the minimum bias sample (Tab. 1) is used. The results are shown in Fig. 6. Note that the x and y component of the impact parameter have very similar resolution, therefore only σ_{IP_x} is shown here. The impact parameter resolution is very similar for both technologies. The momentum resolution is worse in the GPU framework, with a maximum absolute resolution difference of 0.15 – 0.2% at low momenta. This difference is caused by a suboptimal tuning of the parametrization used to derive the momenta of the particles in the GPU algorithm. However, since the computational cost of this parametrization is negligible compared to the track finding itself, this difference in performance can be recovered without any significant change in its throughput.

5.3 Primary vertex reconstruction efficiency and resolution

A simulated primary vertex is defined as reconstructed if a primary vertex is found within 2 mm of its true position. The primary vertex efficiency as a function of the number of its reconstructed Velo tracks and of the primary vertex z position is shown in Fig. 7. All plots in this section are obtained on a sample of minimum bias events (Tab. 1).

The primary vertex resolution in x , y and z are studied as a function of the number of Velo reconstructible particles associated to the primary vertex and as a function of the z position of the primary vertex. The results on minimum bias events are shown in Fig. 8. The x resolution of the primary vertex is very similar to the y resolution, thus only one of them is shown. The average resolution of the minimum bias data set is $\sigma_x = \sigma_y \sim 14 \mu m$ and $\sigma_z \sim 87 \mu m$ for both technologies. The performances in terms of efficiency as well as in terms of resolution are close to identical for the both implementations.

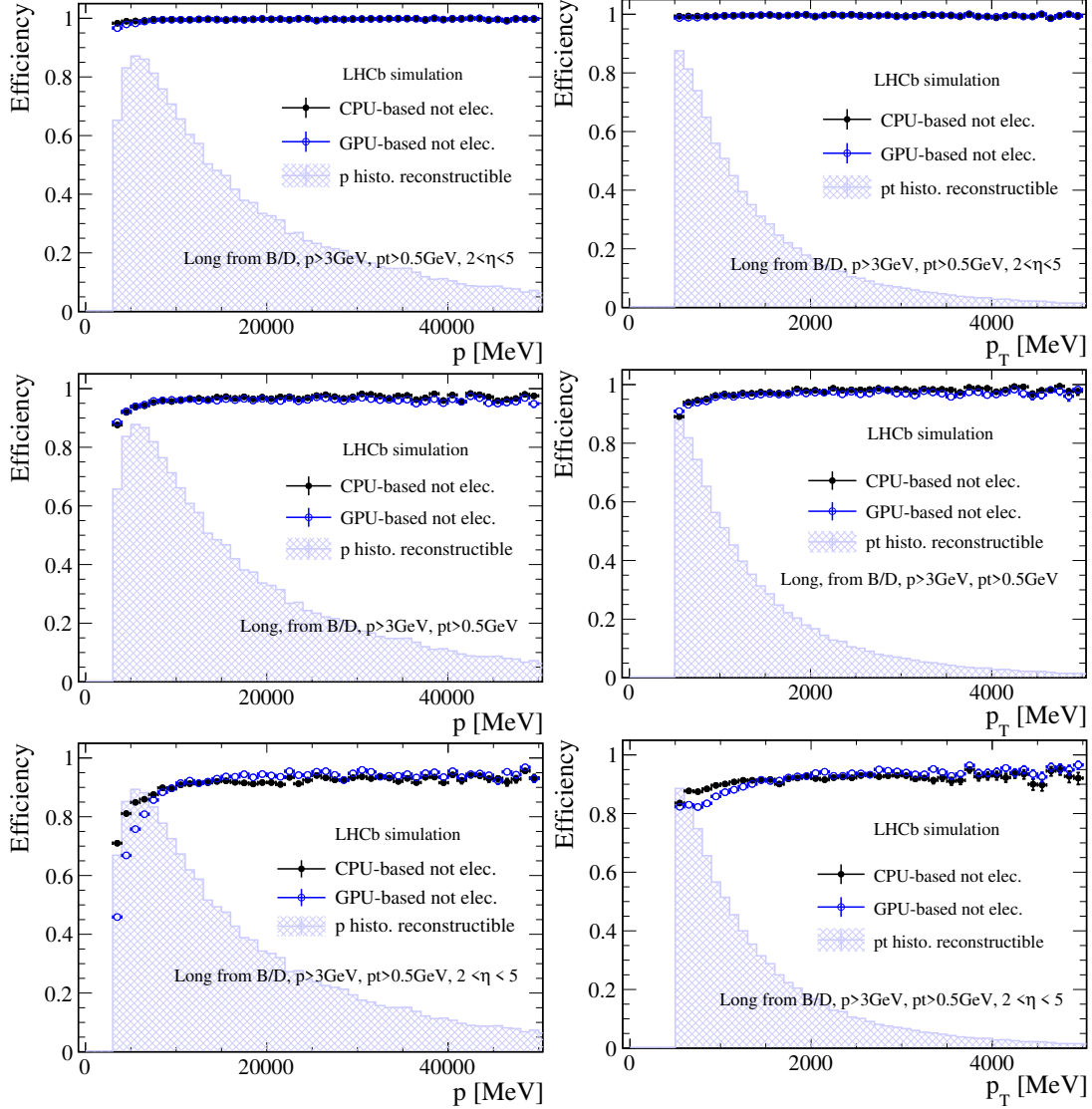


Figure 3: Track reconstruction efficiency versus momentum, p (left) and transverse momentum, p_T (right) for long reconstructible particles from B and D decays within $2 < \eta < 5$. The reconstruction efficiency in the Velo (top), as Velo-UT track (middle) or as long track (bottom) are shown.

5.4 Muon ID efficiency

The efficiency for the muon identification has been measured using the $J/\Psi \rightarrow \mu^+\mu^-$, $Z \rightarrow \mu^+\mu^-$ and $B^0 \rightarrow K^{*0}\mu^+\mu^-$ samples (Tab. 1). The denominator of the efficiency term counts muons with a minimum momentum of $p \geq 3$ GeV/c and a minimum transverse momentum of $p_T \geq 0.5$ GeV/c in the pseudorapidity range $2 \leq \eta \leq 5$. Furthermore, they must be reconstructed as long tracks. The numerator requires additionally these tracks to be identified as a muon in the detector. The efficiency is displayed in Fig. 9 as a function of the momentum of the muon and as function of its pseudorapidity. A second performance criterion is the purity of the muon reconstruction. For this we count which fraction of pions with a minimum momentum $p \geq 3$ GeV/c and a minimum transverse momentum of $p_T \geq 0.5$ GeV/c in the pseudorapidity range $2 \leq \eta \leq 5$ are misidentified

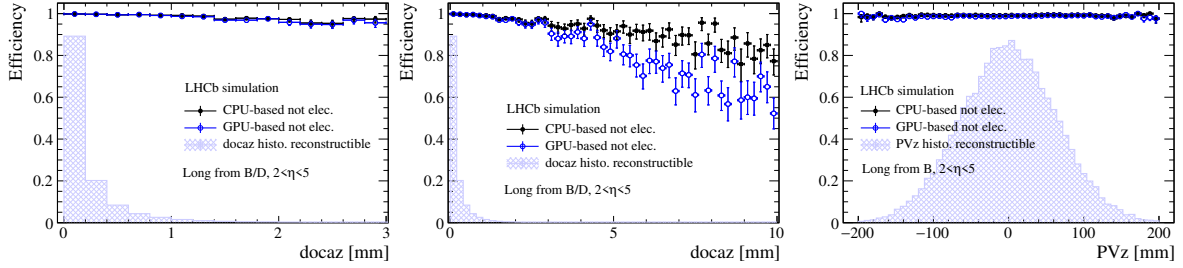


Figure 4: Reconstruction efficiency for long reconstructible particles from B and D decays within $2 < \eta < 5$ in the Velo as a function of distance of closest approach to the beam axis ($docaz$) and the z position of the primary vertex (PVz). The $docaz$ and z distribution of the long reconstructible particles are shown as well.

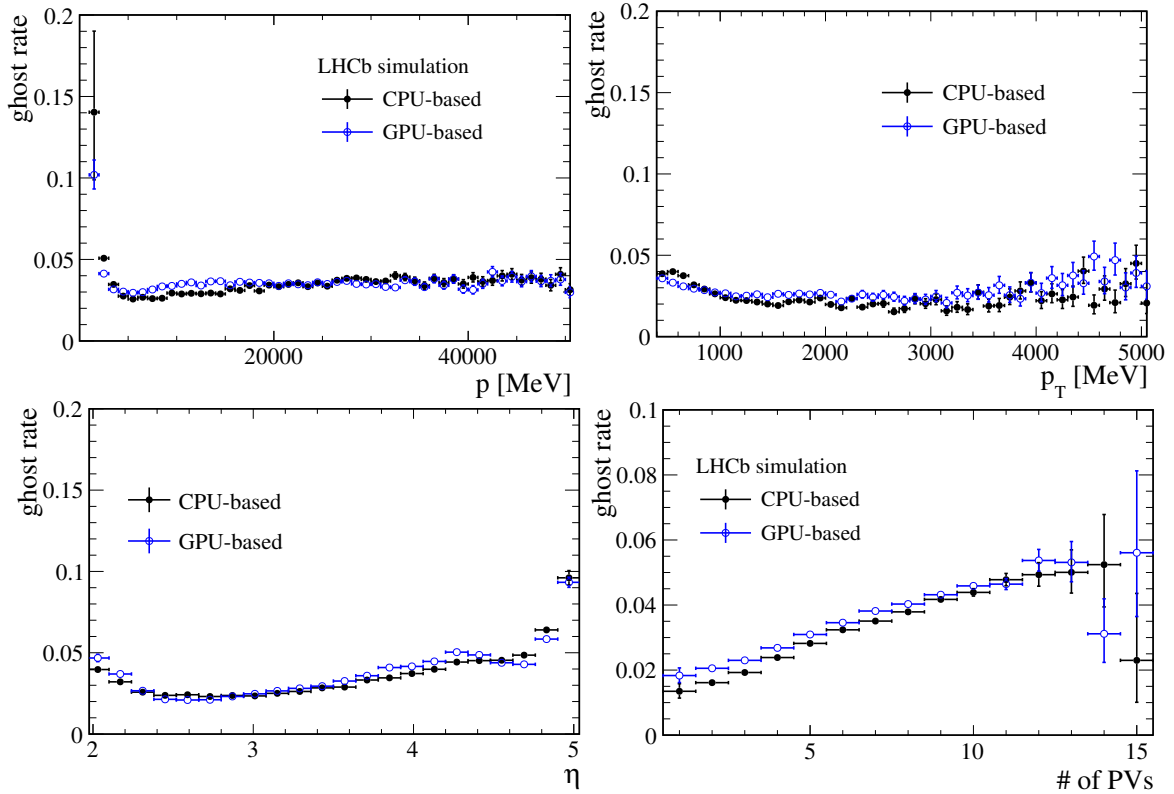


Figure 5: Ghost rate of long tracks as a function of momenta, p , transverse momenta, p_T , and pseudo-rapidity, η , of the reconstructed tracks and the number of primary vertices, $\#$ of PVs, in the event.

as muons by the reconstruction algorithm (Fig. 10). The GPU implementation of the muon identification gives better performance than the CPU version, with an absolute efficiency improvement of up to 10% at small pseudorapidities. The misidentification rate is similarly better by up to 5% absolute at small pseudorapidities and momenta for the GPU implementation, while it is a couple of percent better at higher momenta for the CPU implementation. It is expected that it is a matter of throughput-neutral tuning to obtain similar results for the two implementations.

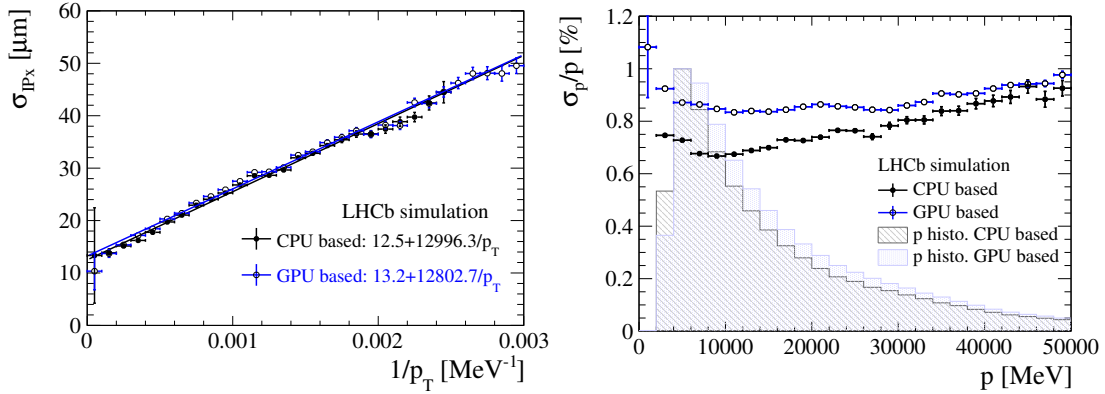


Figure 6: Resolution of the x projection of the impact parameter, σ_{IP_x} as function of the inverse transverse momentum, $1/p_T$, and the relative momentum resolution σ_p/p as function of the momentum, p .

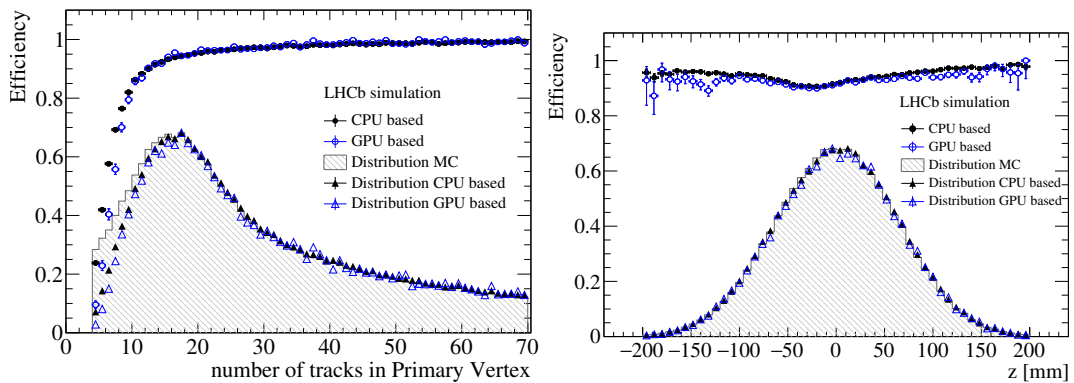


Figure 7: Efficiency to reconstruct primary vertices as function of the number of reconstructed Velo tracks associated to the simulated primary vertex and as function of the true vertex z position.

5.5 HLT1 Efficiency

In this part of the study, the full HLT1 sequence of both the CPU and the GPU implementations are run on the simulated signal samples listed in Tab. 1, and the trigger efficiencies are then compared. The efficiency denominator contains all signal candidates for which all stable charged daughter particles are reconstructible as long tracks within the detector acceptance of $2 < \eta < 5$. Thus, the trigger efficiency is the number of selected events divided by the total number of events. We choose four generic trigger selections to perform the comparison: a one-track MVA line (**TrackMVA**), a two-track MVA line (**TwoTrackMVA**), a low-mass dimuon line (**DiMuonLowMass**) and a single, high- p_T muon line (**SingleHighPtMuon**). The former two are inclusive selections designed to select generic B and D decays, the third is specialized for semileptonic B decays, whereas the latter is more suitable for high- p_T electroweak physics such as $Z \rightarrow \mu^+ \mu^-$. The selections that each trigger line performs are as follows:

- **TrackMVA**: requires a track with a minimum transverse momentum and a minimum significance of its impact parameter with respect to any primary vertex. These requirements are the same for the CPU- and GPU-based implementation, thus very

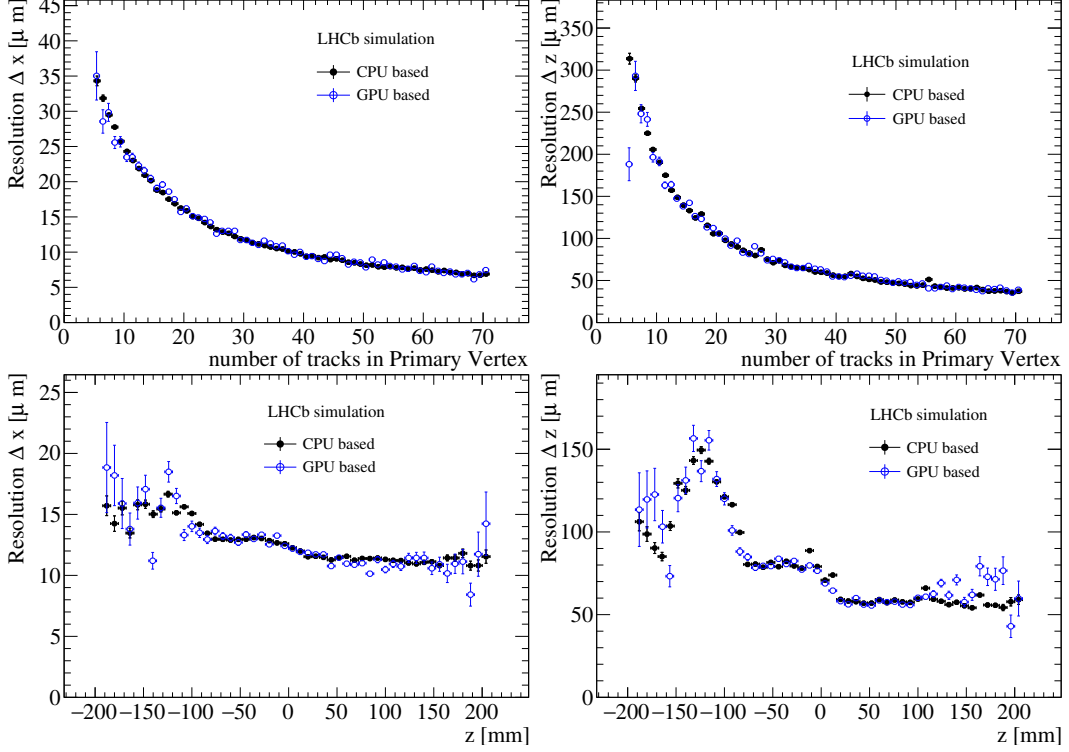


Figure 8: Resolution in x and z of all reconstructed primary vertices as function of the number of reconstructed Velo tracks associated to the simulated primary vertex (top row) and as function of the true vertex z position (bottom row).

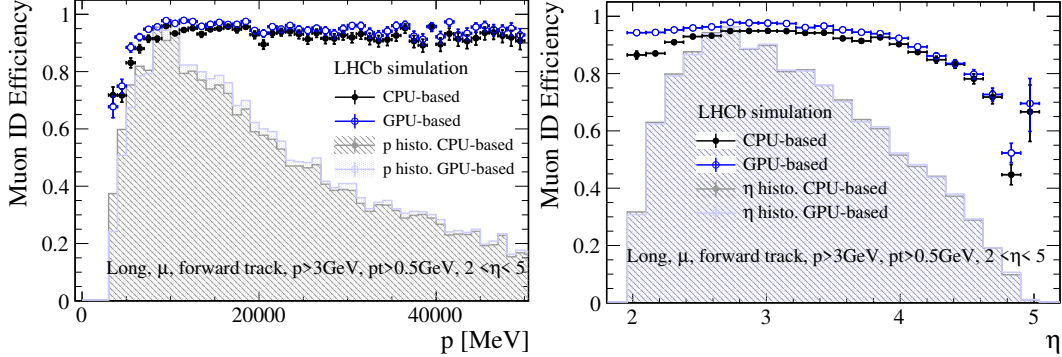


Figure 9: Efficiency to identify true muons reconstructed as long muon tracks as muons as function of the muon momentum and pseudorapidity.

similar performance is expected for this line.

- **TwoTrackMVA**: combines two displaced tracks of good quality with a minimum momentum to a vertex. In the CPU-based implementation a multivariate analysis (MVA) classifier is trained to identify those vertices which originate from a decay of a long-lived particle such as a beauty or charmed hadron. An MVA is not yet in place for the GPU-based implementation, so percent level differences in the signal efficiencies are expected for a given background retention.
- **DiMuonLowMass**: requires a two track combination of displaced tracks with

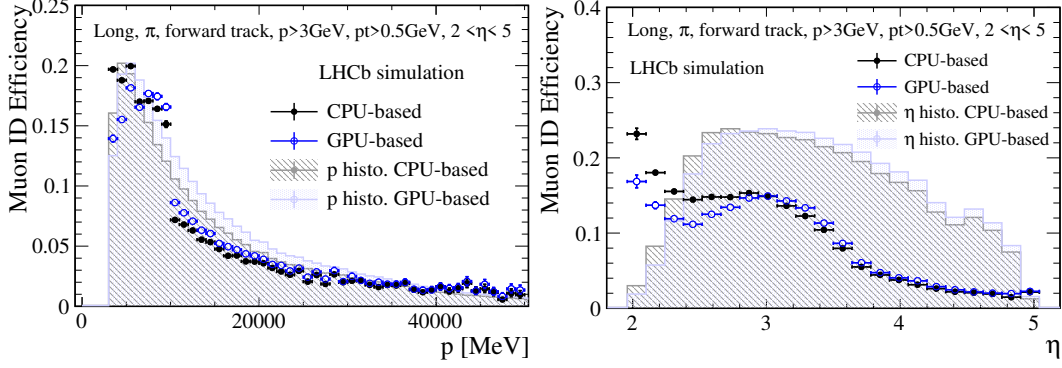


Figure 10: Efficiency to falsely identify true pions reconstructed as long muon tracks as muons as function of the pion momentum and pseudorapidity.

Signal	TrackMVA		TwoTrackMVA	
	GPU	CPU	GPU	CPU
$B_s^0 \rightarrow \phi\phi$	0.340(14)	0.332(14)	0.606(15)	0.621(15)
$J/\psi \rightarrow \mu^+\mu^-$	0.034(4)	0.031(3)	0.049(4)	0.042(4)
$B^0 \rightarrow K^{*0}e^+e^-$	0.276(10)	0.278(10)	0.439(12)	0.473(12)
$B^0 \rightarrow K^{*0}\mu^+\mu^-$	0.391(10)	0.385(10)	0.554(10)	0.582(10)
$D_s \rightarrow KK\pi$	0.076(5)	0.073(5)	0.178(8)	0.193(8)
$Z \rightarrow \mu^+\mu^-$	0.051(6)	0.040(6)	0.024(4)	0.028(5)

Table 2: Comparison of trigger efficiencies integrated over the kinematic phase space of the candidates, for each of the six simulated signal samples for the **TrackMVA** and **TwoTrackMVA** selections. Statistical uncertainties are indicated in parentheses.

very low momentum and transverse momentum. This trigger line is implemented identically for the CPU- and GPU-based HLT1.

- **SingleHighPtMuon**: selects tracks which are identified as a muon and fulfill a minimum momentum and transverse momentum requirement. This trigger line is also implemented identically for both architectures.

In Fig. 11 and 12, the trigger efficiencies of the CPU- and GPU-based HLT1 implementations are plotted against the parent transverse momentum, p_T , and parent decay time, τ (where applicable), respectively, for four signal samples. For ease of interpretation, only the efficiency of one suitably chosen trigger line per sample is shown. The trigger efficiencies of the two implementations are found to be comparable. In Tab. 2 and 3 the trigger efficiencies, integrated across the kinematic phase space of the samples, are compared for all four selections and various simulated signal samples. They are found to be comparable for the GPU and CPU-based implementation.

5.6 HLT1 rates

HLT1 rates are calculated in a similar way to the HLT1 efficiencies in the previous subsection. Both the CPU- and GPU-based implementation are run over 10k minimum bias events, and the positive decisions on each of the four selections defined in the

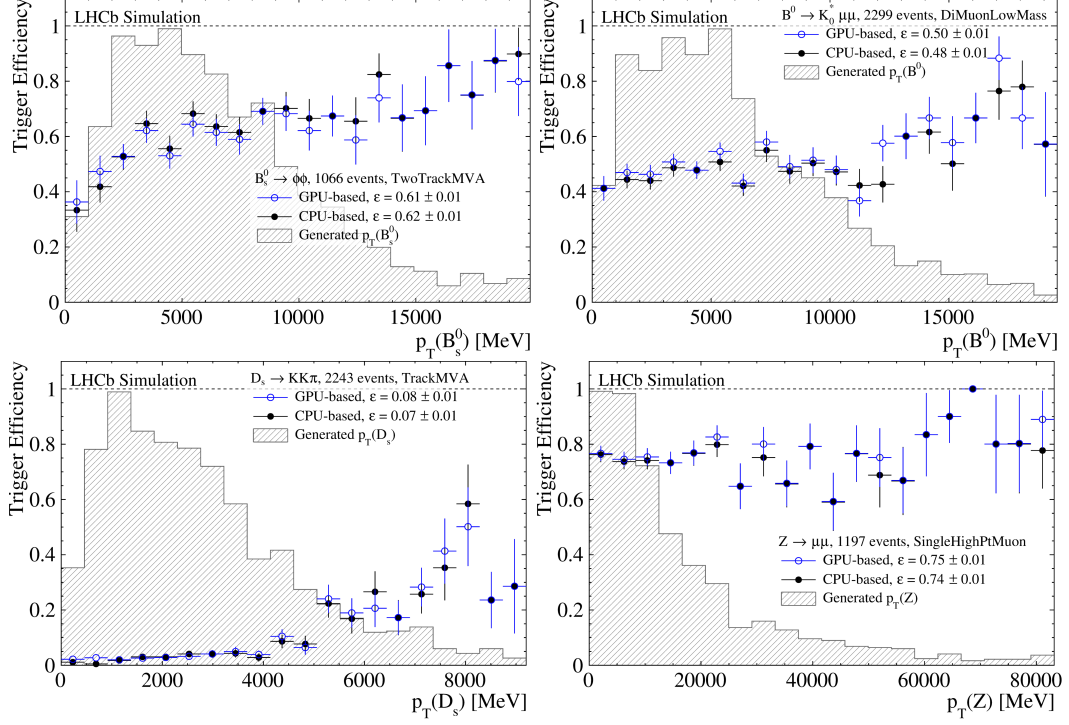


Figure 11: Trigger efficiencies for CPU-based and GPU-based HLT1 as a function of parent transverse momentum. Results are shown for the `TwoTrackMVA` (top left), `DiMuonLowMass` (top right), `TrackMVA` (bottom left) and `SingleHighPtMuon` (bottom right) selections firing on the $B_s^0 \rightarrow \phi\phi$, $B^0 \rightarrow K^{*0}\mu^+\mu^-$, $D_s \rightarrow KK\pi$ and $Z \rightarrow \mu^+\mu^-$ signal samples, respectively. The generated parent transverse momentum distribution is also shown for all events passing the denominator requirement.

Signal	DiMuonLowMass		SingleHighPtMuon	
	GPU	CPU	GPU	CPU
$B_s^0 \rightarrow \phi\phi$	0.025(5)	0.024(5)	0.005(2)	0.004(2)
$J/\psi \rightarrow \mu^+\mu^-$	0.078(5)	0.067(5)	0.048(4)	0.045(4)
$B^0 \rightarrow K^{*0}e^+e^-$	0.024(4)	0.027(4)	0.0011(8)	0.0011(8)
$B^0 \rightarrow K^{*0}\mu^+\mu^-$	0.502(10)	0.482(10)	0.091(6)	0.088(6)
$D_s \rightarrow KK\pi$	0.018(3)	0.019(3)	0.0013(7)	0.0013(7)
$Z \rightarrow \mu^+\mu^-$	0.033(5)	0.036(5)	0.749(12)	0.740(13)

Table 3: Comparison of trigger efficiencies integrated over the kinematic phase space of the candidates, for each of the six MC signal samples and the `DiMuonLowMass` and `SingleHighPtMuon` selections. Statistical uncertainties are indicated in parentheses.

previous subsection are counted. The rate for each line is defined as the number of events that fire that line, divided by the number of minimum bias events that are sampled, multiplied by the LHCb non-empty bunch crossing rate (30 MHz). Note that, contrary to the efficiency studies on MC signal samples of the previous subsection, no preselection (or “denominator” requirement) is applied. The comparison of the HLT1 rates is shown in Fig. 13. The `TrackMVA`, `DiMuonLowMass` and `SingleHighPtMuon` selections have comparable rates, although there is a discrepancy in the rates of the respective

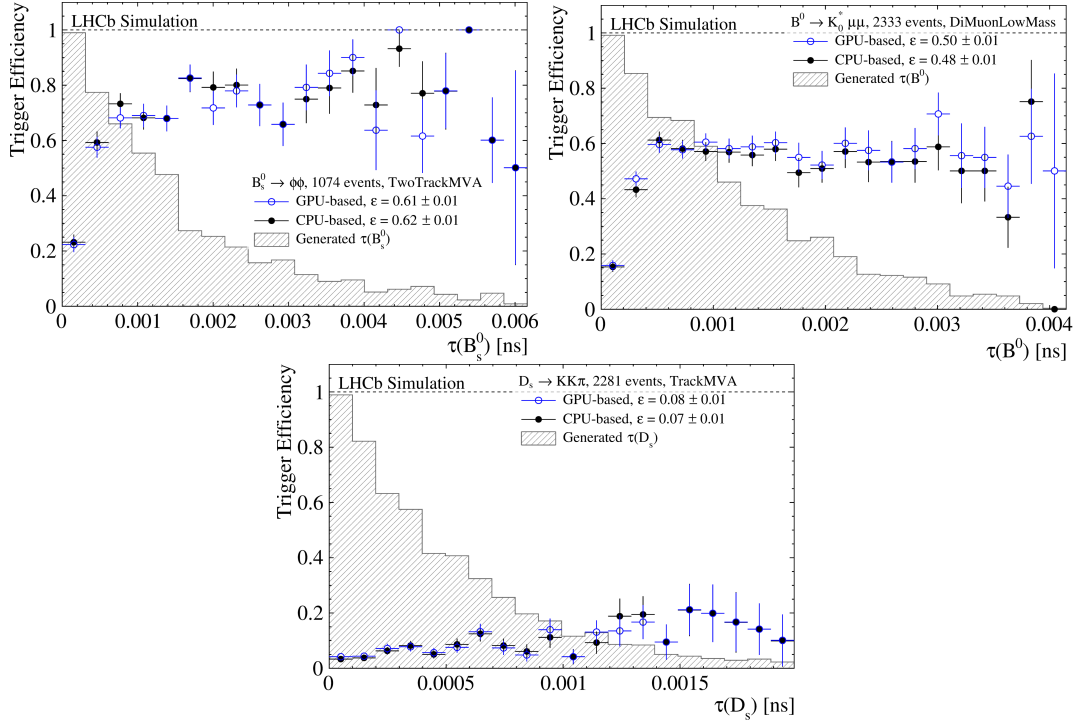


Figure 12: Trigger efficiencies for CPU-based and GPU-based HLT1 as a function of parent decay time. Results are shown for the `TwoTrackMVA` (top left), `DiMuonLowMass` (top right) and `TrackMVA` (bottom) selections firing on the $B_s^0 \rightarrow \phi\phi$, $B^0 \rightarrow K^{*0}\mu^+\mu^-$ and $D_s \rightarrow KK\pi$ signal samples, respectively. The generated parent decay time distribution is also shown for all events passing the denominator requirement.

`TwoTrackMVA` selections. This discrepancy can be explained by the different implementation of this line across the two projects, as detailed in the previous subsection. The inclusive rate for these four selections is found to be 912 ± 52 kHz for the GPU-based implementation, and 798 ± 48 kHz for the CPU-based one, largely due to the different implementation of the `TwoTrackMVA` trigger line. These numbers are well within the requirements of HLT1 to output between 1 – 2 MHz of events for further processing.

5.7 Global Event Cut efficiency

Both HLT1 implementations apply an identical global event cut (GEC) requiring fewer than 9750 total SciFi and UT clusters. Consequently, the GEC efficiencies for the GPU- and CPU-based implementation are found to be identical across all samples. The integrated efficiency of $Z \rightarrow \mu^+\mu^-$ is found to be 0.75 ± 0.01 . The other B and D decay samples under study have GEC efficiencies of about 85%, with statistical uncertainties of $\sim 1\%$. The efficiency on minimum bias events is 0.931 ± 0.003 .

5.8 Summary of the physics performance

The GPU- and CPU-based implementations presented here result in close to identical performance in all aspects. The observed differences are considered to be a matter of tuning the algorithms to balance between efficiency and fake rate, or misidentification rate,

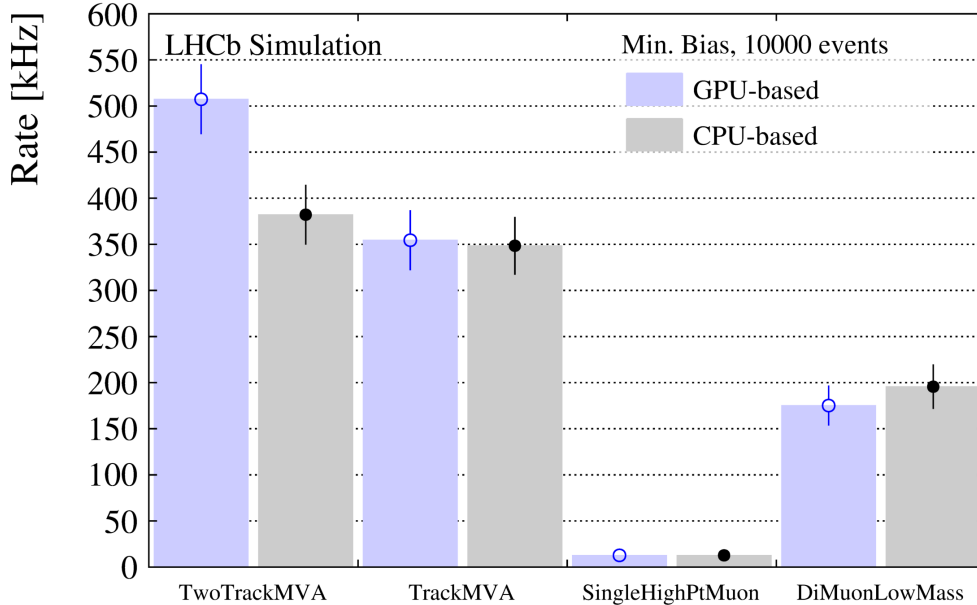


Figure 13: Trigger rates of the CPU- and GPU-based implementation for the four trigger selections of interest. The difference in the rate of the TwoTrackMVA selections can be explained by the differing implementation, which is detailed in Section 5.5.

and are expected to be neutral in terms of throughput. We therefore conclude that only the economic costs and the costs in terms of developer time need to be further considered in the cost-benefit calculation for the two architectures.

6 Cost benefit analysis

In this section the information of the previous sections is brought together to estimate the cost required for each architecture to process 30 MHz of data at the nominal instantaneous luminosity of $2 \cdot 10^{33} \text{ cm}^{-2}\text{s}^{-1}$. As computing prices are in rapid flux, and any real purchase would be subject to tendering, we give here relative “costs” in units of the reference “Quanta” CPU server node used for the HLT during Run 2 datataking.

6.1 Nominal cost-benefit

The nominal cost-benefit analysis is based on the assumption that the LHC will run at 30 MHz with full luminosity from the start of data taking. It is summarised in Tab. 4. The difference in capital expenditure costs for the CPU-only and the hybrid scenario is the sum of the difference in HLT1 compute costs and the HLT1 network costs. For the HLT1 compute costs of the CPU-only scenario, it should be noted that additional usage of the EFF is foreseen, namely for processing simulation when the LHC is in a technical stop or end-of-year shutdown (“out-of-data-taking”) and for processing HLT2 when the LHC is between fills (“out-of-fill”), respectively. While the out-of-fill use of the CPU nodes for HLT2 processing is taken into account in the following studies, the usage for simulation out-of-data taking is not.

Item	CPU-only	hybrid	Difference
Event Builder nodes	1000	1000	0
HLT1 network	275	25	250
HLT1 compute	450	125	325
Storage for 1 MHz output	575	575	0
Sub-total	2300	1725	575
Storage add. cost 2 MHz output	575	575	0
Total	2875	2300	575

Table 4: Indicative overall cost of the HLT1 implementations including contingency in units of the reference “Quanta” CPU server node used for the HLT during Run 2 data-taking. Numbers have been rounded to reflect inevitable order(10%) fluctuations in real-world costs depending on the context of any given purchase.

6.2 Disk buffer simulation

Based on the boundary conditions outlined in Section 2, studies of the expected maximum HLT1 output rate that can be processed within the budget envelope have been performed. The studies are performed as follows:

1. Fill lengths and inter-fill gaps are randomly sampled from 2018 data, with three machine development and technical stop periods and an average machine efficiency of 50%;
2. The fills and gaps are grouped into pairs before sampling to capture potential correlations between them;
3. A thousand such toys are randomly generated, and are shown to have a residual distribution compatible with that of the 2018 machine efficiency and a width of $\pm 1.7\%$;
4. The disk buffer required to ensure sufficient space for 95% of these 1000 toys is then determined using as input an event size of 100 kB, a chosen HLT1 output rate, and a chosen HLT2 throughput in and out of fill;
5. These values are scanned, generating 1000 toys per datapoint, over a range of HLT1 output rates and HLT2 throughput rates resulting in a 3 dimensional distribution of required disk buffer as a function of HLT1 and HLT2 rates.

For each datapoint in the distribution, the cost of the HLT2 throughput and cost of the disk buffer are determined. For combined costs greater than that of the overall budget in the CPU-only and hybrid scenarios, the datapoint is rejected. This leaves a distribution of valid points for which LHCb could purchase the necessary resources. The optimal working point is the one which maximises the HLT1 output rate. The inputs to this procedure are summarized in Tab. 6 and described in more detail in the following sections.

The simulation studies assume that in both the CPU-only and hybrid scenarios, any remaining budget after attaining a 30 MHz throughput at HLT1 and a sufficient buffer is used to buy CPU to provide additional HLT2 throughput. This corresponds to an increased HLT1 output rate allowing more, or more efficient selections at HLT1.

6.2.1 Quantification of HLT1 and HLT2 throughputs

Node type	HLT1 Throughput/node	HLT2 Throughput/node	minimum nodes
AMD 7502	171 kHz	471 Hz	175
Quanta	n/a	134 Hz	2050

Table 5: Summary of HLT1 and HLT2 throughputs for the EFF nodes described in Sec. 4. A minimum of 175 AMD 7502 nodes would be required to sustain 30 MHz of HLT1 throughput, corresponding to 82.5 kHz of out-of-fill HLT2 throughput. This is supplemented by 2050 Quanta equivalent nodes which are only used for HLT2 in-fill and out-of-fill, providing an additional 275kHz of HLT2 throughput.

The planned EFF nodes are equivalent to dual-socket AMD 7502 servers. The HLT1 throughput for this node was given in Section 4 requiring 175 nodes to sustain 30 MHz. As described in Section 4 these nodes are also capable of an HLT2 throughput of 471 Hz. The baseline scenario then results in an out-of-fill HLT2 throughput of 82.5 kHz and the remaining funds can be used to purchase disk or additional CPU. In the GPU scenario the purchase of these nodes for HLT1 is not necessary, but as a result no additional out-of-fill rate is available for HLT2 beyond the free resources described in the next subsection. The throughputs and quantities of these nodes before cost-benefit optimisation are presented in Table 5.

6.2.2 Additional free computing resources for HLT2

The Run 2 EFF Quanta nodes currently sustain a throughput of 134 Hz for HLT2, as documented in Section 4. In addition to the most recently purchased Run 2 EFF nodes, older nodes are available which correspond to 1450 quanta-equivalents, with up to 1200 additional nodes. These nodes will only be used for HLT2. In this study we take the 1450 nodes and assume a conservative 50% of the additional nodes, totalling 2050 quanta equivalents, or 275 kHz of in- and out-of-fill HLT2 throughput.

6.2.3 Results

The two scenarios are scanned according to the method described in Sec. 6.2. The results are shown in Fig. 14. In these figures, any combination of additional CPU and disk buffer that is cheaper than the total budget for the two scenarios is shown, defining a region of affordability, with the necessary disk buffer capacity represented by the z-axis colour. Unsurprisingly the maximum HLT1 throughput that can be sustained arises when the buffer is fully used and the remaining resources are spent exclusively on HLT2.

Scenario	CPU-only	hybrid
Maximum usable disk (PB)	28	28
In-fill min. HLT2 rate (kHz)	275	275
Out-of-fill min. HLT2 rate (kHz)	358	275

Table 6: Input parameters to the toy studies used to determine the maximum HLT1 output rate in each scenario.

Scenario	Maximum HLT1 output rate (kHz)
CPU-only	660
hybrid	790

Table 7: Results of the toy studies indicating the maximum affordable throughput in the both scenarios. The disk buffer is used to its maximum of 28 PB.

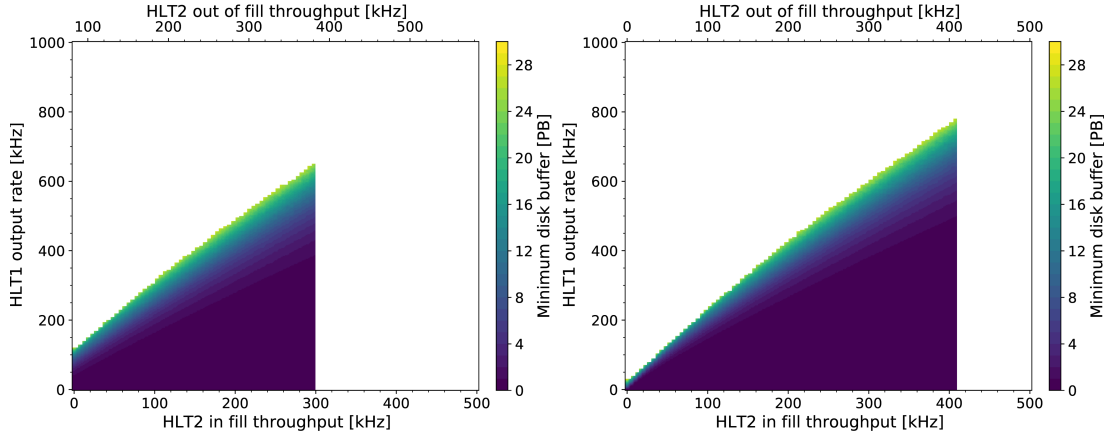


Figure 14: Region of affordability for the CPU-only scenario (left) and the hybrid scenario (right). The maximum affordable HLT1 throughput is 660 kHz and 790 kHz respectively.

The maximum HLT1 sustainable throughput in these scenarios is provided in Tab. 7, corresponding to the maximum y-axis extent of the region of affordability.

6.2.4 Results assuming a factor two performance increase in HLT2

HLT2 has not been optimised to the same extent as HLT1, therefore it is expected to improve its performance. To highlight the importance of doing so, the previous studies are repeated assuming a factor two improvement in the HLT2 throughput using the same architectures. This is motivated by preliminary studies using optimisations to the reconstruction sequence that are pending validation.

The two scenarios with double the HLT2 performance result in the regions of affordability shown in Fig. 15. As before, the maximum HLT1 throughput that can be sustained arises when the buffer is fully used and the remaining resources are spent exclusively on HLT2. The maximum HLT1 sustainable throughput in these scenarios is provided in Tab. 8.

Scenario	Maximum HLT1 output rate (kHz)
CPU-only	1185
hybrid	1425

Table 8: Results of the toy studies indicating the maximum affordable throughput in the four scenarios in which HLT2 throughput has increased by a factor of two. In each case the maximum throughput arises when the full remaining spend is used on HLT2 and the disk buffer is used to its maximum of 28 PB respectively.

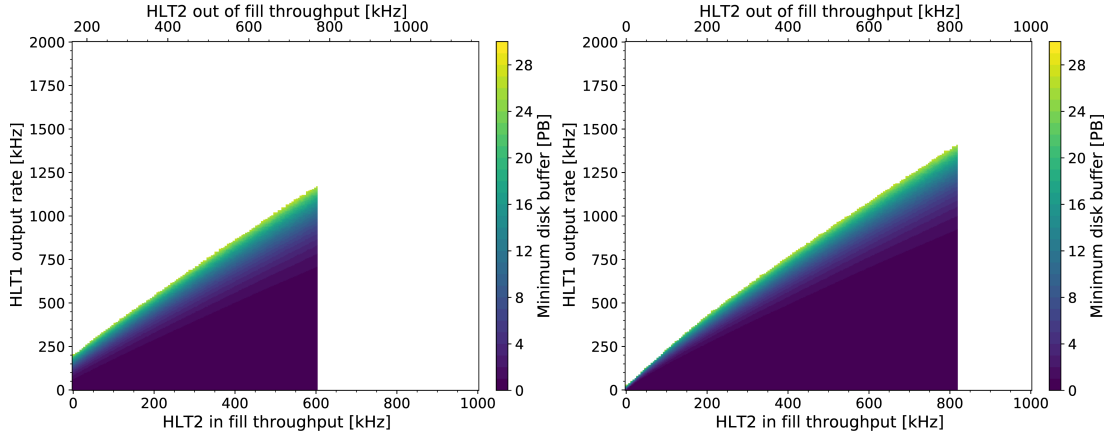


Figure 15: Region of affordability for the CPU-only scenario (left) and the hybrid scenario (right) when HLT2 throughput is two times higher than in the current best performance. The maximum affordable HLT1 throughput is 1185 kHz and 1425 kHz, respectively.

6.3 Deferred purchasing

It is expected that the first year of Run 3 will be commissioning year, with stable pp running at lower luminosities than nominal early in the year. This is expected to be followed by a period of running at nominal luminosity. There is general agreement that if possible hardware purchasing should be deferred to make best use of the collaboration’s financial resources. For this reason, the plan is to buy a system for 2022 which can handle half the expected nominal processing load. As the throughput of both the considered HLT1 architectures scales linearly (or worse) with detector occupancy, this implies that buying half the number of HLT1 processing units is sufficient. Many of the relevant costs from Tab. 4 can therefore be divided by two.

In the case of the CPU-only architecture, the disk storage required for 1 MHz output involves a great deal of additional mechanical and cabling work if purchased in two parts, so it is preferable to buy the full amount already for the commissioning year. However it is technically possible to divide this purchase also, in case there are good financial reasons to expect that it saves the collaboration a significant amount of money. Similar arguments also apply to the network between the EB and the CPU HLT1.

In the hybrid case, deferred purchasing allows HLT1 to be implemented with one GPU card per EB node in the first year. This results in maximum flexibility to benefit from the next generation of GPU cards which are expected to be released over the next two years, and mitigates the risk arising from the limited number of EB slots available to implement the GPU HLT1.

7 Conclusion

LHCb has investigated two viable proposals for its first level trigger to be deployed for Run 3 onwards. They are characterised by similar performances: The choice between x86-based and GPU-based technologies is far from obvious when viewed solely in terms

of raw throughput and physics performance, and as a result the decision over which option to use has come down to a detailed cost/benefit assessment, presented in this paper.

In addition to certain immediate cost savings the GPU-based solution offers more opportunities for future performance gains, as a more forward-looking longer-term investment; it is anticipated that GPU performance will advance at a more rapid pace, as evidenced by the new generation of GPU cards released while this paper was being prepared. This increased growth rate presents an opportunity to do more physics in Run 3 and beyond; and provides LHCb with the expertise and increased flexibility to exploit fully heterogeneous computing environments in future upgrades. The overhead of maintaining different technologies on HLT1 and HLT2 is not negligible and has been taken into account in the decision making process as well. On balance however, as a result of all these considerations, the collaboration has ultimately decided to implement a Hybrid HLT with a GPU-based first stage, in Run 3. The studies presented in this paper show that GPUs are increasingly a viable alternative to CPUs as general purpose, quasi-standalone processors.

Data availability statement

The datasets generated during and/or analysed during the current study are not publicly available as LHCb collaboration policy is to make data fully available 5 years after it has been analysed. They are available from the corresponding author on reasonable request.

Conflict of interest statement

No author of this paper has a financial or personal relationship with a third party whose interests could be positively or negatively influenced by the article's content.

Acknowledgements

This version of the article has been accepted for publication, after peer review, but is not the Version of Record and does not reflect post-acceptance improvements, or any corrections. The Version of Record is available online at: <https://doi.org/10.1007/s41781-021-00070-2>.

We thank the technical and administrative staff at the LHCb institutes. We acknowledge support from CERN and from the national agencies: CAPES, CNPq, FAPERJ and FINEP (Brazil); MOST and NSFC (China); CNRS/IN2P3 (France); BMBF, DFG and MPG (Germany); INFN (Italy); NWO (Netherlands); MEiN and NCN UMO-2018/31/B/ST2/03998 (Poland); MEN/IFA (Romania); MSHE (Russia); MICINN (Spain); SNSF and SER (Switzerland); NASU (Ukraine); STFC (United Kingdom); DOE NP and NSF (USA). We acknowledge the computing resources that are provided by CERN, IN2P3 (France), KIT and DESY (Germany), INFN (Italy), SURF (Netherlands), PIC (Spain), GridPP (United Kingdom), RRCKI and Yandex LLC (Russia), CSCS (Switzerland), IFIN-HH (Romania), CBPF (Brazil), PL-GRID (Poland) and NERSC (USA). We are indebted

to the communities behind the multiple open-source software packages on which we depend. Individual groups or members have received support from ARC and ARDC (Australia); AvH Foundation (Germany); EPLANET, Marie Skłodowska-Curie Actions and ERC (European Union); A*MIDEX, ANR, IPhU and Labex P2IO, and Région Auvergne-Rhône-Alpes (France); Key Research Program of Frontier Sciences of CAS, CAS PIFI, CAS CCEPP, Fundamental Research Funds for the Central Universities, and Sci. & Tech. Program of Guangzhou (China); RFBR, RSF and Yandex LLC (Russia); GVA, XuntaGal and GENCAT (Spain); the Leverhulme Trust, the Royal Society and UKRI (United Kingdom). The authors would like to thank the LHCb computing and simulation teams for their support and for producing the simulated LHCb samples used to benchmark the performance of RTA software.

References

- [1] LHCb collaboration, A. A. Alves Jr. *et al.*, *The LHCb detector at the LHC*, JINST **3** (2008) S08005.
- [2] R. Aaij *et al.*, *Performance of the LHCb trigger and full real-time reconstruction in Run 2 of the LHC*, JINST **14** (2019) P04013, [arXiv:1812.10790](#).
- [3] S. Akar *et al.*, *Review Document: Low Level Trigger (LLT)*, , CERN, Geneva, 2014.
- [4] ATLAS Collaboration, *Technical Design Report for the Phase-II Upgrade of the ATLAS TDAQ System*, , CERN, Geneva, 2017.
- [5] CMS Collaboration, *The Phase-2 Upgrade of the CMS DAQ Interim Technical Design Report*, , CERN, Geneva, 2017. This is the CMS Interim TDR devoted to the upgrade of the CMS DAQ in view of the HL-LHC running, as approved by the LHCC.
- [6] P. Buncic, M. Krzewicki, and P. Vande Vyvre, *Technical Design Report for the Upgrade of the Online-Offline Computing System*, , 2015.
- [7] LHCb collaboration, *LHCb Trigger and Online Technical Design Report*, CERN-LHCC-2014-016, 2014.
- [8] R. Aaij *et al.*, *Allen: A high level trigger on GPUs for LHCb*, *Comput. Softw. Big Sci.* **4** (2020) 7, [arXiv:1912.09161](#).
- [9] LHCb collaboration, *LHCb Upgrade GPU High Level Trigger Technical Design Report*, CERN-LHCC-2020-006, 2020.
- [10] for the ALICE Collaboration, S. Gorbunov *et al.*, *ALICE HLT high speed tracking on GPU*, *IEEE Trans. Nucl. Sci.* **58** (2011) 1845.
- [11] LHCb Collaboration, *Performance of the GPU HLT1 (Allen)*, LHCb-FIGURE-2020-014, 2020.
- [12] R. Aaij *et al.*, *Tesla: an application for real-time data analysis in High Energy Physics*, *Comput. Phys. Commun.* **208** (2016) 35, [arXiv:1604.05596](#).

- [13] R. Aaij *et al.*, *A comprehensive real-time analysis model at the LHCb experiment*, JINST **14** (2019) P04006, [arXiv:1903.01360](#).
- [14] LHCb collaboration, *Computing Model of the Upgrade LHCb experiment*, CERN-LHCC-2018-014, 2018.