**OPEN ACCESS**

# The LHCb Upgrade I

To cite this article: R. Aaij *et al* 2024 *JINST* **19** P05065

View the article online for updates and enhancements.

THE LARGE HADRON COLLIDER AND THE EXPERIMENTS FOR RUN 3 —
ACCELERATOR AND EXPERIMENTS FOR LHC RUN3

# The LHCb Upgrade I

**The LHCb collaboration**

*E-mail:* giovanni.passaleva@fi.infn.it

ABSTRACT: The LHCb upgrade represents a major change of the experiment. The detectors have been almost completely renewed to allow running at an instantaneous luminosity five times larger than that of the previous running periods. Readout of all detectors into an all-software trigger is central to the new design, facilitating the reconstruction of events at the maximum LHC interaction rate, and their selection in real time. The experiment's tracking system has been completely upgraded with a new pixel vertex detector, a silicon tracker upstream of the dipole magnet and three scintillating fibre tracking stations downstream of the magnet. The whole photon detection system of the RICH detectors has been renewed and the readout electronics of the calorimeter and muon systems have been fully overhauled. The first stage of the all-software trigger is implemented on a GPU farm. The output of the trigger provides a combination of totally reconstructed physics objects, such as tracks and vertices, ready for final analysis, and of entire events which need further offline reprocessing. This scheme required a complete revision of the computing model and rewriting of the experiment's software.

# Contents

# 1 Introduction

The LHCb experiment [1] is one of the four large detectors at the Large Hadron Collider (LHC) accelerator at CERN, and its primary purpose is to search for new physics through studies of CP-violation and decays of heavy-flavour hadrons. Although LHCb was designed primarily for precision measurements in heavy-flavour physics, the experiment has demonstrated excellent capabilities in many other domains ranging from electroweak physics to heavy ion and fixed target physics. The LHCb Upgrade experiment has been designed with this wider physics programme in mind as a general purpose experiment covering the forward region. LHCb has been successfully operated from 2010 to 2018 during the LHC Run 1 (2010–2012) and Run 2 (2015–2018) data-taking periods with excellent performance [2], collecting a total of $9\,\mathrm{fb}^{-1}$ of proton-proton ($pp$) data, about $30\,\mathrm{nb}^{-1}$ of lead-lead and $p$-lead collisions and about $200\,\mathrm{nb}^{-1}$ of fixed target data.

Notwithstanding this considerable data set, the precision on many of the key flavour physics observables studied and measured by LHCb remains statistically limited, as discussed in detail in ref. [3], thus requiring significantly larger data sets to probe the Standard Model at the level of precision achieved by theoretical calculations and obtain the required sensitivity to observe possible new physics effects.

While originally designed to take data at a maximum instantaneous luminosity $\mathcal{L} = 2 \times 10^{32}\,\mathrm{cm}^{-2}\,\mathrm{s}^{-1}$ to keep the average number of visible primary $pp$ interactions (*pile-up*) close to unity [4], LHCb has been successfully operated for most of Run 1 and Run 2 at $\mathcal{L} \sim 4 \times 10^{32}\,\mathrm{cm}^{-2}\,\mathrm{s}^{-1}$, demonstrating the capability to run and to produce excellent physics results at higher luminosity and with a pile-up larger than initially foreseen. A proposal for a major upgrade to operate LHCb at substantially larger instantaneous luminosity than Run 1-2 was thus formalised in a Letter of

Intent [5] and detailed in a Framework TDR [6]. The physics motivations for this upgrade have been discussed in great detail in refs. [3, 5, 6], assuming an expected total luminosity of $\sim 50\,\mathrm{fb}^{-1}$ integrated by the end of LHC Run 4.

The LHCb Run 1-2 system design would not allow a significant increase in statistics, especially for fully hadronic final state decays, the main limitation coming from the maximum allowed output rate of the first trigger stage, the L0, implemented in hardware [7]. The simple inclusive selection criteria implemented in the L0 trigger stage, based essentially on particle transverse momentum, would result in an effective loss of efficiency with increasing luminosity, especially for the most abundant processes with hadrons in the final state, and in the saturation of the event yield, as clearly visible in the left panel of figure 1.



**Figure 1.** Left: relative trigger yields as a function of instantaneous luminosity, normalised to $\mathcal{L} = 2 \times 10^{32}\,\mathrm{cm}^{-2}\,\mathrm{s}^{-1}$. Reproduced from [5]. CC BY 4.0. Right: rate of decays reconstructed in the LHCb acceptance as a function of the cut in $p_\mathrm{T}$ of the decaying particle, for decay time $\tau > 0.2\,\mathrm{ps}$. Reproduced from [8]. CC BY 4.0.

In addition, inclusive flavour physics signals have relatively large cross sections and, at the upgrade luminosity, every event in the LHCb acceptance will contain on average two long-lived hadrons not containing heavy quarks [8, 9]. Therefore, simple cuts based on displaced vertices or on $p_\mathrm{T}$ would be either not effective in rejecting background or, once enough purity is reached, would amount to downscaling the signal as shown in the right panel of figure 1. Profiting from a higher luminosity to collect significantly more data is therefore only possible by removing the L0 trigger stage and introducing selections that are more discriminating than simple inclusive criteria. In particular a full-software trigger discriminating signal channels based on the full event reconstruction has been deemed essential for this strategy.

Based on these considerations the LHCb upgrade has been designed to run at a nominal instantaneous luminosity $\mathcal{L} = 2 \times 10^{33}\,\mathrm{cm}^{-2}\,\mathrm{s}^{-1}$ and to collect events at the LHC crossing rate of 40 MHz. The events are discriminated by an all-software trigger reconstructing in real time all events at the visible interaction rate of $\sim 30$ MHz. By increasing the instantaneous luminosity by a factor of five and improving the trigger efficiency for most modes by a factor of two [9], the annual yields in most channels will be an order of magnitude larger than for the previous LHCb experiment. A total integrated luminosity (including Run 1 and runtwo) of around $50\,\mathrm{fb}^{-1}$ is expected by the end of Run 4 of the LHC.

The new trigger strategy, the higher luminosity and correspondingly higher pile-up required a complete renewal of the LHCb detectors and readout electronics that are now able to read events at

the 40 MHz LHC bunch crossing rate and cope with the larger event multiplicity thanks to a higher granularity. A full revision of the experiment's software and of the data processing and computing strategy was also necessary to deal with the expected large increase in data volume.

This paper describes the design and construction of the upgraded LHCb experiment providing details on all the new subdetectors, on the trigger and online systems and on the software and data processing frameworks.

# 2 The LHCb detector

## 2.1 Detector layout

LHCb is a single-arm forward spectrometer covering the pseudorapidity range $2 < \eta < 5$, located at interaction point number 8 on the LHC ring. Figure 2 shows the layout of the upgraded detector. The coordinate system used throughout this paper has the origin at the nominal $pp$ interaction point, the $z$ axis along the beam pointing towards the muon system, the $y$ axis pointing vertically upward and the $x$ axis defining a right-handed system. Most of the subdetector elements (with the notable exception of vertex and Cherenkov detectors) are split into two mechanically independent halves (the *access side* or *Side A* at $x > 0$ and the *cryogenic side* or *Side C* at $x < 0$), which can be opened for maintenance and to guarantee access to the beam pipe.
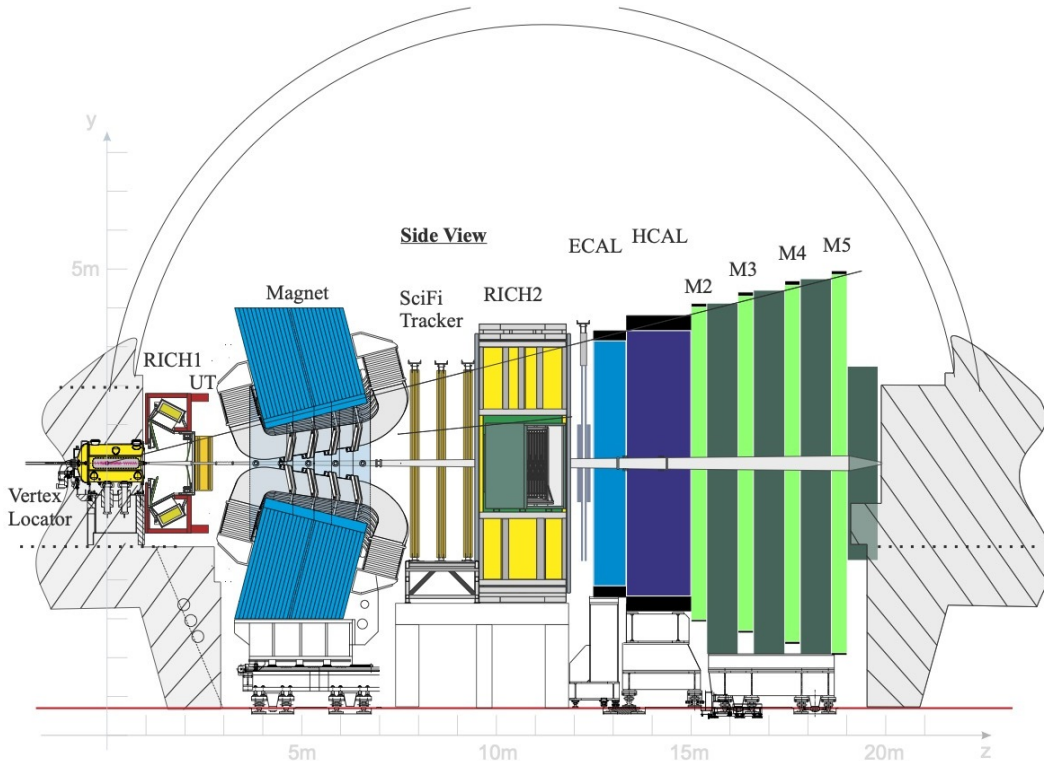


**Figure 2.** Layout of the upgraded LHCb detector.

The particle tracking system comprises an array of pixel silicon detectors surrounding the interaction region called vertex locator (VELO), the silicon-strip upstream tracker (UT) in front of the large-aperture dipole magnet, and three scintillating fibre tracker (SciFi Tracker) stations downstream

of the magnet.[1] All three subsystems were designed to comply with the 40 MHz readout architecture and to address the challenges associated with the increased luminosity. The upgraded VELO, based on hybrid silicon pixel detectors, is described in section 3, and the UT is described in section 5. The SciFi Tracker, which replaces both the straw-tube Outer Tracker and silicon-strip Inner Tracker systems used in the downstream tracking stations in the original LHCb experiment, is described in section 6.

The particle identification (PID) is provided by two ring imaging Cherenkov detectors (RICH1 and RICH2) using $C_4F_{10}$ and $CF_4$ gases as radiators, a shashlik-type electromagnetic calorimeter (ECAL), an iron-scintillator tile sampling hadronic calorimeter (HCAL), and four stations of muon chambers (M2–5) interleaved with iron shielding.[2] The Scintillating Pad Detector and Pre-Shower, which were part of the previous calorimeter system, as well as the most upstream muon station, have been removed due to their reduced role in the full software trigger compared to the former hardware L0. The upgraded ring imaging Cherenkov detectors (RICHs) are described in section 7, the calorimeters are described in section 8, and the muon system is described in section 9.

The data acquisition system (DAQ) comprises the front-end (FE) and back-end (BE) electronics connected by long-distance optical links, the event-builder and the event-filter farms, both described in section 10.

## 2.2 Magnet

The spectrometer's dipole magnet has been maintained unchanged with respect to Run 1-2. It provides a vertical magnetic field with a bending power of $\simeq 4$ Tm. It consists of two identical, saddle-shaped coils, which are mounted symmetrically inside a window-frame yoke. To match the detector acceptance, the pole gap increases both vertically and horizontally towards the downstream tracking stations. Detailed descriptions of the magnet design can be found in refs. [10–12].

Each coil is made from five triplets of aluminium pancakes and is supported by cast aluminium clamps fixed to the yoke.

The initial magnetic field map was determined based on a set of measurement campaigns (prior to Run 1) complemented by finite-element simulations. Subsequent measurements for limited regions inside the magnet, were carried out in 2011, 2014 and 2021, and were used to apply corrections to the field map.

During data taking, the magnet polarity is reversed regularly (every few weeks) to collect data sets of roughly equal size with the two field configurations.

## 2.3 Electronics architecture

The architecture of LHCb Upgrade I is designed to transmit data collected from every bunch crossing all the way to the event-builder computing farm. To implement this architecture, LHCb maximised the use of common building blocks to benefit from a unified approach. Common developments for the LHC experiment upgrades, such as radiation-tolerant optical links, have proven to be vital enabling technologies for the new LHCb detector. The general architecture is shown in figure 3. The FE electronics amplify and shape the signals generated within the particle detectors.

---

[1]Upstream and downstream are intended in the direction of increasing $z$.

[2]The muon detector consisted of five stations of which the first (M1) has been removed — see text. For historical reasons the remaining stations kept their original names.

**Figure 3.** Electronics architecture of the upgraded LHCb experiment. Reproduced from [14]. CC BY 4.0.

These signals are digitised and optically transmitted off the detector. All components of the FE electronics are located on or close to the detector, and are therefore exposed to beam-induced radiation. Hence, they are all radiation tolerant by design and/or qualified for the radiation environment in which they operate. The BE electronics are situated in a data centre on the surface and are connected to the FE in the cavern by 250 m-long optical fibres. The BE electronics preprocess and format the data for transmission to the event builder. The data centre is a radiation-free environment and commercial components have been used for the implementation of the BE. Clocks and fast, beam-synchronous commands are distributed by a timing and fast signal control (TFC) system. The experiment control system (ECS) configures and monitors the BE and FE. The ECS implements also the slow controls like for example high voltage (HV), low voltage (LV) and temperature monitoring. The TFC and ECS systems are described in sections 10.3 and 10.4. LHCb requires a much larger bandwidth for data than for TFC and ECS signals. Even though the radiation-tolerant optical links were conceived to provide data transmission, TFC and ECS functionality in the same link, LHCb has separated these functions to maintain a strict independence between data and controls. This has allowed a minimisation of the number of links and construction of a modular system with clear boundaries between functions. Hence, data from the detectors are transmitted on dedicated unidirectional links whilst the TFC and ECS communications to the FE electronics are merged onto a much smaller number of separate bidirectional links. The FE electronics are a mixture of customised components for each subdetector and common components used across all systems. The analog electronics connected to the detectors have all been implemented as application specific integrated circuits (ASICs) and profit from the intrinsic radiation tolerance of deep-submicron CMOS technology. Data are digitised and compressed either in the ASIC or, if the radiation levels allow, in commercially available field-programmable gate arrays (FPGAs). All such FE digital electronics have been implemented to resist single event effects (SEEs) by using techniques such as triple-modular redundancy. Data compression was introduced into the architecture to minimise the number of data links, and algorithms ranging from simple zero-suppression to hit clustering have been successfully implemented by different detectors. However, compression offers little advantage when the channel occupancy is high and so has not been used in some regions of the experiment. The relatively modest radiation level in many parts of LHCb has

allowed the widespread use of FPGAs although only after careful qualification procedures. These have brought many advantages such as shorter design time and flexibility, as well as relaxing the demand for specialised personpower and tools required for ASIC implementations. Efficient data compression comes at the cost of variable latency, thus the data transport system in LHCb is asynchronous. To allow the proper reconstruction of event fragments downstream in the system, data packets are tagged as early as possible with a unique time-stamp based on the bunch-crossing identifier (BXID). Simplex data links, running at 4.8 Gbit/s, are constructed from the gigabit transceiver (GBT) serialiser/deserialiser ASIC [13] and the versatile twin transmitter (VTTx) opto-electrical converter [15]. The only exception is the VELO, where the serialiser is embedded directly in the FE ASIC. The duplex TFC and ECS links are constructed with the GBT, versatile link transceiver (VTRx) and slow control adapter (SCA) [16]. This standardisation across LHCb has had major benefits in easing both development and deployment, as well as the sharing of experience across the wider LHC community. Power distribution has followed a similar common approach with the use of FEASTMP DC-DC converters [17] for local power regulation. The specific implementations of the FE architecture by each subdetector are described in subsequent sections. The BE electronics consist of custom PCI-express modules mounted in PC servers in the data centre. This module, known generically as PCIe generic back-end board (PCIe40), contains arrays of optical transmitters and receivers connected to a powerful FPGA. The PCIe40 was conceived and designed to fulfil the functionality required for both data acquisition and controls. Hence, by the choice of the FPGA firmware, the module can be configured for either data acquisition or controls. The role of the PCIe40 board for data acquisition (TELL40) is to decode and process data, and then build multievent packets for transmission to the event-builder. The PCIe40 board for controls (SOL40) is the ECS interface used to configure the FE electronics and transmit TFC commands to the FE and TELL40s. The PCIe40 also plays the role of interface to the LHC machine timing when configured as a readout supervisor board (SODIN) board. More details on the PCIe40 and its functions are given in section 10.2.1.

## 2.4 Infrastructure

A significant part of the Run 1-2 LHCb infrastructure has been completely refurbished to comply with the upgraded detectors and modernise old equipment.

### 2.4.1 Power distribution

The power supplies providing low and high voltages are housed in electronics racks in counting rooms on Side A of the LHCb cavern, which is separated from the experimental area by a concrete shielding. During LHC long shutdown 2 (LS2), some optical fibres were removed to free space in the long-distance cable trays and additional copper cables, with a total length of ~ 36 km, were installed between the counting rooms and the detector.

Most of the LHCb equipment required for operating the experiment is fed from the same electrical network (the so-called *machine network*) used by the LHC. Equipment such as lighting and overhead cranes is fed from the separate, general services network. In case of an outage, a change-over from one network to the other can be performed. Systems with particularly stringent up-time requirements, such as the detector safety system (DSS) or the magnet safety system (MSS), are connected to an uninterruptible power supply (UPS) network, with back-up provided by a diesel generator. The electricity consumption of the upgraded experiment is dominated by the dipole magnet with 4.6 MW

of electric power and a yearly consumption of $\sim 20\,\mathrm{GWh}$. This is followed by the surface data centre, which draws $\sim 2\,\mathrm{MW}$, and the detectors' electronics and cooling systems which consume $\sim 200\,\mathrm{kW}$. The electrical infrastructure has been partially renewed for the upgrade. In particular, new power distribution lines have been installed for the upgraded cooling system (see section 2.4.3) and two new 18 kV high-voltage cells and two 3.15 MVA tranformers were added for the data centre at the LHCb site.

### 2.4.2 Neutron shielding

As explained in section 6, the performance of the silicon photomultiplier (SiPM) arrays used as photon detectors in the SciFi Tracker is significantly impacted by radiation damage. The dominant contribution to the fluence at the location of the SiPM arrays comes from high-energy neutrons produced in showers in the calorimeter. During LS2, a dedicated neutron shielding has been installed upstream of the calorimeter, taking the space formerly occupied by muon station M1 and reusing its support structure. The shielding, made from polyethylene ($C_2H_4$) with a 5% admixture of boron, has an inner region ($2 \times 2\,\mathrm{m}^2$) with a thickness of 300 mm and an outer region ($5 \times 5\,\mathrm{m}^2$) with a thickness of 100 mm. From simulations with FLUKA [18, 19], it is expected to reduce the 1 MeV neutron equivalent ($n_{\mathrm{eq}}$) fluence at the location of the SciFi Tracker SiPMs by a factor 2.2–3.0. Polyethylene was selected as material for the shielding since it efficiently moderates fast neutrons by elastic scattering, while boron reduces the activation due to the resulting thermal neutrons because of its high cross-section for thermal neutron capture.

### 2.4.3 Detector cooling

In order to further mitigate radiation effects, the SciFi Tracker SiPM arrays will be kept at a temperature of $-40°C$ using a monophase liquid cooling system. Initially, the system has been operated with the perfluorocarbon $C_6F_{14}$ as cooling fluid. However, work is underway to validate alternative fluids which possess similar thermal properties and radiation tolerance as $C_6F_{14}$ but feature a significantly lower global warming potential such as the hydrofluoroether $C_4F_9OCH_3$ or the fluoroketone $C_6F_{12}O$.[3]

Both SciFi Tracker and RICH cooling plants have been constructed by CERN which is also responsible for the demineralised water cooling system used for the SciFi Tracker electronics.

The VELO and UT use evaporative $CO_2$ cooling for the thermal management of their silicon sensors and front-end ASICs. Two identical cooling plants based on the 2-phase accumulator controlled loop (2PACL) concept [20], and each having a cooling capacity of 7 kW at $-30°C$, have also been constructed by CERN.

The SciFi Tracker and VELO/UT cooling plants use a shared primary chiller, with a capacity of 24 kW at $-56°C$. All cooling plants are located in the LHCb cavern Side A and are connected to the detector via 50–80 m long transfer lines. Primary cooling is provided by the chilled water (at $\sim 6°C$) and mixed water (at $\sim 14°C$) circuits.

### 2.4.4 Data centre

A new modular data centre has been built on the surface of the experimental site to accommodate all the computing resources needed for the upgraded readout system and event-filter farm. The new data centre comprises six 22-rack modules with a total power capacity of $\sim 2\,\mathrm{MW}$ of computing

---

[3]Novec 7100™ and Novec 649™. These fluids have global warming potential (GWP) $\sim 1$ and $\sim 300$, respectively, compared to GWP $\sim 9300$ for $C_6F_{14}$.

**Figure 4.** The readout system located in the modular data centre and the front-end electronics in the underground cavern are connected through long-distance optical fibres installed in the PM85 shaft.

equipment. The two central modules house the event-builder servers, connected to the front-end electronics underground via 190 thousand OM3 multimode optical fibres over a length of $\sim 250$ m (figure 4). The remaining four modules host the servers of the event-filter farm. The data centre uses a highly efficient cooling system, based on a combination of indirect free air cooling and evaporative water cooling, with a power usage effectiveness smaller than 1.1.

## 2.5 Beam pipe

The vacuum beam pipe and its support structure have been optimised to reduce background occupancy in the nearby tracking detectors [1, 21]. The conical shape of the beam pipe in LHCb leads to unbalanced forces in the axial direction due to the atmospheric pressure, which must be counterbalanced with mechanical restraints. In the aperture of the dipole magnet a support system consisting originally of eight stainless steel wires and rods, provided enough stiffness in all transverse directions. The two wire systems were attached to aluminium collars, connected to the beam pipe by graphite-reinforced polyimide-based plastic[4] rings.

During Run 1 this support system was identified as a significant source of scattering in the experiment and was redesigned as part of the upgrade programme for the LHCb vacuum system. The improved support system was installed in 2014 during the LHC long shutdown 1 (LS1).

---

[4]Vespel$^{TM}$.

In the new support system, the volume of the collars has been reduced and materials with longer radiation length have been selected for all components. Carbon fibre reinforced plastic tubes and synthetic ropes have replaced the stainless steel rods and cables. This led to a remarkable increase of more than 90% in material transparency for these components, while retaining sufficient stiffness. The aluminium collars were redesigned and remade with beryllium, which has led to a material transparency improvement of more than 85% [22].[5] A range of innovative materials was selected: aramid[6] for the ropes; thermoset carbon high-modulus fibres[7] for the tubes; and polybenzimidazole[8] for interface rings. A thorough qualification process of these materials has been carried out by the CERN vacuum group considering mechanical strength, creep and radiation tolerance.

### 2.6 Background and luminosity monitors

A set of detectors has been developed to monitor the machine-induced background conditions around the interaction point. To reduce systematic uncertainties and facilitate the reconstruction, the instantaneous luminosity delivered to LHCb is kept constant throughout a fill using a procedure known as *luminosity levelling* [23].[9] Every few seconds during nominal operation, LHCb publishes a measurement of the average number of visible $pp$ interactions per beam-beam crossing, denoted by $\mu_{\text{vis}}$, which is used by the LHC control system to adjust the offset of the two beams in the direction perpendicular to the nominal crossing plane. As the number of visible $pp$ collisions per beam crossing follows a Poisson distribution, the average of the distribution (i.e. $\mu_{\text{vis}}$) can be determined by measuring the probability that a beam crossing has no observable activity in the detector, $P_0 = \exp(-\mu_{\text{vis}})$; this is the *logZero* method [24]. During Run 1 and Run 2, the real-time luminosity measurement was based on information available in the L0 hardware trigger. For the upgraded experiment, a dedicated luminosity subdetector, dubbed probe for luminosity measurement (PLUME), has been installed, see section 2.6.1.

Excursions in luminosity or elevated machine-induced background are not only detrimental to the quality of the collected data but can also cause damage to the detector. For the safe operation of the experiment, it is essential to quickly detect and react upon anomalous beam conditions. Such protection is ensured by the beam conditions monitor system (BCM), described in section 2.6.2. Additional monitoring of the beam environment is provided by a metal foil detector called the radiation monitoring system (RMS), see section 2.6.3. Figure 5 shows a picture of the RMS, PLUME, and the upstream BCM station as installed in the cavern.

### 2.6.1 PLUME

PLUME [25] is a luminometer measuring Cherenkov light produced by charged particles crossing a quartz radiator. Its basic detection element, shown in figure 6, is a photomultiplier tube (PMT)[10] with a 1.2 mm thick quartz entrance window and a photocathode with a diameter of 10 mm. To increase the amount of Cherenkov light, a 5 mm thick quartz tablet is placed in front of the PMT window.

---

[5]Here, the material transparency is defined as $I = t/X_0$ where $X_0$ is the radiation length of the material and $t$ is the thickness of material seen by particles when traversing the various supporting elements.

[6]Teijin Technora™.

[7]Torayca M46J™.

[8]Celazole PBI U-60™.

[9]In the LHC jargon, a *fill* is the full beam cycle from injection to beam dump. Experiments normally divide the data taking part of a fill in *runs*, which correspond to samples of data taken at constant conditions.

[10]Model R760 by Hamamatsu Photonics K.K.™, Hamamatsu City, Japan.

**Figure 5.** Left: RMS (under the mylar protection foil at the left side) and PLUME (inside the scaffolding). The arrow indicates the position of BCM which is hidden behind the PLUME scaffolding. Right: upstream BCM detectors inside their kapton-insulated support surrounding the beam pipe.



**Figure 6.** Schematic view of a PLUME elementary detection module. The module is 153 mm long and has a diameter of 24 mm. Reproduced from [25].. CC BY 4.0.

The detector, located upstream of the $pp$ interaction region (between $z = -1900$ mm and $z = -1680$ mm), is a hodoscope consisting of two stations, each of which comprises 24 PMT modules arranged in a star-shaped structure around the beam pipe. The modules are placed at radial distances between 157 and 276 mm with respect to the beam line (corresponding to a pseudorapidity range $2.4 < \eta < 3.1$) and are angled such that the PMT axes point to the nominal interaction point. The FE, based on components developed for the upgraded LHCb calorimeters, and the LED monitoring system (described below) are located at a distance of $\sim 20$ m from the PMTs to reduce the level of radiation to which they are exposed.

The detector design was optimised with simulations using PYTHIA8 [26] and GEANT4 [27, 28], and was validated using test beam measurements. A summary of the studies can be found in refs. [24, 25]. Figure 7 (left) shows a test setup used in a 5.4 GeV electron beam at DESY.[11] It consists of two PMT modules placed one behind the other in the beam line and a trigger scintillator at the rear. An example of the charge spectrum measured in the first PMT (operated at 1000 V) is shown in figure 7 (right). For charged-particle tracks that produced a signal in both PMT modules and the trigger scintillator, the average collected charge was 12 pC and a time resolution of 0.6 ns was found.

As is the case for the other LHCb detectors, PLUME is fully integrated in the ECS, DAQ and TFC systems. A particular feature of PLUME is that it needs to be read out even if the other LHCb

---

[11]Deutsches Elektronen-Synchrotron, Hamburg, Germany.

**Figure 7.** Left: PLUME test beam setup. Right: charge collected by the first PMT in the test beam for all events (black line) and events producing a simultaneous signal in the two PMTs and the trigger (red line). The scale is in nVs where 1 nVs = 20 pC.

detectors are off or the event-builder system is not available. Such circumstances occur during injection and focusing of the LHC beams, when feedback on the instantaneous luminosity is vital for the optimisation of collisions at the LHCb interaction point.

The online luminosity calculation is performed in the firmware of the BE electronics. For every pair of PMT modules, the PLUME TELL40 counts the total number of bunch-crossings 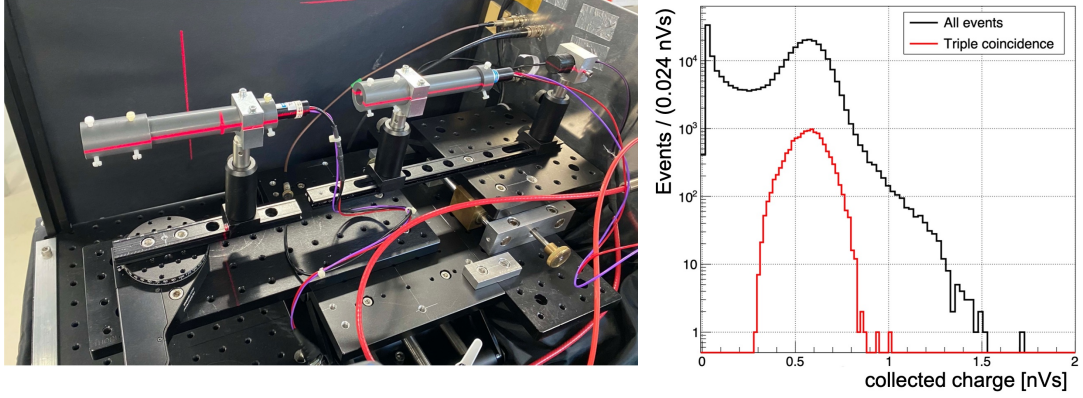$N$ and the number of those with a signal below threshold, $N_0$. Using the *logZero* method, the average number of visible interactions per bunch-crossing, $\mu_{\text{vis}}$, is given by [25]:

$$\mu_{\text{vis}} = -\log P_0 = -\log \frac{N_0}{N} - \frac{1}{2}\left(\frac{1}{N_0} - \frac{1}{N}\right), \tag{2.1}$$

where the second term accounts for second-order bias corrections to the Poisson statistics. In the 3 s interval between counter resets, the value of $\mu_{\text{vis}}$ is expected to be stable and deviations from the Poisson statistics of eq. (2.1) can be neglected. The estimated statistical uncertainty on the average luminosity is given by $5\%/\sqrt{n_{bb}}$, where $n_{bb}$ is the number of colliding bunch pairs, and becomes negligible for $n_{bb} \gg 1$.

The systematic uncertainty depends, among other things, on the stability of the PMT response which can change with time due to variations in temperature or occupancy, or because of ageing. The monitoring and calibration system, which was developed based on experience from the LHCb calorimeters [29] and the ATLAS LUCID detector [30], is therefore an integral part of PLUME. At regular time intervals (using suitable gaps in the LHC filling scheme), the LED calibration system, located next to the FE, sends light pulses over ∼ 20 m long quartz fibres to the front face of each PMT. The stability of the LED light pulses is monitored by PIN photodiodes located in the same rack as the LEDs. Finally, the degradation in the transparency of the quartz fibres due to radiation damage is monitored with dedicated fibres looped back to the LED position and read out by PMTs placed next to them. Based on the PMT response to the injected light, the high voltage of the PMTs is adjusted in steps of $\Delta V = 0.5$ V, which corresponds to a ∼ 2% change in gain. Tracks reconstructed in upstream VELO stations[12] and passing through PLUME can be used to cross-check the reliability of the calibration and monitoring system.

---

[12]These are VELO stations placed upstream of the interaction point.

### 2.6.2 BCM

The beam conditions monitor system (BCM) comprises two stations, one upstream (at $z = -2131$ mm) and one downstream (at $z = +2765$ mm) of the interaction region. A detailed description of the system can be found in ref. [31]. Each station consists of eight poly-crystalline chemical vapour deposition (pCVD) diamond pad sensors arranged symmetrically around the beam pipe. For each of the diamond sensors, the average current, integrated over periods of $40 \,\mu$s, is measured and a beam abort is requested through the LHC beam interlock system if three adjacent diamond sensors exhibit a current above threshold for two consecutive periods. In addition, running sums over 32 consecutive measurements are computed and a dump of the LHC beams is triggered if the average of the running sums in one station exceeds a given threshold. For monitoring purposes the BE also calculates average and maximum values over intervals of few seconds, which are read by the ECS and used for calculating the normalised background figures of merit, which are made available to the LHC control system. The system has been operating successfully throughout Run 1 and Run 2. During LS2 the diamond sensors were replaced, the support structures rebuilt, and the BE have been upgraded to be compatible with the new readout architecture.

### 2.6.3 RMS

The upgraded radiation monitoring system (RMS) consists of four metal-foil detector modules located upstream (at $z \sim -2200$ mm) of the nominal interaction point, at a radial distance of $\sim 30$ cm from the beam line. Each module houses two five-layer stacks of copper foils, with the central $50 \,\mu$m thick foil serving as the sensor. The detector concept exploits the phenomenon of secondary-electron emission at the metal surface due to charged particles crossing the foil. The readout electronics, located $\sim 80$ m away in the accessible part of the LHCb cavern, convert the resulting current to a frequency. The RMS is integrated in the ECS and the measurements are displayed in the LHCb control room. A similar system was used during Run 1 and Run 2 to monitor the charged particle fluence [32].

## 3 Vertex locator

### 3.1 Overview

The VELO detects tracks of ionising particles coming from the beam collision region and thereby measures the location of interaction vertices, displaced decay vertices and the distances between them. VELO tracks seed the reconstruction algorithm of the LHCb spectrometer and provide discriminatory information for event selection. The VELO has been redesigned [33] to be compatible with the luminosity increase and the trigger-less 40 MHz readout requirement of the upgraded experiment. It must continue to provide pattern recognition within an acceptable CPU budget, whilst maintaining the highest track-finding efficiency. The core technology of the new VELO is pixelated hybrid silicon detectors, which are arranged into modules and cooled by a silicon microchannel cooler. Of the mechanical structures, only the principal vacuum vessel and motion services remain from the version that was in operation until 2018. In particular, the RF boxes, the enclosures that interface the detector to the LHC beams, were entirely redesigned reducing both material and the inner radius of the VELO along the beam line. Furthermore, a new structure, a *storage cell*, is fitted immediately upstream of the VELO detector in the beam vacuum, see section 4. A summary of the changes is shown in table 1.

**Table 1.** Specifications of the upgraded VELO compared to those of the original version.

| | 2009–2018 | 2022 |
|---|---|---|
| RF box inner radius (minimum thickness) | 5.5 mm (300 µm) | 3.5 mm (150 µm) |
| Inner radius of active silicon detector | 8.2 mm | 5.1 mm |
| Total fluence (silicon tip) [$n_{eq}/\text{cm}^2$] | $4 \times 10^{14}$ | $\sim 8 \times 10^{15}$ |
| Sensor segmentation | $r - \phi$ strips | square pixels |
| Total active area of Si detectors | $0.22 \,\text{m}^2$ | $0.12 \,\text{m}^2$ |
| Pitch (strip or pixel) | 37–97 µm | 55 µm |
| Technology | n-on-n | n-on-p |
| Number of modules | 42 | 52 |
| Total number of channels | 172 thousand | 41 million |
| Readout rate [ MHz ] | 1, analogue | 40, zero suppressed |
| Whole-VELO data rate | 150 Gbit/s | $\sim 2$ Tbit/s |
| Total power dissipation (in vacuum) | 800 W | $\sim 2$ kW |

The combination of the pixel geometry, a smaller distance to the first measured point and reduced material means the performance of the VELO is significantly improved. However, the closer proximity to the LHC collisions and the step-change in design luminosity means the design must prepare for hit rates and radiation doses that are an order of magnitude higher than those experienced by the earlier VELO.

### 3.2 Design requirements

The principal metric for a vertex detector design is impact parameter resolution $\sigma_{\text{IP}}$, the precision with which the perpendicular distance of a track to a point is measured. This metric is a function of: track transverse-momentum, $p_{\text{T}}$; the average axial distance of the material before the second measurement, $r_1$; the distances from the point to the first and second measurements, $\Delta_i$ ($i = 1, 2$); and the position uncertainties of those measurements, $\sigma_i$. In the VELO case, it can be approximated as [34],

$$\sigma_{\text{IP}}^2 \approx \underbrace{\left(\frac{r_1}{p_{\text{T}}[\,\text{GeV}/c]}\right)^2 \left(0.0136\,\text{GeV}/c\sqrt{\frac{x}{X_0}}\left(1 + 0.038\ln\frac{x}{X_0}\right)\right)^2}_{\text{multiple scattering}} + \underbrace{\frac{\Delta_2^2\sigma_1^2 + \Delta_1^2\sigma_2^2}{\Delta_{12}^2}}_{\text{extrapolation}} \qquad (3.1)$$

where $x/X_0$ is the fraction of radiation length traversed before the second measurement. The first term describes the degradation induced by multiple scattering. The second term is the extrapolation error, which is dominated by detector geometry: pixel size and lever arm, $\Delta_{12}$, between the first and second measured points. The upgraded VELO design was optimised to achieve, within the nominal LHCb acceptance, a performance at least as good as that of its predecessor VELO, in terms of both $\sigma_{\text{IP}}$ and track-finding efficiency, despite the increased instantaneous luminosity.

#### 3.2.1 LHC interface

The VELO performance, as described by eq. (3.1), improves by reducing the radius of the first pixel hits, though this must be balanced against the limitations of proximity to the beam line. The minimal VELO aperture allowed by the requirements of the LHC collimation and protection depends on

several factors: the maximum expected separation of the counter-rotating beams; their transverse sizes ($\beta$ functions and transverse emittances); the beam direction relative to the longitudinal axis of the VELO (crossing angle); the mechanical accuracy and stability of the RF boxes. A reduced, but still conservative, radial clearance of 3.5 mm is chosen [35, 36] for the RF boxes, the structures that directly interface with the LHC beam environment. This allows the silicon sensors to be arranged such that the radius of the closest active pixel edge is 5.1 mm from the beam line. The decision takes into account the intended luminosity, beam crossing schemes, luminosity levelling and the requirements of special running scenarios such as van der Meer scans [37].

The electromagnetic fields of the two LHC beams, pulsing at radio frequencies, must also be taken into account. The principle of electrical continuity is maintained for the VELO upgrade with the reimplementation of flexible wakefield suppressors at the entrance and exit of the VELO vacuum vessel. The RF box shape was optimised to reduce the beam impedance while maintaining a good impact parameter resolution. Simulation and measurements with a full-size mock-up of the new RF boxes show the longitudinal and transverse impedance to have good, broadband behaviour when the VELO is closed with 3.5 mm inner radius [38]. In the open position, simulation of the cavity predicts several resonance modes but the total beam power loss due to the impedance presented by the whole VELO is a tolerable 14 W [39, 40]. To prevent possible beam-stimulated electron emission, which can lead to instabilities, the beam-facing surfaces of the RF boxes are coated with a material with a low secondary electron yield (SEY); a low activation temperature non-evaporable getter (NEG) is chosen [41].

### 3.2.2 Mechanics, vacuum and cooling

The concept of separating primary (beam) and secondary (detector) vacua is preserved for the VELO upgrade. The RF boxes must be leak-tight with a tolerance on the pressure difference between vacuum volumes of 10 mbar. The detector components inside the secondary vacuum must be constructed of materials with minimal outgassing and bespoke vacuum-tight solutions are needed to route high-speed data, power and high voltages cables in and out of the vacuum. The vacuum vessel, which is integrated into the LHC beam pipe, remains from the original vertex detector, as well as the large rectangular bellows and detector supports, so that the total allowed detector length of about 1 m, the size and location of access ports and the mechanisms for the horizontal and vertical movement of the detectors are unchanged.

The power dissipation of the FE ASICs operating in the detector vacuum must be removed by a cooling system. Moreover the sensors must be maintained at low temperatures (typically $< -20°C$) for the entire life of the detector, including shutdown periods. Bi-phase $CO_2$ cooling is chosen, following the same principle as in the predecessor VELO. However, the system is entirely redesigned. In the secondary vacuum, the $CO_2$ flows in microchannels within a silicon cooler to which the active components are glued. The risk of a cooling system rupture in the secondary vacuum is mitigated by additions to the mechanical design. A tertiary vacuum volume, the *isolation* volume, is added to house the local distribution of the $CO_2$ supply, including a fast-response bypass valve system. The preserved vacuum vessel is shown in figure 8 (left) contrasted against the parts that have been added or upgraded.

### 3.2.3 Detector geometry and layout

With the LHCb acceptance unchanged, the optimised layout of VELO is similar to its predecessor. Active elements and their services are assembled into a series of identical modules, populated with
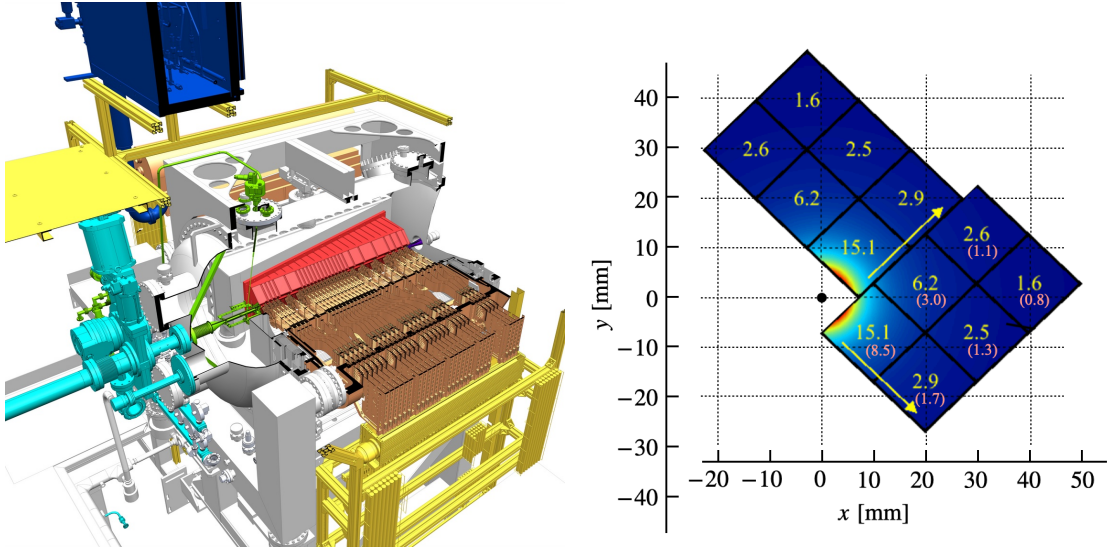
**Figure 8.** Left: a 3D view of the upgraded VELO, with cut-out. Some of the new items are highlighted, such as the Side C pixel modules and readout electronics (brown), the Side A RF box (red), the internal gas target system with a storage cell (green), the upstream beam pipe with a sector valve (cyan). Right: data rate per pixel ASIC in Gbit/s for the most active module. The numbers in parenthesis are the number of traversing tracks per LHC bunch crossing for an average number of interactions per crossing equal to 7.6. Arrows indicate the readout direction.

pixelated ASICs, arranged perpendicular to the beam line. The decision to use identical modules throughout greatly simplifies the production process and quality control. The distribution of the modules must cover the full pseudorapidity acceptance of LHCb ($2 < \eta < 5$) and ensure that most tracks from the interaction region traverse at least four pixel sensors, for all azimuthal directions [42]. With the chosen sensor arrangement shown in figure 9 (left), 52 modules are necessary to satisfy these requirements, including the modules placed upstream of the interaction region whose purpose is to improve the unbiased measurement of primary vertices.

The modules are arranged into two movable halves, the Side C and Side A. Except for a shift, the distribution in $z$ (parallel to the beam line) is identical for the two sides. The minimal, nominal spacing between modules is 25 mm and the Side A modules are displaced in $z$ by +12.5 mm relative to the Side C modules to ensure the two sides overlap when closed to provide a complete azimuthal coverage.

The rectangular pixel detectors are arranged in a rotated 'L' shape, as shown on figure 9 (right). The purpose of the 45° rotation around the $z$ axis is to minimise any risk of the detectors grazing the RF box during installation.

### 3.2.4 Expected particle fluxes and irradiation

The most-occupied 2 cm$^2$ ASIC will experience 8.5 charged particles in every bunch crossing. The LHCb upgrade expects an average bunch-crossing rate of 27 MHz, with a peak rate of 40 MHz. Particles traverse detectors at relatively high angle and on average, given the pixel size of 55 µm × 55 µm, 2.6 pixels will record the passage of an ionising particle. For the busiest ASIC, this implies a peak pixel-hit rate of ~ 900 million/s.

Section 3.3.1 describes the dedicated ASIC developed for the VELO upgrade, which has digital logic that groups hits into super-pixel packets (SPPs) encoded by 30 bits. The busiest ASIC records
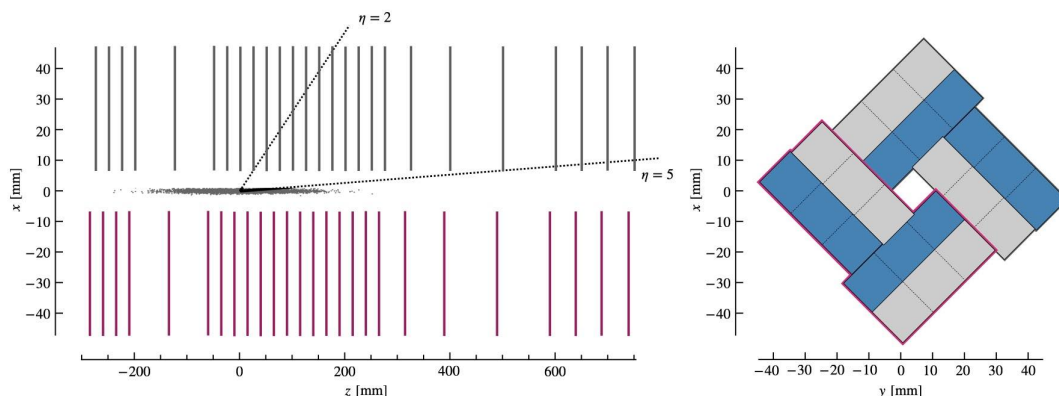
**Figure 9.** Left: schematic top view of the $z - x$ plane at $y = 0$ (left) with an illustration of the $z$-extent of the luminous region and the nominal LHCb pseudorapidity acceptance, $2 < \eta < 5$. Right: sketch showing the nominal layout of the ASICs around the $z$ axis in the closed VELO configuration. Half the ASICs are placed on the upstream module face (grey) and half on the downstream face (blue). The modules on the Side C are highlighted in purple on both sketches.

hits in $\sim 1.5$ SPPs per traversing particle, giving a maximum SPP rate of 520 million/s, or 15.1 Gbit/s from the most central ASIC, see figure 8 (right). The peak total data rate out of the whole VELO may reach 2.85 Tbit/s and the readout scheme is designed accordingly. The power needed for such FE processing is significant and performant on-detector cooling is vital.

The pixel ASIC and silicon sensors are designed to tolerate a high and non uniform fluence, which ranges from $5 \times 10^{12}$ to $1.6 \times 10^{14} n_{eq}/cm^2$ per 1 fb$^{-1}$ of integrated luminosity exposure. With 50 fb$^{-1}$, it is expected that some ASICs accumulate an integrated flux of $8 \times 10^{15} n_{eq}/cm^2$. With this dose, leakage currents of around 200 µA/cm$^2$ ($\sim 7$ nA per pixel) are expected with 1000 V of bias voltage at $-25°$C. In terms of total ionising radiation dose, the ASICs must remain fully operational up to 4 MGy.

## 3.3 The pixel tile

The VELO pixel *tile* is composed of a pixelated, planar silicon sensor and three pixelated ASIC chips. Known as VELO pixel chip (VeloPix) [43], these bespoke ASICs provide analogue signal processing and digitisation. They are bonded to the sensor by an array of solder bumps (SnPb) to form each of the four tiles composing a module.

### 3.3.1 VeloPix

The VeloPix is an ASIC based on the Timepix3 [44] developed by the Medipix/Timepix consortia. It has an active matrix of $256 \times 256$ pixels, each 55 µm $\times$ 55 µm in size, giving a sensitive area of 1.98 cm$^2$. Each ASIC chip is thinned from 700 µm down to 200 µm after fabrication. On three sides, the distance between the edge of the pixel matrix and the physical edge of the device is 30 µm. On the fourth side the ASIC extends by 2.55 mm and contains common digital processing and wire-bond pads.

The ASIC is fabricated in 130 nm complementary metal-oxide semiconductor (CMOS) process,[13] a technology which has proven radiation hardness above 4 MGy. In addition, VeloPix is designed with protection against single event upset (SEU) with the use of dual interlocked storage cells. The

---

[13]By TSMC$^{TM}$ Taiwan Semiconductor Manufacturing Company.

space for this extra logic is obtained by removing some of the functionality present on Timepix3 from the pixel cell, such as the fine-time measurement (640 MHz clock). The main commonalities with Timepix3 are the fast analogue FE with a time walk < 25 ns and zero-suppressed readout using a data-push scheme. Whenever a pixel-hit is recorded, it is time-stamped, labelled with the pixel address and sent from the ASIC immediately, without the need for a trigger signal. One of the main differences between Timepix3 and the VeloPix is the hit rate capability. Timepix3 is limited to a maximum hit rate of 80 million hits/s while the VeloPix can handle 900 million hits/s. Several modifications have been necessary to achieve this capability. The time-of-arrival information is removed so VeloPix records only the occurrence of the hit (*binary* readout). The time-stamp granularity increases from 1.56 to 25 ns, the FE is optimised for a negative input charge and the power budget of the VeloPix is raised to facilitate an increased throughput. An additional data reduction ($\sim 30\%$) comes from grouping $2 \times 4$ neighbouring pixels into a SPP, thereby removing duplication of the time stamp and address fields.

With a 40 MHz acquisition rate, up to 20.48 Gbit/s can flow from one ASIC [45] via four, highly optimised serial links each running at 5.12 Gbit/s. A custom serialiser, gigabit wireline transmitter (GWT), has been designed for this purpose [46]. This bandwidth is significantly greater than the 15.1 Gbit/s anticipated from the busiest ASIC. A notable achievement of the VeloPix development is the power consumption: 1.2 W typical, 1.9 W maximum, which is 65% of the original expectation. Table 2 lists the key VeloPix features.

**Table 2.** Summary of the VeloPix capabilities.

| | |
|---|---|
| Technology | TSMC 130 nm CMOS |
| Radiation hardness | > 4 MGy, SEU tolerant |
| Pixel size (analogue part) | 55 µm × 55 µm (55 µm × 14.5 µm) |
| Peak rate per ASIC (per pixel) | $9 \times 10^8$ hits/s ($5 \times 10^4$ hits/s) |
| Maximum of charge distribution | $16\,000\,e^-$ |
| Minimum threshold | $500\,e^-$ |
| Timing resolution (range) | 25 ns (9 bits) |
| Super-pixel data size | 30 bits |
| Maximum data rate per ASIC | 20.48 Gbit/s |
| Power consumption per ASIC | $\sim 1.2$ W (spec. 3 W) |

### 3.3.2 Sensors

The 208 silicon pixel sensors are each 200 µm thick and 43.470 mm × 14.980 mm large, including 450 µm wide inactive edges in which lie the guard rings. They are manufactured[14] using a float-zone p-bulk with n-type implants insulated between pixels by p-stops. The quoted bulk resistivity is 3–8 kΩ cm. The sensors are designed to provide charge collection efficiency greater than 99% and signals of at least $6000\,e^-$ after 4 MGy and 1000 V applied bias voltage. Key characteristics are listed in table 3.

Each sensor comprises $768 \times 256$ pixel implants, matching the pixel arrays of three ASICs. The sensors are delivered with under-bump metallisation. On three sides the sensor dimensions are larger than the ASICs because of the guard ring. On one side the ASICs extend beyond the sensor, this is

---

[14]Hamamatsu Photonics K.K.™, Hamamatsu, Shizuoka 435-8558, Japan.

**Figure 10.** Top left: microscope image showing the elongated sensor pixels above the inter-ASIC region. Top right: image highlighting the ion-etched round corner of the sensor. Bottom: schematic of the sensor tile, showing the overall dimensions of the sensor and ASIC. The pixel layout is shown only under ASIC 2. There are $256 \times 256$ active bonded pixels (only every fourth pixel is shown in the figure). An additional row of pixels identified in dark red provides a connection between the ASIC ground and the innermost guard ring of the sensor. Three corners, encircled in red, are shown in detail on the left (A, B, C).

**Table 3.** VELO sensor specifications.

| | |
|---|---|
| Bulk material thickness | 200 µm n-on-p silicon |
| Most probable unirradiated signal charge | $16\,000\,e^-$ |
| Minimum end-of-life signal charge | $6\,000\,e^-$ |
| Maximum operational voltage | 1000 V |
| Required charge collection efficiency | $> 99\%$ |

the periphery region with the wire-bonding pads. The metallisation process deposits an additional row of solder bumps on the innermost guard on one side of the sensor and connects the guard ring to the ground row of the ASIC. This serves to tie the guard ring to the ASIC grounds, on the periphery side.

For even distribution of the bias voltage to the backside, a 1 µm aluminium layer is applied on top of the sensor backside. The corners of the sensor are ion-etched to a curved shape to increase the clearance to the RF box and thus minimise risk of contact damage during insertion. To provide sensitivity in the gap between ASICs, sensor pixels that span the distance between two ASICs are elongated to 137.5 µm. These features are highlighted in figure 10 (top).

### 3.3.3 Tile production and quality control

After production by the manufacturer, the VeloPix wafers, each containing 91 ASIC chips, are individually quality-controlled using a semiautomatic probe station[15] with 140 tungsten-rhenium needles. A probe card[16] routes signal from the VeloPix to a SPIDR readout board [47]. Needle probes are placed in contact with the VeloPix ASICs pads to provide power, control and readout during testing. In a first round, the digital functionalities are tested such as power-up, matrix readout behaviour, the calibration of digital to analogue converters and the response to trigger and control commands. In a second round, the quality of the output links and the behaviour of the analogue part of each ASIC are tested. The eye diagram of each output link is verified on an oscilloscope while the pixel noise and equalisation are checked. The sensors and the ASIC wafers are sent to a specialist firm[17] for ASIC wafer thinning to 200 µm, deposition of solder bumps on the ASIC pads, dicing and, finally, tile production by bump-bonding three ASICs and one sensor. Figure 10 (bottom) shows a schematic of the bump-bonded tile and the pixel matrix layout.

The quality control programme continues by checking that the ASICs remain functional and that the sensor can withstand a biasing up to $-1000$ V. The tiles are held on a vacuum jig that allows one to position ten tiles with a 40 µm mechanical precision such that they can be tested with the semiautomatic probe station. While the ASIC pads are connected to the probe card via the tungsten needles, the tile sensor is biased to $-140$ V, the sensor depletion voltage before irradiation. The pixel discriminator response is equalised and the single pixel noise is measured. Analogue test pulses are fired on individual pixels to detect cross-talk, which would indicate a short in the solder bumps. Finally, after equalisation of the discriminators, a strontium source is placed on top of the probe card while a $1000\,e^-$ threshold is set on all pixels. Empty pixels indicate missing bonds. Across the whole tile production, no shorted bumps were found. Tiles with missing bumps were returned to the bump-bonding firm for reworking.

## 3.4 Microchannel cooling

The silicon detectors operate in vacuum and the significant heat generated must be dissipated by conduction. There are 12 ASICs per module, each with a power budget of 3 W. With the power consumption of the FE electronics and the ohmic heating from the irradiated sensors as well, the cooling is designed to extract up to 40 W per module, which is $\sim 2$ kW from the whole VELO. Moreover, to mitigate the effect of the radiation damage, the silicon must be permanently cooled to below $-20\,°C$. The chosen solution is an evaporative cooling system with bi-phase $CO_2$ flowing through microchannels within a silicon substrate approximately 100 cm$^2$ in area, herein referred to as a *cooler* [48].

Carbon dioxide is attractive as a coolant because it is inert, inexpensive and has a large latent heat of evaporation which is exploited in a bi-phase system. By circulating the $CO_2$ in microchannels inside the cooler on which heat-generating components are glued, a thermal performance of 2–4 K cm$^2$ W$^{-1}$ was demonstrated during construction, where the range depends primarily on the 50–100 µm glue thickness between cooler and heat-generating ASICs. The maximum temperature difference between the $CO_2$ exit temperature and the sensor tip, which overhangs the cooler by 5 mm, is 6 °C with nominal ASIC operation (100 µm glue layer).

---

[15]Karl Suss[TM] PA200.

[16]Manufactured by Technoprobe[TM] S.p.A., 23870 Cernusco Lombardone, LC (Italy).

[17]ADVACAM Oy[TM], Tietotie 3, FI-02150 Espoo, Finland.

The $CO_2$ cooling system operates at typical pressures of 14 bar at $-30\,^{\circ}C$ and up to 62 bar at room temperature. For detector safety, the integrity of the entire cooling network, including every cooler, must be verified at three times the maximum operational pressure i.e., 187 bar.

### 3.4.1 Microchannel cooler design and fabrication

The 20 microchannels within each 500 µm-thick cooler range in length from 271 mm to 332 mm. The channels stretch from the inlet manifold, pass under the detector elements before returning to the outlet manifold located adjacent to the inlet. Each channel is 200 µm wide and 120 µm deep. For the first 40 mm of every channel, the width and depth is reduced to 60 µm × 60 µm to provide a uniform flow restriction. This ensures an even distribution of coolant despite the varying heat load; the $CO_2$ boils in a uniform manner as the cross section increases by a factor $\sim 7$ at the end of the restrictions. An illustration of the composite device is shown in figure 11 beside an X-ray image that reveals the microchannels inside the cooler. The $CO_2$ passes into and out of the silicon through slits machined in fluidic connector that match the position and size of the microchannel manifolds. The fluidic connector is made of Invar36 whose thermal expansion closely matches that of silicon. Two vacuum-brazed, $^1/_{16}$ inch pipes with inner diameter of 0.57 (0.87) mm service the inlet (outlet) of the connector, which is attached to the silicon by a fluxless soldering assembly.



**Figure 11.** Left: illustration of the silicon microchannel coolers and fluidic connector. Right: the parallel lines represent the etched microchannels, which can be seen in the X-ray image.

The microfabrication[18] of the silicon was performed using deep-reactive ion etching and direct bonding techniques. The process starts with double-side polished silicon wafers in which the microchannel patterns are etched. The microchannels are closed by applying a second silicon wafer and using hydrophilic bonding. Afterwards, the second wafer is thinned to 240 µm and the first wafer to 260 µm achieving a final thickness of 500 µm. A soldering pad comprising three layers of metallisation, Ti (200 nm), Ni (350 nm) and Au (500 nm) is deposited around the microchannel manifold. Alignment marks are also deposited during this step. Last, the inlet and outlet slits are opened by ion-etching before plasma-dicing cuts the silicon into the characteristic shape.

---

[18]CEA-Leti, 38054 Grenoble, France.

**Figure 12.** Left: image of a fluidic connector aligned in the horizontal plane with the corresponding solder layer on the silicon cooler, indicated with the red shadow. Right: X-ray of the solder joint attaching the fluidic connector to the microchannel cooler. No solder has entered the two inlet regions, and no large voids are seen in the solder layer.

### 3.4.2 Cooler assembly and quality control

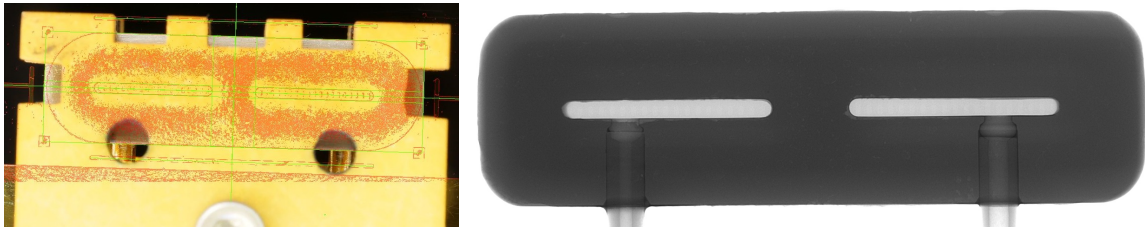In order to attach the fluidic connector to the silicon in a reliable and reproducible manner, a novel, fluxless soldering technique was developed. The technique involves many steps, including polishing, cleaning (acetone, ethanol, ultrasound bath, plasma cleaning), outgassing, and pretinning solder on each piece. Pretinning is performed in a reducing atmosphere of nitrogen with 3% concentration of formic acid vapour. Subsequently, the opposing silicon and Invar surfaces are aligned and the solder is reflowed in vacuum. During this procedure, nitrogen at atmospheric pressure is reintroduced to minimise the size of voids forming in the molten solder. After soldering, a rectangular piece of silicon is glued behind the manifold, where it serves to reinforce the area where the blow-out force is largest. The full procedure is described in ref. [48].

Due to the risks associated with this new technology, an exhaustive quality control programme was developed. An X-ray of the solder joint, e.g. figure 12 (right), checks for voids in the solder layer and verifies that no solder has entered the fluidic pathway. Each cooler is checked for leak tightness by helium detection to below $10^{-9}$ mbar l/ s. Tests are done in two ways: in ambient air with 60 bar helium inside the cooler and with the cooler placed in a 1 bar helium atmosphere whilst pumping on the microchannels. The silicon surface is then checked for flatness with a tolerance on planarity of $100\,\mu$m. To check for robustness, the cooler is pressurised with nitrogen to 130 bar for 45 minutes (the burst-disk release pressure of the final system) as well as a 15 minute, 187 bar stress test. Finally, a visual inspection requires the silicon surface to be clean of any residue and the shape of the pipe loops to be within tolerance for the module construction. Around 30 full-size prototypes were trialled before 81 installation-quality coolers were produced with an acceptance yield of 87%.

### 3.5 The module

The VELO module brings together the silicon detectors, their cooling, powering, readout and mechanical support into a single, repeating unit. This section describes the design, construction, and quality control of these objects.

### 3.5.1 Anatomy

A VELO pixel module is a double-sided detector structure with a microchannel cooler at its core. The cooler is glued at the Invar cooling connector to the *mid-plate*, a carbon fibre plate which sits atop two carbon fibre legs anchored on an aluminium foot. On each cooler face, a pair of tiles are glued at right
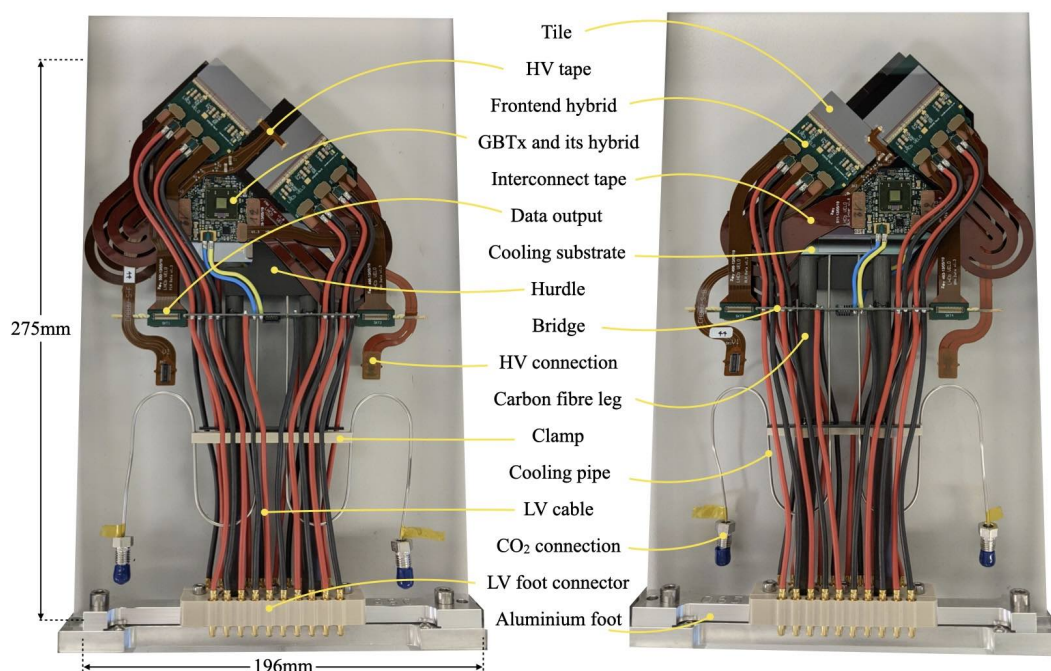
**Figure 13.** Photo of the (left) upstream and (right) downstream faces of a fully-assembled VELO module.

angles such that they approach the LHCb beam line along their inner edge. Along their periphery, each triplet of VeloPix ASICs are wire-bonded to the FE electronics *hybrid* that transmits control signals and routes their data out to a PCB flex cable. The 2 FE hybrids (per module side) are connected to a third hybrid housing a gigabit transceiver (GBTx) control ASIC, see section 3.6.2. Power is delivered through an assembly of 20 silicone-coated copper cables and a PCB transition bridge. The bridge is used to change to thinner cables near the FE in order to reduce the material budget close to the beam line. It also provides mechanical support for the various cables. The high voltage (HV) for the silicon sensors is delivered by PCB flex cables which are bonded to the top surface of the sensors. Figure 13 shows these components on both sides of an assembled module.

### 3.5.2 Assembly and quality control

The module production was performed in clean rooms at two production sites. All steps followed preagreed instructions to ensure consistency across the two sites. Every process and qualification was recorded in an online database, which automatically aggregated the results to provide immediate feedback. Following several years developing the assembly process, 53 identical, installation-quality modules were built over the course of 16 months.

Different stages during the assembly of a module are shown in figure 14. First, two carbon fibre legs are glued[19] to the aluminium foot, and at the other end to the carbon fibre mid-plate. This rigid structure is glued to the cooler at the Invar fluidic connector, using the same glue. The cooling pipes are clipped into a clamp attached to the legs, which minimises the risk of transmitting stress to the cooler from manual handling of the pipes. The four tiles are glued[20] to the microchannel cooler

---

[19]Using Araldite[TM] 2011.
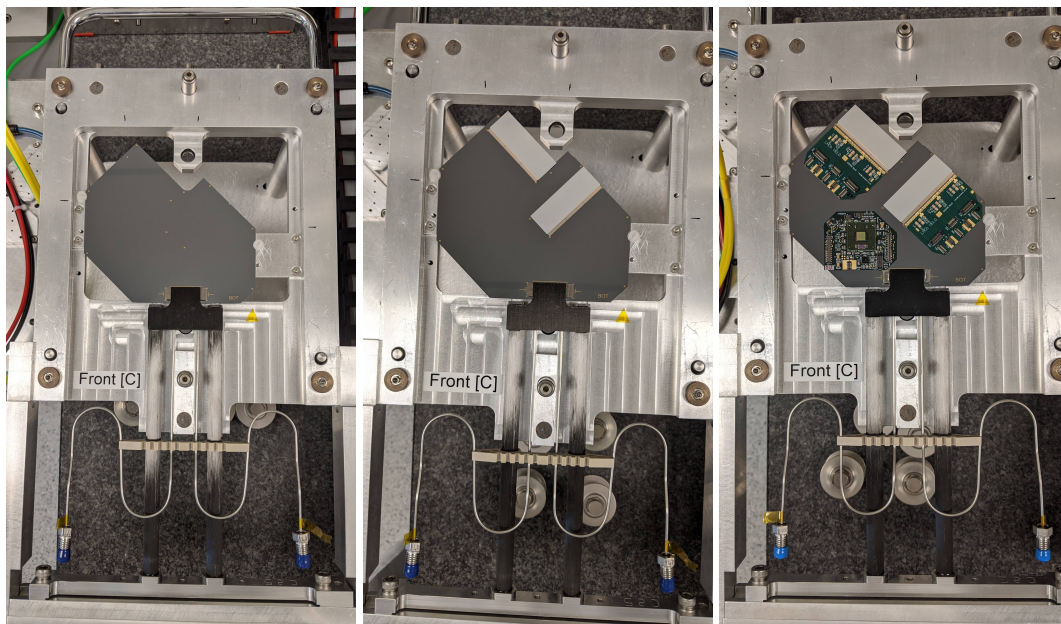[20]Using Stycast[TM] FT2850 and 23LV catalyst.

**Figure 14.** From left to right: bare module; module with tiles; then with hybrids too.

with precision-made jigs. The placement of the tiles is performed with an in-plane precision relative to the foot of better than 30 μm. In addition, the thickness of the glue layer is a critical parameter of the module quality, as it affects not only the mechanical properties of the attachment, but also the cooling performance of the module. In order to meet all the requirements on the glue interface, a uniform layer with thickness in the range of 50–100 μm is required. The placement of the FE and GBTx hybrids is performed in a similar manner, although the placement tolerance is relaxed to 100 μm. A silicone adhesive[21] is used for this process as its flexibility ensures minimal stress on the cooler given the different coefficients of thermal expansion on either side of the glue layer. The HV cables are connected to the cooler with a fast-curing glue[22] and manually bonded to the sensor surface. Finally the module is placed in a wire-bonding machine, where an automated programme carries out this process one side at a time, for a total of 1680 bonds. The last step consists in attaching all interconnecting and power cables to the module.

Once built, the functionality of each module is verified in vacuum with cooling close to $-30\,°C$, both before and after a thermal cycle. The two-way communication with the GBTx chips and the 12 VeloPix ASICs is verified as well as the response of each VeloPix to fast trigger signals and the equalisation of the pixel matrix. A bit error rate test of the 20 output links is performed and the sensor leakage current is checked with a HV scan.

### 3.5.3 Metrology

Since the nominal distance of the module tips to the RF box is as small as 0.8 mm, direct knowledge of the position of the innermost tiles is vital to confirm suitability for installation. The relative position of the sensors on the module is also an important input to the track-based alignment algorithm. The metrology is divided in two parts: orientation and position in the plane of the sensors ($xy$ alignment);

---

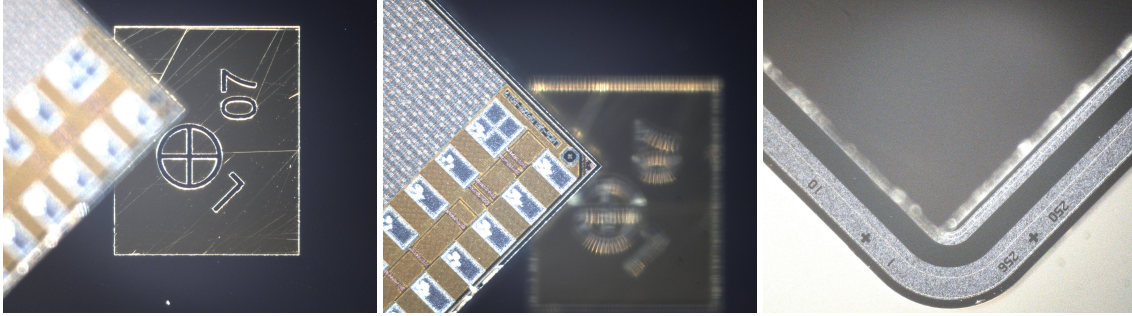[21]Loctite™ SI 5145.

[22]Araldite™ 2012.

**Figure 15.** Cross-shaped markers used for module metrology. Left: on the microchannel cooler; centre: on the VeloPix ASIC; right: on the ASIC-side of the sensor.



**Figure 16.** Tile metrology results showing $x$ and $y$ positions and angles for each module.

and position perpendicular to that plane, parallel the beam axis ($z$ alignment). The $xy$ alignment is based on measurements of markers (see figure 15). The absolute positions of these markers are measured relative to the dowel pin on the module foot. These measurements extract deviations of the $x$ and $y$ positions and rotation of a sensor around the $z$ axis in the nominal module frame. In the $z$ direction, the relative position and orientation of the sensors is affected by any curvature of the cooler and variations in the thickness of the glue layer.

Each production site developed its own method for metrology of the modules, however both achieved a resolution of a few microns in all coordinates. The results in $z$ reveal thicknesses of the glue between tile and cooler within the range 30–110 µm, with a mean value of 80 µm. The results in $xy$, shown in figure 16 show excellent tile positioning in the $y$ direction, which is most critical in terms of clearance to the RF foil.

### 3.6 Electronics and readout chain

The main role of the VELO electronics system is to transport data from the VeloPix ASICs to the off-detector processing units. The system also delivers clock and control signals to the modules, as well as low voltage to power the electronics and high voltage to bias the sensors.

### 3.6.1 System architecture

The electronic components of the system are located in three places: on the detector module, immediately outside the VELO vacuum vessel and off detector, with dedicated cabling running between them. An overview of the system is shown in figure 17.

**Figure 17.** Block diagram showing the main parts of the VELO electronics system.

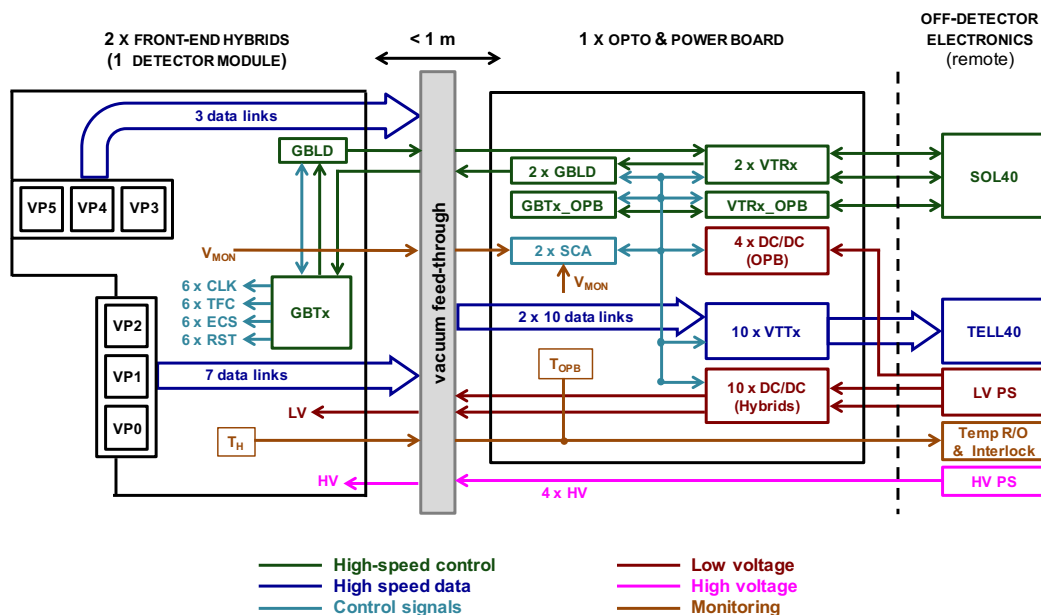The VeloPix ASIC [45] reads the analogue signals from the sensor and sends binary hit data in SPP (see section 3.3.1) over serial links at rates up to 5.12 Gbit/s per link. Serial data routing as well as distribution of clock, control and power is managed on the FE by hybrid circuits (section 3.6.2). The serial data from the hybrids is transmitted out of the secondary vacuum on high-speed serial links, through a vacuum feedthrough board to the opto- and power board (OPB), see section 3.6.5, mounted on the exterior of the vacuum vessel. The control signals to and from the FE are transmitted on identical serial links. The low and high voltage is supplied through the same vacuum feedthrough board to separate cables. The temperature monitoring (section 3.6.6) is routed through the data flex cables and the OPB.

Through optical links, the OPB transmits the data to the TELL40 data acquisition cards (section 3.8) whilst receiving control signals from the SOL40 readout supervisor. The TELL40 and SOL40 boards are located in the DAQ server rooms in the data centre, which requires over 300 m of optical fibre. The high and low voltage supplies are located in the electronics barracks in the LHCb cavern, requiring about 60 m of cable.

### 3.6.2 Front-end circuits (hybrids)

On the module, the distribution of ASIC power, clock and control signals and the routing of outward-bound data is provided by hybrids. These are four-layered, flexible printed circuits interconnected with two-layered polyimide cables. The hybrids have a total thickness of 390 µm and come in two types, as shown in figure 18. The first type provides the FE electronic interface to each VeloPix where wire bonds are used to connect the hybrid to the ASIC periphery. The second type houses the GBTx chip and distributes timing signals and fast control instructions to the VeloPix ASICs via the FE hybrids. Slow controls are routed through the GBTx hybrid as well as monitoring of the bias voltage.

On each module face there are two FE hybrids and one GBTx hybrid. The FE electronics are packaged into these three pieces rather than one larger hybrid circuit in an effort to minimise stress
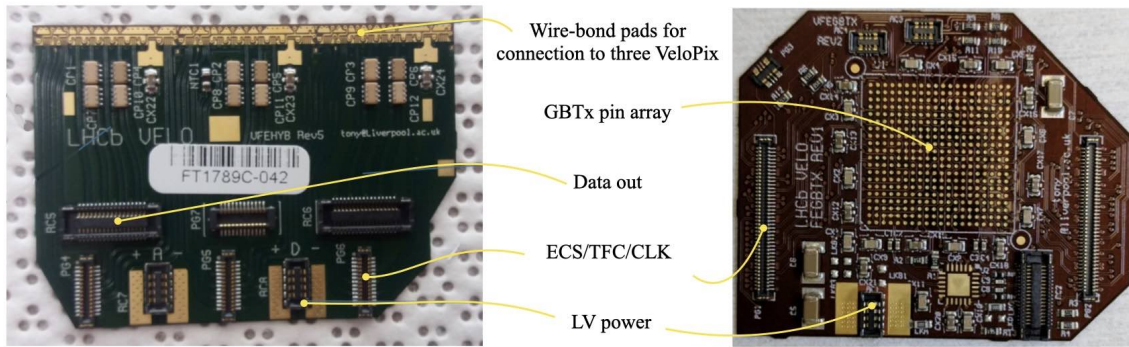
**Figure 18.** The FE (left) and GBTx (right) hybrids.

on the module. The necessary thickness of these hybrids is driven by the cross section of copper necessary to transfer a suitable current to the ASICs, and leaves them relatively rigid. As described in section 3.5.2, a sufficiently flexible glue was chosen to attach the hybrids to the cooler and absorb contraction differences when cooling to $-30\,°C$.

### 3.6.3 High speed serial cables

All module readout and control signals are routed inside the vacuum using low-mass, high speed PCB flex cables [49]. Polyimide microstrip technology is used on the module where low-mass materials are of critical importance. From the module to the vacuum wall, four flex cables each carry up to seven 5.12 Gbit/s serial links. The dielectric used for these flex cables has to be radiation tolerant and provide low dissipation loss for high-frequency signals. In addition, high reliability and yield are required for the impedance control. The chosen material for this purpose is an all-polyimide thick copper-clad laminate[23] offering high signal integrity, reliability and has wide use in medical and aerospace applications.

### 3.6.4 Vacuum feedthroughs

The routing of such a dense number of high speed signals through a vacuum barrier represented a challenge which could not be solved by any commercially available solution. A bespoke solution was developed using a PCB with edge metallisation glued into a vacuum flange with epoxy. The board is a 12-layer printed circuit (2 mm thick) through which the high speed data signals are transmitted. The low voltage power supply to the detector electronics ($\sim 2.5\,A$), the bias voltage ($< 1000\,V$) and the temperature data also pass through this board. Tests show that sealing with epoxy[24] provides better leak tightness than standard O-ring vacuum sealing. The feedthrough boards are grouped in sets of six or four, depending on position in the main vacuum flanges, which are sealed by O-rings to the VELO vacuum vessel. An image of the vacuum feedthrough board is shown in figure 19 with an image of a flange populated with six boards.

---

[23]DuPont Pyralux™ AP-*PLUS*.

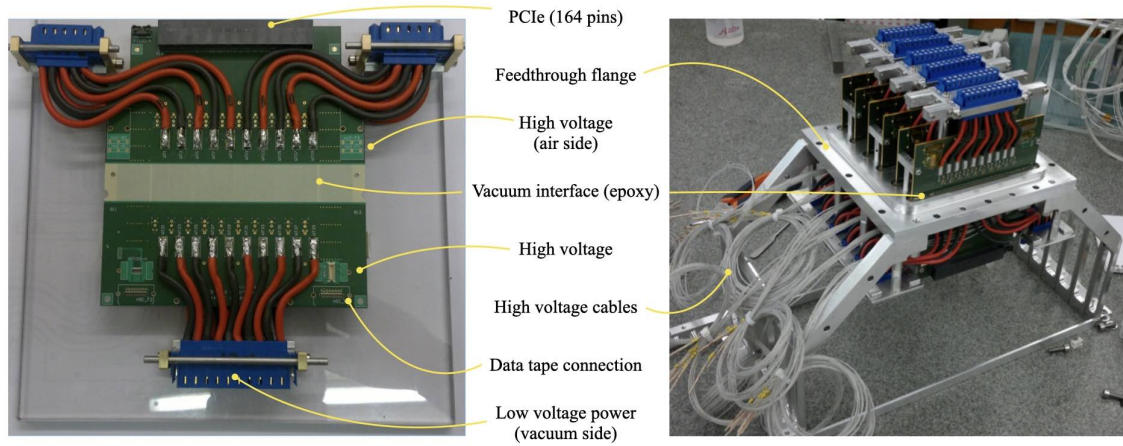[24]Araldite™ 2011 (2020) on the air (vacuum) side.

**Figure 19.** Left: custom-made vacuum feedthrough board. Right: a populated vacuum flange.

### 3.6.5 Opto- and power board

The opto- and power board (OPB) [50] are the interface between the detector modules and the off-detector electronics. The OPB has 14 DC/DC converters [17, 51] that transform the five input voltages to the ten supplies needed by the detector module and the four voltages required by the OPB itself. It performs the electrical-to-optical conversion of the 20 high-speed data links from each module and the bidirectional electrical-to-optical conversion for the three control links, one for each side of the detector module and one for the OPB itself. The OPB hosts a GBTx ASIC [52, 53] that decodes the high-speed control links and interfaces to two gigabit transceiver slow control adapter (GBT-SCA) ASICs [16]. These latter two devices provide ADC channels that monitor the supplied and regulated voltages as well as the received optical power. They interface with inter-integrated cicuit (I2C) buses that control the GBLD laser drivers [54, 55] on the optical transceivers, and general-purpose I/O signals used for local control. The connections to the temperature sensors (NTC thermistors) on the detector module and OPB are available on the front panel of the OPB.

The OPB is manufactured with eight metal layers and has an overall size of $40 \times 15 \, \text{cm}^2$, plus a 10 mm protrusion which connects with the peripheral component interconnect express (PCIe) connector on the vacuum feedthrough board. The layers that surround the high-speed electrical signals are made of low-loss dielectric[25] and the remaining dielectric layers are made of standard FR4 glass fibre laminate. The OPBs are vertically mounted on the exterior wall of the VELO vacuum vessel in a dedicated mechanical frame integrated with the vacuum feedthrough flange. The frame acts as a partial heat sink for the OPBs. Forced vertical air flow through the crates provides additional cooling.

### 3.6.6 Temperature and voltage monitoring

Across each VELO half there are 730 temperature sensors split between monitoring and safety readout systems. Within each cooling loop, which contains two modules and one shared isolation valve, 50 temperatures are monitored (54 for the seventh module pair), which are broken down into the following: 10 Pt100 probes attached to the cooling pipes; 24 NTC sensors from the hybrids; 8 band gap measurements implemented in the VeloPix, one per tile; and 8 more measurements, one per

---

[25]Isola<sup>TM</sup> I-Tera MT40.

**Figure 20.** The Side C half, with 26 modules, ready for the installation into the vacuum vessel.

GBTx and four on the OPB boards. There are a further 22 Pt100 sensors distributed over the RF box, module support base, and in the isolation vacuum. In addition, for each of the 13 module pairs there are 10 low voltage readings and 4 independent readings of the high voltage, which gives 182 voltages to be monitored.

### 3.7 Mechanical design

The VELO mechanical design retains the concept from the original vertex detector of two movable halves, retracted from the beam line at all times other than when stable beams are circulating. Each VELO half moves independently in the horizontal direction from a −29 mm retraction from the beam line to +4 mm overclosure. The halves have common vertical motion and may be moved ±4.7 mm above or below the beam line. The VELO detector halves, their support structures and the RF boxes are replaced to accommodate the requirements of the new pixel detector.

The central structure, onto which 26 modules are mounted, is the aluminium module support base, visible in figure 20. Whereas the modules are designed to have minimal radiation length, the bases are built for precision and rigidity. Any distortion of the bases moves the module tip towards the RF box with a lever-arm equal to the module height. To avoid thermal distortions, they are maintained to 20 °C (the manufacturing temperature) by several adhesive heating pads. Once installed, the bases are bolted to the detector support which is, in turn, fixed to large, rectangular bellows that provide a flexible barrier between the primary and secondary vacua. All electronic and cooling services run

**Figure 21.** Illustration of the VELO halves showing modules on the module support bases and the LHCb acceptance as a transparent pyramid. On the left, the flexible electronic cables are shown leading to the vacuum feedthrough boar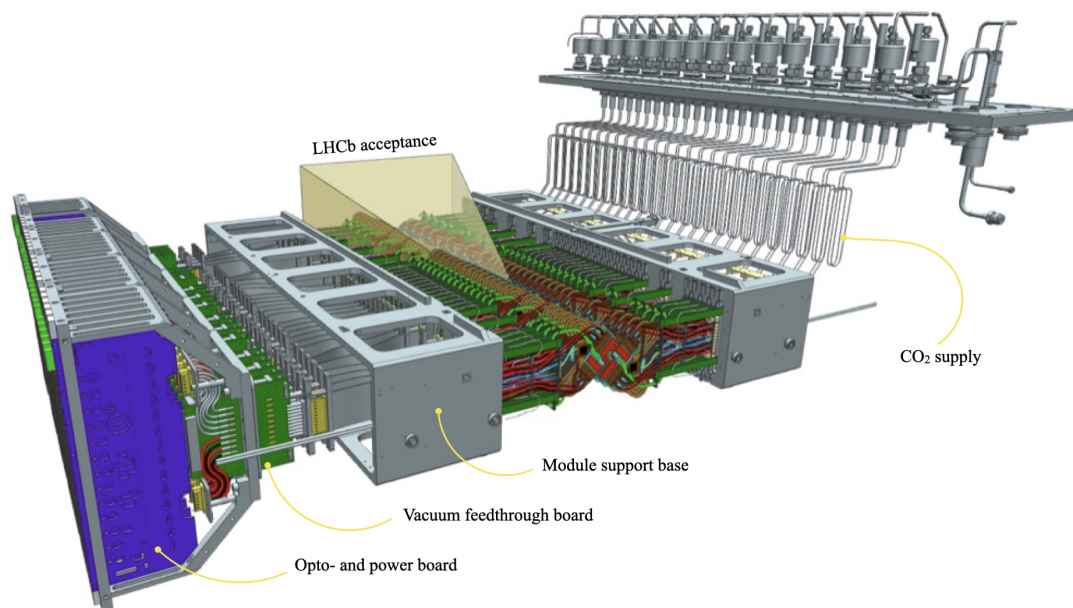ds and OPB boards in their custom frame. On the right, the flexible construction of long cooling loops is shown as well as the interface between the secondary and isolation vacua, in which sits an array of valves.

from the movable bases to the fixed *detector hood*, which is the large flange that seals the detector volume on the external wall of the vacuum vessel. The ∼ 3 cm travel of the halves is absorbed by flexible power and data cables running between the module foot and the vacuum feedthrough. These details are shown in figure 21. For the $CO_2$ supply, an elongated cooling loop, incorporated into every pipe running to/from each module, absorbs the movement. The cooling lines are connected to a series of valves located in the tertiary vacuum, the isolation volume.

### 3.7.1 RF boxes

The RF boxes are the thin-walled corrugated enclosures that provide the barrier between the primary (beam) vacuum and the secondary (detector) vacuum and interface the VELO detector halves to the LHC beams. They are made from aluminium, a light and electrically conductive material. Their complex shape accommodates overlaps between sensors of opposing halves, while maintaining electromagnetic effects to an acceptable level. In order to compensate for the reduced spatial resolution of the pixel detector, compared to that of the innermost microstrips of the previous VELO sensors, the distance of approach to the beams was reduced from 8.2 to 5.1 mm. The beam aperture, as defined by the inner surface of the RF boxes, reduces from 5.5 to 3.5 mm. Special blocks of AlMg4.5Mn0.7 alloy were forged to obtain a homogeneous material with small grain size and without cavities. The initial blocks had dimensions $1200 \times 300 \times 300 \, \text{mm}^3$ and were milled to the desired shape with a 5-axis milling machine.

The RF box fabrication procedure included several steps, such as verification of the block quality, rough milling of the outside and inside shape, stress-relieving annealing, final milling to
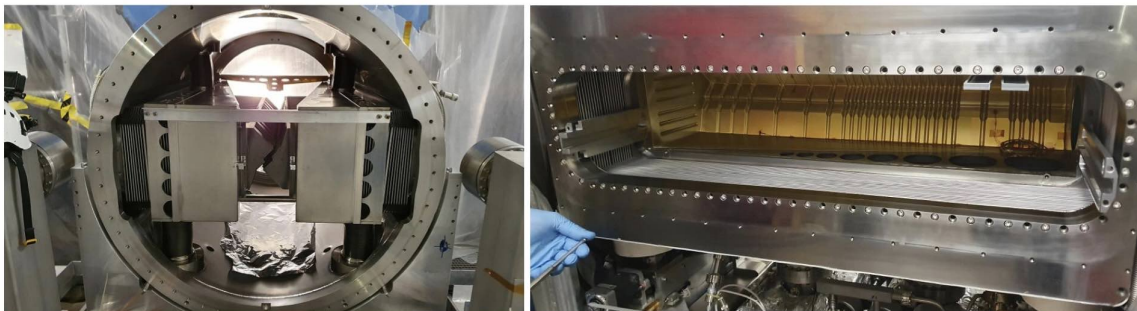
**Figure 22.** Left: the open VELO vessel (seen from upstream) during the installation of the RF boxes. Right: view inside the Side A RF box showing the module slot structure.

the nominal 0.25 mm thickness with use of special moulds, supported by wax-filling techniques, and interspersed thickness measurements to achieve the desired thickness and geometry with the required precision in a reproducible manner.

A vacuum test of each RF box was performed by closing the volume with a flat flange and filling with helium at 1–5 mbar pressure inside a large vacuum vessel. The leak rate was measured with a mass spectrometer for different pressures to know the background level; no leaks were detected. The metrology was done with the RF box mounted on its side. RF box deformation studies at ±10 mbar under- and over-pressure showed a maximal wall displacement of 0.4 mm in the central region (slot 22).

In total, two pairs of RF boxes (one for installation and one spare) were fabricated. The RF boxes to be installed were further thinned down by chemical etching along the central spine within ±25 mm of the beam line. The etching proceeded in 20 minute steps of 0.02 mm using a painted mask on the aluminium to stop etching before the wall could become locally too thin. The final result is a typical thickness reduction from $0.30 \pm 0.07$ mm to $0.25 \pm 0.10$ mm, with a minimum intended thickness of 150 μm.

To protect against electrical breakdown between the RF box and VELO sensors, the inside of the RF boxes is sprayed with polyamide[26] prepared in an aqueous solution, to about 10 μm thickness. This procedure included the attachment of Pt100 temperature probes on the RF boxes. Finally, the RF boxes were again cleaned before applying a NEG coating to the beam-facing surfaces. Figure 22 (left) shows the two RF boxes installed in the VELO vacuum vessel and the module side of one RF box with its characteristic corrugated shape to accommodate modules.

### 3.7.2 Wakefield suppressors

Wakefield suppressors at both ends of the VELO vacuum vessel are used to smooth the transition from the wide (50–56 mm) cylindrical beam pipe to the narrower apertures of the storage cell and RF boxes. The wakefield suppressors must be flexible enough to accommodate the motion of the VELO halves, and provide a continuous electrical path for the mirror currents of the LHC beams. They are made from a 50 μm thick, corrugated copper-beryllium sheet with a deposited gold finish to reduce the secondary electron yield. The upstream connection attaches to the storage cell and is further described in section 4. On the downstream end of the vacuum vessel, the suppressor is clamped to the exit window on a 0.2 mm thick gold-plated stainless-steel tension ring, as done previously.

---

[26]Torlon™ from Solvay Specialty Polymers.

The attachment to the RF boxes is made by four keyholes in the wakefield suppressor and matching mushroom features formed on the ends of the RF boxes during their manufacturing.

### 3.7.3 Vacuum safety system

The risk of a microchannel rupture causing a sudden rise of pressure in the secondary vacuum and damage to the RF boxes is mitigated by continuous supervision of the pressure inside that volume. Measurements of pressure over the range of 1 bar to $10^{-9}$ mbar are made using three different types of pressure sensors. If a sudden pressure rise is detected, the pneumatic shut-off valve system is activated and coolant flow is diverted from all modules to bypasses. This action immediately reduces the amount of the coolant injected in each cooling loop and in the case of a microchannel rupture, would minimise the amount of $CO_2$ released. Each cooling loop serves two microchannel coolers and contains two shut off valves (inlet and outlet) and an additional, small, safety volume. It is thermally connected with the detector structure to ensure a temperature higher than the coolant. Consequently, in case of activation of the safety system, coolant is trapped in between two shut off valves, will expand and reduce in temperature following the typical $CO_2$ saturation curve. The 26 pneumatic bypass valves, two for each two-module cooling loop, are installed in the isolation volume, an entirely new feature of the VELO mechanics.

## 3.8 ECS and DAQ

The final component in the VELO read-out chain are the PCIe40 cards, hosted in PCs in the data centre. These cards are hardware-wise identical for all subdetectors, including VELO, while the firmware is partially subsystem-specific. Detector data, control and monitoring is performed within the LHCb common and centrally managed framework, based on several platforms, see sections 10, 11 and 12. Within these, a number of VELO specific tools are developed and used.

### 3.8.1 Firmware

Cards with the SOL40 firmware are used to control and monitor thirteen VELO detector modules, over a total of 39 bidirectional optical fibre links. This firmware, which is largely shared with the other LHCb subdetectors, includes the trigger and fast control commands and the protocols for communication with and via the GBTx [52, 53] and GBT-SCA [16] chips. For the complete VELO detector, a total of four SOL40-flavoured PCIe40 cards and 52 TELL40-flavoured cards are used.

Cards with the TELL40 firmware are used for processing the twenty high-speed GWT data serial links of a single module. This firmware is unique to the VELO detector because of the different clock frequency with which the VeloPix serialiser protocol (GWT) works compared to the GBT protocol used by all other subdetectors for their data taking. The firmware, described in detail in ref. [56], processes the data in the following way. First it performs the descrambling of the GWT serialiser output frame, splitting it back into the SPP format that was encoded on the VeloPix. Time ordering of all hits within the last 512 clock cycles is done. This corresponds to the 9 bit resolution of the timestamp inside the SPP. Any SPP arriving outside the 512 clock cycle time window is dropped. The overall LHC clock timing information is added to each hit. The firmware also performs pattern recognition of cluster of pixels to recognise particle hits and evaluate their coordinates and the cluster topology, dropping the raw pixel data [57, 58]. This accelerates the track reconstruction process in both high level trigger first stage (HLT) and high level trigger second stage (HLT).

### 3.8.2 Detector control

Communication with the hardware components is done via the LHCb ECS software framework using both fast and slow control protocols. The project manages the configurable elements and monitors low-level parameters of the hardware, including the DAQ system, the OPBs, the GBTx chips and the VeloPix ASICs. It also collects, logs and displays information on the VELO environment (temperatures, voltages, pressures, etc.). The calibration is handled using fast readout, using the external DIM writer library to catch the data points from the detector. The control application sends the calibration data, in a raw format, to a dedicated database. The raw data are analysed in C++ calibration software (Vetra). Since the calibration analysis is CPU-consuming, it is moved to a separate machine. When the control application receives back the calibration result, it configures the detector accordingly.

### 3.8.3 Calibration data processing

The Vetra software package, written in Gaudi framework, is dedicated to handling standard LHCb data, and special VELO-only data collected for specific purposes (e.g. IV scans). The core part of Vetra is written in the C++ and it can be used as a quasi-online monitoring application to rapidly analyse raw data produced by the VELO. The overall processing pipeline comprises data decoders, processing algorithms and monitoring algorithms. The VELO is monitored with this software tool and the result is fed back to the detector control system where necessary. For example, occupancies, dead-pixel maps, charge-collection efficiency studies, alignment and other calibrations. This software is also used for track and vertex reconstruction, providing beam position information to the control algorithm that ensures safe closure of the VELO halves. A significant part of the data handling, processing and visualisation will be done outside Vetra. For this purpose, a database system, called Storck, and visualisation framework, Titania, have been created. Storck is optimised to manage the calibration data files and their child objects (e.g., calibration parameters, monitoring histograms), whilst Titania provides a means for rapid data exploration and trending. These tools are vital to the operation of the subdetector and quality control of the VELO data.

## 4 Internal gas target

Fixed-target physics with LHC beams was pioneered in the LHCb experiment during Run 2 thanks to the availability of an internal gas target. A gas injection system, called system for measuring the overlap with gas (SMOG) [59], originally conceived and implemented for precise colliding-beams luminosity calibration, was used to inject light noble gas into the VELO vacuum vessel. This produced a temporary local pressure bump peaking at around $10^{-7}$ mbar over the length of the vessel (about 1 m) and decaying down to the LHC background level ($\sim 10^{-9}$ mbar) over the 20 m LHCb beam pipe sections on each side of the interaction point. The resulting beam-gas interactions were used for precise imaging of the beam profiles [60]. SMOG also gave the unique opportunity to operate the LHCb experiment in fixed target mode. Gaseous targets of different nuclear size (He, Ne and Ar) were used in combination with proton and lead beams at (nucleon-nucleon equivalent) centre-of-mass energies of up to 115 GeV, with negligible effect on LHC operation. Encouraged by first results and future prospects, an upgrade of SMOG (also called SMOG2) was proposed and implemented [61].

The core idea of the SMOG upgrade is to inject the gas directly into a so-called storage cell and benefit from the increased areal density at an identical injected flux, as has been done in the past at

other accelerators [62]. The principle is sketched in figure 23. The open-ended cylindrical tube has an inner diameter $D$ and a length $L$. Gas is injected via a capillary at the storage cell centre at a flow rate $\Phi$ from a gas feed system (GFS), resulting in an approximately triangular density distribution $\rho(z)$ with maximum $\rho_0 = \Phi/C_{\text{tot}}$ at the centre ($z = 0$). Here, $C_{\text{tot}}$ is the total flow conductance of the tube from the centre outwards and is given by the conductance of two parallel tubular conductances of length $L/2$ in the molecular flow regime [62], and thus amounts to (in $1\,\mathrm{s}^{-1}$) $C_{\text{tot}} \approx 7.62\sqrt{T/M}\,D^3\,(0.5\,L + 1.33\,D)^{-1}$, where $L$ and $D$ are in cm, the storage cell temperature $T$ in K, and $M$ is the molecular mass number of the injected gas. The areal density seen by the beam is $\theta = \rho_0\,L/2$.
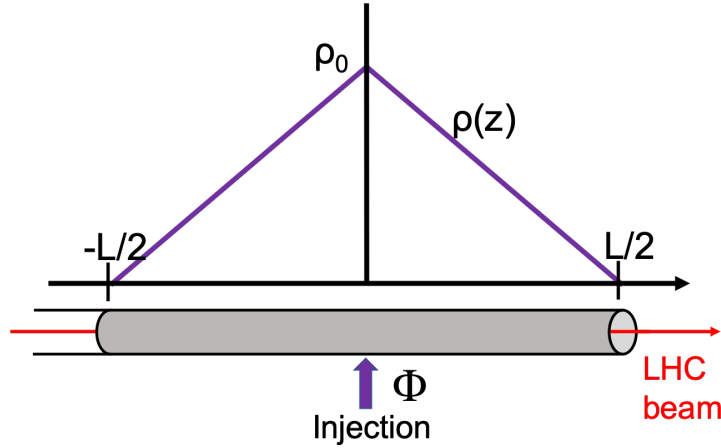


**Figure 23.** Sketch of a tubular storage cell of length $L$ and inner diameter $D$. Gas is injected at the centre with flow rate $\Phi$, giving a triangular density distribution $\rho(z)$ with maximum $\rho_0$ at the centre. Reproduced with permission from [61].

The VELO upstream connection in the LHC vacuum between the VELO detector boxes and the beam pipe has been modified to accommodate the integration of a storage cell composed of two cylindrical halves, each attached to the upstream end of one of the VELO detector halves and moving together with them. Thus, when the VELO is brought into a closed position, the two halves form an open-ended tube coaxial with the LHC beam axis. It is expected that the SMOG upgrade will facilitate fixed-target runs with an effective gas areal density higher by a factor of about 10 for He at the same flow rate, and even higher gain factors for heavier gases.

The SMOG upgrade introduces other important improvements. First, the determination of the target density (and beam-gas luminosity) is significantly more precise because the target is confined to the storage cell, whose conductance is well known and can be combined with an accurate measurement of the injected gas flow rate from the GFS. Second, it will be possible to select among several gas species without intervention (including non-noble gases such as $H_2$, $D_2$, $O_2$, etc.). Finally, the beam-gas interaction region is much better defined and well separated from the beam-beam collision region, which also opens the possibility to have concurrent beam-gas and beam-beam collisions.

The SMOG upgrade is composed principally of two systems: the storage cell assembly, mounted inside the beam vacuum, and the GFS, located on the "balcony", a platform near the detector inside the experimental cavern. Given its vicinity to the LHC beams, the design of the storage cell assembly must fulfil several requirements derived from aperture considerations, RF or impedance related aspects, and dynamic vacuum phenomena, as already discussed for the VELO RF boxes (see section 3.7).
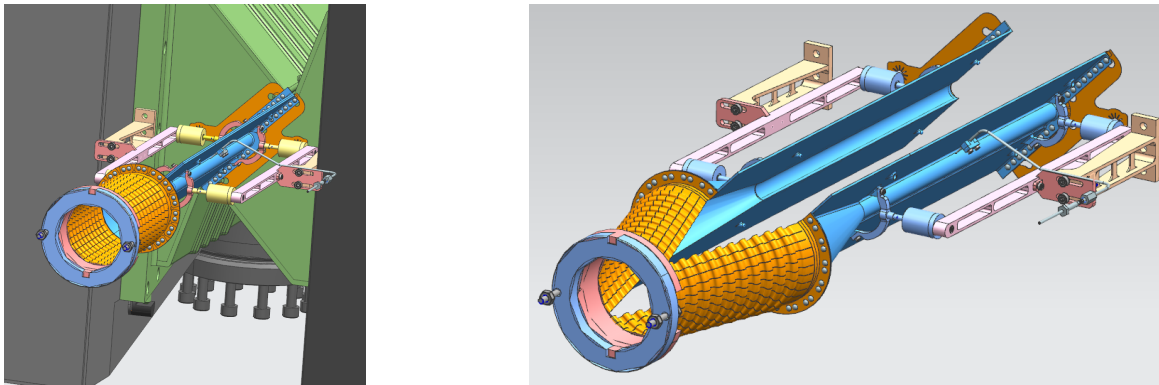
**Figure 24.** Left: view of the storage cell (blue) supported from the VELO RF box flanges (in green) in the closed VELO position. Two flexible wakefield suppressors (orange) provide the electrical continuity. Right: storage cell in the open position (without showing VELO elements). Reproduced with permission from [61].

## 4.1 The storage cell system

The storage cell and its arrangement inside the VELO vessel is visible in figure 24. The assembly fits into the limited space available inside the existing VELO vessel, upstream of the VELO detector. In order to leave sufficient beam aperture for beam operations (injection, energy ramp, squeeze, etc) the assembly is split in two opposing halves each attached to its respective VELO RF box. The assembly is composed of a flexible wakefield suppressor split in two halves, two opposing storage cell shapes containing a half cone, a half tube and side wings, a short wakefield suppressor which connects to the VELO detector box, and two arms to support the storage cell halves from the VELO RF box flanges. The conical shape allows for a smooth transition from the 56 mm diameter of the upstream beam pipe to the 10 mm diameter of the storage cell tube. The tubular part is 20 cm long. Gas is injected from the GFS into the tube centre via a flexible line ended with a 0.8 mm inner diameter stainless steel capillary pressed into a hole in the Side C half of the storage cell. All parts are sufficiently light to keep the beam-induced background due to the material in the proximity of the beams at a negligible level.

The cell opens and closes together with the VELO detector boxes to which it is mounted by two cantilevers rigidly attached to the flange of the VELO RF boxes. Because the VELO design foresees the possibility to operate the detector at a slightly retracted position (by ∼ 0.1 mm) relative to the nominal closed position, the Side C half of the storage cell is rigidly fixed to its cantilever, while the Side A half is coupled via a spring system that allows to always reach the final closed position, guided by the rigid half. This spring system allows to reach the storage cell nominal closure (thus, sufficient gas tightness along the longitudinal slit) even if the VELO halves are not completely closed, within a range of up to 1 mm. The minimum allowed aperture over the length of the storage cell is imposed by the van der Meer scan configuration and amounts to 3 mm [36], well below the chosen radius of the storage cell.

The surfaces surrounding the beam are made of electrically conductive materials, in order to shield the chamber from the beam RF fields, prevent excitation of RF modes and provide electrical continuity for any position of the VELO halves. The cell structure is made from a 99.5% pure Al block, milled with an accuracy of ∼ 20 μm. The cone, tube and wings are milled to a final thickness of 1.2, 0.2, 1.2 mm respectively. Before completion, the cell was heat-treated to 290°C for one hour to allow for stress release. All screws are silver-coated and perforated in compliance with standards for high vacuum systems. Particular attention was given to the transitions to the RF

boxes (downstream side) and to the beam pipe (upstream side). They are made from 0.075 mm thick Cu-Be foil. The upstream wakefield suppressor contains slits and corrugated strips which ensures adequate flexibility to the wakefield suppressor. It is attached to the upstream support by screws and to the storage cell by Al rivets. The downstream wakefield suppressor is made of short curved fingers. It is connected on one end to the storage cell, in the same manner as the upstream wakefield suppressor. The fingers press on the RF box, while the wakefield suppressor is locked to the RF box on the mushroom shapes mentioned in section 3.7.2.

The storage cell is coated with amorphous carbon to present to the beams a surface with a SEY around ~ 1.0 [63]. This precaution avoids the formation of a beam-induced electron cloud and the possible onset of beam instabilities. The coating is applied by sputtering. First, a 50 nm thick Ti adhesion layer is applied, then a $1 - 10$ nm thick layer of amorphous carbon. From simulation studies it has been concluded that such SEY is largely sufficient to prevent electron cloud build-up, even when taking into account that a higher residual gas pressure can favour such phenomenon.

The storage cell is equipped with five 0.34 mm outer diameter K-type thermocouples (with a precision of about 0.1 K) insulated with a nickel-based super alloy[27] and terminated with a ceramic connector for use in ultrahigh vacuum. The temperature measurements are needed both for determining the areal density $\theta$ ($\sqrt{T}$ dependence) and for monitoring a possible temperature increase by beam-induced effects.

The electromagnetic compatibility of the SMOG-VELO assembly was validated using frequency and time domain electromagnetic field simulations that were benchmarked with RF measurements on a 1:1 scale mockup with the wire method [38, 39]. The wakefield suppressor robustness has been demonstrated on a prototype with a fatigue test (15 000 open/close cycles), equivalent to about 15 years of nominal operation in the experiment. No sign of fatigue has been observed. The installation of the storage cell into the VELO vessel has been successfully completed in the summer of 2020, see figure 25. A detailed alignment survey of the storage cell found no misalignment in excess of 0.25 mm.

## 4.2 Gas feed system

The GFS allows one to choose the gas type to be injected among those available in the four installed reservoirs. The amount of injected gas can be accurately set and measured in order to precisely compute the target densities from the storage cell geometry and temperature. The GFS consists of four assembly groups, as shown in figure 26, and is based on precise absolute thermo-stabilised gauges that cover four decades of pressure reading: absolute gauge 1 (AG1), which covers a pressure range from 1100 to 0.1 mbar, and absolute gauge 2 (AG2), which covers the range from 1.1 to $10^{-4}$ mbar. The absolute gauges measure and monitor the pressure in the main volume, and are therefore used in determining the stability of the injected flow. The flow is obtained by setting a nominal pressure and keeping it constant by means of a thermo-regulated valve (DVS). After stabilising the injection pressure, another thermo-regulated valve (DVC) is set to the appropriate value depending on the gas type and the gas flow rate chosen (typical values will be in the range $0.5$–$8 \times 10^{-5}$ mbar l/ s). The connection of the GFS to the VELO vessel is performed by a 10 mm inner diameter bakeable stainless steel pipe with a length of 15 m terminated by two feed lines (with valves CV and VV), one feeding directly into the VELO vacuum vessel and one into the storage cell centre. A full range gauge (FRG)

---

[27]Special Metals Corporation Inconel$^{\text{TM}}$.

**Figure 25.** Picture of the SMOG storage cell system installed into the VELO vessel.



**Figure 26.** The four assembly groups of the SMOG GFS are the gas supply (with 4 reservoirs), the main table, the pumping station (PS) and the feed lines to the VELO vacuum vessel. The various components are described in the text. Reproduced with permission from [61].

monitors the pressure just upstream of these two valves. A rest gas analyser (RGA) is employed to analyse the composition of the injected gas in the main volume.

Preliminary studies using the MOLFLOW simulation code [64] show that a total systematic uncertainty of around 2–3% in the determination of the target gas areal density is achievable (slightly worse for light gases compared to the heavier ones).

# 5    Upstream tracker

## 5.1    Overview

The UT is located between the RICH1 detector and the dipole magnet. It is used for charged-particle tracking and is an integral component of the first processing algorithm in the software trigger [9]: the UT hits are combined with the VELO tracks and, exploiting the magnetic field between the interaction region and the UT, a first determination of the track momentum $p$ with moderate precision ($\sim 15\%$) is obtained. A momentum and charge estimate is performed only for tracks with $p_T > 0.2\,\text{GeV}/c$ and this information is used to speed up the matching with the SciFi Tracker hits. These two features provide significant improvement of the speed of this matching algorithm. Moreover, UT hit information reduces the rate of fake tracks created by mismatched VELO and SciFi Tracker segments. Lastly, it provides measurements for particles decaying after the VELO, e.g. long lived $K_S^0$ and $\Lambda$ particles.

### 5.1.1    Detector requirements

The physics goals and environmental conditions dictate the following requirements:

- Acceptance: in order to fulfil its function in the trigger algorithm and to be effective in suppressing fake tracks, the coverage of the UT detector in the nominal LHCb acceptance should not have any gaps.

- Single-hit efficiency: the detector should have a high enough hit efficiency to ensure that more than 99% of the charged particles traversing the detector within the acceptance leave hits in at least three planes.

- Hit purity: spurious hits due to noise or signal shape must be minimised. This is further specified in the FE ASIC section.

- Radiation damage: the UT detector needs to maintain its performance up to an integrated luminosity of at least 50 fb$^{-1}$, taking into account that the radiation dose has a strong radial dependence. The sensitive elements near the beam pipe need to withstand fluences up to $4 \times 10^{14} n_{\text{eq}}/\text{cm}^2$, while the maximum fluence is less than $2 \times 10^{13} n_{\text{eq}}/\text{cm}^2$ for most of the detector area, with the outermost sensors receiving less than $10^{12} n_{\text{eq}}/\text{cm}^2$. In addition, the near detector electronics need to withstand a radiation level of the order of 1 kGy.

- Occupancy: the charged particle density follows a similar radial trend as the radiation fluence. The detector segmentation must be finer near the beam pipe in order to keep the occupancy below a few percent.

- Material budget: the detector and its enclosure must be designed with the goal of obtaining a significant overall reduction of material in the forward region of the acceptance when compared to the LHCb detector of Run 1 and Run 2.

A silicon microstrip detector technology was chosen to fulfil these requirements.

### 5.1.2 Geometry overview

The UT detector comprises four planes of silicon detectors organised in two stations, as shown in figure 27. The circular hole in the middle provides clearance for the beam pipe. The silicon strip pitches and lengths are matched to the expected occupancy. The silicon sensors, shown as coloured boxes in the image, are described in section 5.4. They are arranged in vertical units, called *staves*, described in section 5.3. The first station (labeled 'a') is composed of an *x*-measuring layer (UTaX) with vertical strips and a stereo layer (UTaU) with strips inclined by 5°. Both layers are made of 16 staves each. The second station ('b') is similar, with first a stereo layer (UTbV) with opposite inclination, and a layer with vertical strips (UTbX). Both layers contain 18 staves each. The two pairs of stations are symmetrically arranged around $z = 2485$ mm. There is a gap of 205 mm between the nominal $z$ positions of UTaU and UTbV, and of 55 mm between the UTaX and UTaU or UTbV and UTbX.

To ensure full coverage in the vertical direction, the sensors are arranged on both sides of the staves such as to obtain a vertical overlap. Similarly, a *z*-staggered arrangement of the staves with horizontal overlaps facilitates a full horizontal coverage. Moreover, special sensors are utilised in the innermost area to maximise the active area near the beam pipe. Compared to its predecessor, the Tracker Turicensis (TT) detector [21] of LHCb during Run 1 and Run 2, these improvements considerably reduce the gaps in the acceptance.

In order to minimise the amount of material seen by particles, all the components in the active volume of the detector have been thinned as much as practical.



**Figure 27.** Drawing of the four UT silicon planes with indicative dimensions. Different colours designate different types of sensors: Type-A (green), Type-B (yellow), Type-C and Type-D (pink), as described in the text. Reproduced from [79]. CC BY 3.0.

## 5.2 Mechanics and near detector infrastructure

An overview of the detector and its associated electronics and mechanical services is shown in figure 28. The staves are enclosed in a thermally insulating, light and air-tight box that fits around the beam pipe. The UT beam pipe section is wrapped in a lightweight thermally insulating blanket.

A cooling manifold distributes the $CO_2$ coolant to the staves. Dry gas is flushed through the box to prevent condensation on the stave components. These items are further described in section 5.8. An off-detector electronics system, described in section 5.6, is mounted close to the staves, just outside the detector box. These electronics convert the signals from the FE chips into optical signals and drive them along optical cables. The box is split into two halves and mounted on rails, such that the detector can be opened for maintenance or during beam pipe bake-out. Four cable chains allow horizontal movement of the two box halves without putting strain on electronic cables, optical fibres and cooling fluid pipes. Four service bays, on the fixed racks, located on the Side A and Side C of the UT detector, host the low voltage boards that power the near detector electronics and the hybrids, and patch panels that distribute the high voltage to the silicon detectors.



**Figure 28.** A 3D view of the UT system.

### 5.3 Staves and modules

A UT stave provides the mechanical support for the sensors and FE electronics, as well as active cooling. An instrumented stave is shown in figure 29 (left) and an exploded view shows the individual components (right). The bare staves are designed to minimise the radiation length in the acceptance region of the detector. Each stave is approximately 10 cm wide, 1.4 m long and is composed of a carbon fibre sandwich, about 3.9 mm thick, that contains a 2.275 mm diameter titanium tube as part of the evaporative $CO_2$ cooling. The tube is embedded in a low density high thermal conductivity carbon foam[28] for good heat conduction. Polymethacrylimide[29] structural foam fills the remaining voids. Four PCB flex cables, called *dataflex*, are glued on the two faces of the bare stave in a process involving stencils to ensure the use of a controlled and optimal quantity of glue. The dataflex provides the

---

[28]Allcomp[TM] K9.

[29]Rohacell[TM].

electrical connectivity from the stave outer edges to the FE electronics. Three types of dataflex cables are needed, Short, Medium and Long. The FE readout electronics and silicon sensors are assembled in *modules*, themselves glued and bonded to the dataflex cable. A UT module is composed of a silicon sensor (see section 5.4), supported on a boron nitride ceramic stiffener, itself glued to a *hybrid* flex circuit (see section 5.6.1) which hosts the readout ASICs (see section 5.5). The modules are glued to both sides of a stave in a staggered manner such that sensors overlap in the vertical direction.



**Figure 29.** Left: a completed stave. At the bottom, the end of the cooling tube is visible with the high voltage and signal connections. In orange are the dataflex cables and in brown the hybrids. The reflective areas are the silicon detectors. Reproduced with permission from [65]. Right: an exploded view of an instrumented stave showing the individual components described in the text.

In order to fulfil the occupancy requirements in different regions of the detector, four sensor types (named A, B, C and D) are utilised with different strip pitch and different strip lengths. Hybrids can host either 4 or 8 ASICs, the former version being used for sensors of Type-A and the latter version for all other sensor types. The modules are also of four types, A, B, C and D, named after the sensor type that it hosts. This, in turn, leads to three different stave types (see figure 29 and 36):

- staves of variant A are used to cover most of the detector acceptance and host 14 sensors of Type-A; they have one Short and one Medium dataflex cable on each face;

- staves of variant B include 10 Type-A and 4 Type-B sensors served by one Short and one Medium dataflex cable on each face;

- staves of variant C are used for the region adjacent to the beam pipe and include 10 Type-A, 2 Type-B, 2 Type-C and 2 Type-D sensors; here, one Medium and one Long dataflex cable are needed on each face.

Table 4 summarises the most salient features and numbers of the UT staves and modules.

**Table 4.** Summary of UT detector components.

| Staves | | | |
|---|---|---|---|
| Variant | Quantity | dataflex cables per stave | Module types per stave |
| A | 52 | 2×Short, 2×Medium | 14×A |
| B | 8 | 2×Short, 2×Medium | 10×A, 4×B |
| C | 8 | 2×Medium, 2×Long | 10×A, 2×B, 2×C, 2×D |
| Modules | | | |
| Type | Quantity | Sensor type | Hybrid type |
| A | 888 | A | 4-chip |
| B | 48 | B | 8-chip |
| C | 16 | C | 8-chip |
| D | 16 | D | 8-chip |

## 5.4 Silicon sensors

The silicon microstrip sensors of the UT need to cope with occupancies and radiation fluences spanning different orders of magnitudes. For this reason, four different sensor designs are used, with n-in-p for the central region and p-in-n in the outer region. Strip isolation in the p-substrate sensors is achieved through p-stop technology. All sensors were produced by Hamamatsu.[30] Type-A sensors constitute the majority of the detector and cover the outermost parts of the instrumented area. These n-substrate sensors have 512 p-type strips with a pitch of 187.5 µm. All other sensor designs use p-substrates and have 1024 n-type strips with a twice smaller pitch of 93.5 µm. Type-B sensors are located at approximately 10 cm from the beam line. The innermost area is covered by Type-C and Type-D, which have half the strip length. Type-D is characterised by even shorter strips on part of the sensor as its shape includes a circular cut-out that matches the beam pipe. In this way, a minimum distance to the beam line of 34 mm is achieved, as seen in figure 30 (left).

The sensor design and arrangement in the UT detector keeps the maximum occupancy below 1%, being highest in the Type-D sensors. The signal-to-noise performance necessary for the efficiency requirements has been demonstrated in the beam tests discussed in section 5.9. For all sensor types, high voltage is brought to the sensor backplane via a silicon implant embedded along the top-side edges of the sensor. This simplifies the stave construction by allowing high voltage contact to be made via wire bonds on the same side as the connections between sensor strips and readout electronics. The pitch of the readout ASIC channels matches the strip pitch on sensors of Type-B, -C, and -D. The Type-A sensors require a pitch adapter (PA), which is achieved using additional metal layers on the sensor surface, separated from the AC coupled metal strips and the guard rings by a thin layer of silicon dioxide. This embedded PA design is shown in figure 30 (right). While prototype sensors showed signs of signal coupling to these metal layers resulting in efficiency loss and spurious hits (see section 5.9), this effect has been minimised in the final design by locating the bonding pads above the guard ring of the sensor.

The sensor thickness and resistivity have been chosen to ensure full depletion of the sensors over the life of the detector. Before irradiation, full depletion is achieved applying between 200 and 300 V

---

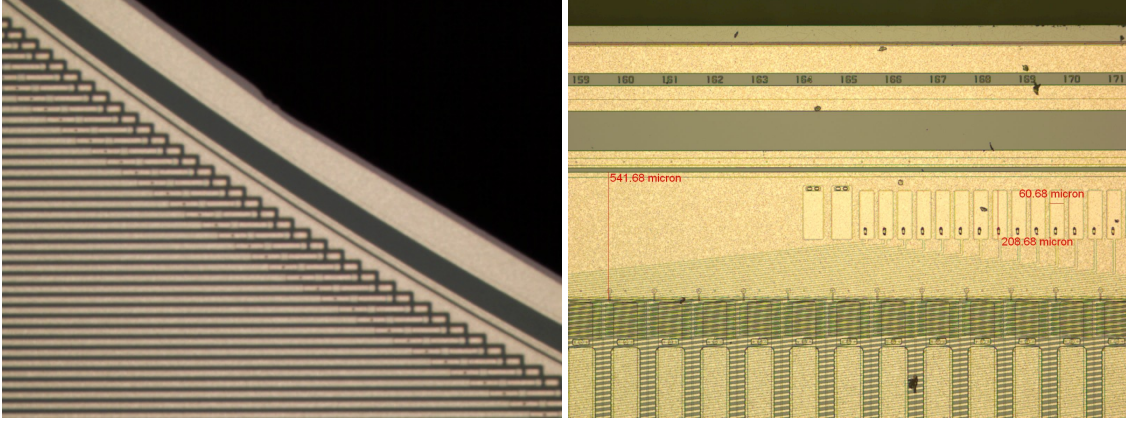[30]Hamamatsu Photonics K.K., Hamamatsu, Shizuoka 435-8558, Japan.

**Figure 30.** Design features of the UT silicon sensors (see text). Left: Type-D sensor cut-out region. Right: embedded pitch adapter. Reproduced with permission from [66].

**Table 5.** Main design parameters for the UT silicon sensors. For Type-D the given length is for outside the cut-out region.

| Design type | Implant/ bulk | Length mm | Width mm | Thickness μm | Pitch μm | Strips | Note |
|---|---|---|---|---|---|---|---|
| A | p/n | 99.50 | 97.50 | 320 | 187.5 | 512 | embedded PA |
| B | n/p | 99.50 | 97.35 | 320 | 93.5 | 1024 | |
| C | n/p | 51.45 | 97.35 | 250 | 93.5 | 1024 | |
| D | n/p | 51.45 | 97.35 | 250 | 93.5 | 1024 | cut-out |

bias voltage to the sensor. It is expected that the most irradiated p-substrate sensors (Type-D) will be fully depleted with less than 500 V at the end of the lifetime of the detector.

The main features of the four sensor designs are summarised in table 5.

### 5.5 SALT readout chip

A 128-channel ASIC, called Silicon ASIC for LHCb Tracking (SALT), was developed in a 130 nm CMOS process[31] to read out silicon strip detectors of the UT. A block diagram is shown in figure 31. The SALT features an analog processor and a low-power (less than 1 mW/channel), fast (40 MSps) 6-bit ADC per channel, followed by a digital signal processor (DSP) block, a data formatting block and a serialiser block.

The analogue FE comprises a charge preamplifier with pole-zero cancellation and a fast 3-stage shaper (with peaking time $T_{peak}$ = 25 ns and a fast recovery) required to distinguish between the LHC bunch crossings at 40 MHz. The FE is designed to work with load capacitances in the range 1.6–12 pF. Each channel contains an 8-bit trimming DAC for baseline equalisation. In the last stage of the analogue FE a single-to-differential block converts a single-ended signal to a differential one.

A fully differential 6-bit successive-approximation-register (SAR) ADC running at 40 MHz converts the analogue signal to the digital domain [67]. In order to achieve highest speed and lowest power consumption (significantly below 1 mW) the SAR logic is asynchronous and dynamic circuitry

---

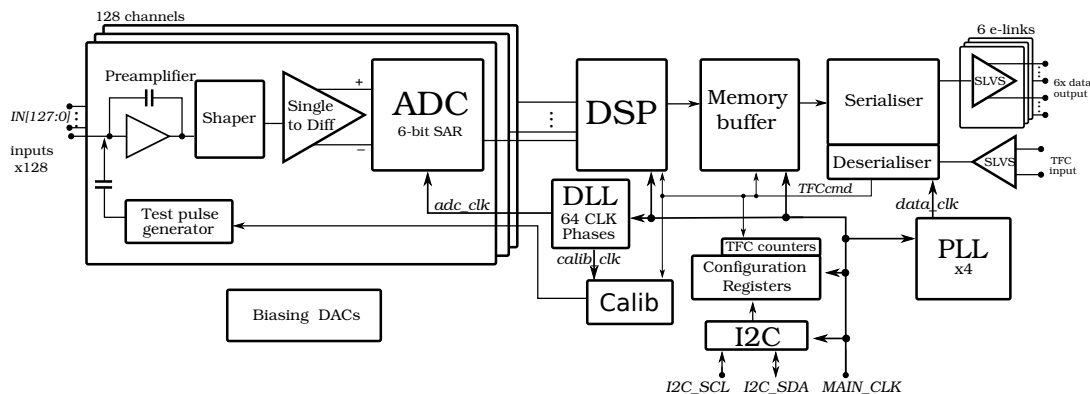[31]By TSMC™ Taiwan Semiconductor Manufacturing Company.

**Figure 31.** Block diagram of the 128-channel SALT ASIC. Reproduced from [70]. CC BY 4.0.

is used in the ADC logic and comparator. To synchronise ADC sampling instances with beam collisions a dedicated ultra-low power (< 1 mW) DLL is used to shift and align an external clock. The ADC samples are signed 6-bit numbers coded as two's complements.

The digital ADC output is processed by a DSP block, which performs a pedestal subtraction, a mean common mode (MCM) subtraction and a zero suppression. For better testability the DSP can also transmit raw ADC data or various combinations of partially processed data. In the DSP calculation, in each place where subtraction is performed, a saturation arithmetic is used (giving numbers inside the range from −32 to 31). In the next step the data packets are created and recorded in a local memory.

After the DSP the data are serialised with 320 Mbit/s rate, obtained by increasing the system clock frequency by a factor four and double data rate (DDR) transmission. An ultra-low power (< 1 mW) phase-locked loop (PLL) is used to generate the 160 MHz clock from the 40 MHz system clock. The data are sent out by an SLVS interface. The ASIC is controlled via the LHCb common protocol consisting of two interfaces: the TFC and the ECS [68, 69]. The TFC interface delivers the 40 MHz clock and other crucial information and commands, synchronised with the experiment clock, while the ECS serves to configure and monitor the ASIC and is realised through an I2C interface. The main specifications of the SALT ASIC are shown in table 6.

The 130 nm CMOS process is generally considered as intrinsically radiation resistant against total ionising dose. For this reason, and given the main clock frequency and the expected charged particle fluence, SEE protection in the SALT is generally limited to SEU protection and monitoring. In a few specific blocks working with very small currents (< 1 µA), like 7-bit baseline DACs, protection with NMOS enclosed layout transistors was applied. To improve the ADC robustness against SEE, the ADC is reset by the sampling signal if the previous conversion is not yet completed. The triple-voting technique is applied to almost every flip-flop in the whole digital part. The combination logic is not triplicated, but clock and reset signals are. As a result, to triplicate the complete clock and reset trees, there are three PLLs and three input reset synchronisers. In addition, the configuration registers have built-in self-correcting circuits that can correct the radiation-induced bit flips (only one of the three copies) in each clock cycle.

Several ASIC design iterations were required to meet the UT requirements. The chip versions used in the UT are the SALTv3.5 for most of the 4-chip hybrids and the SALTv3.9 for the 8-chip hybrids and some of the most exposed 4-chip hybrids. More detailed information on the SALT ASIC can be found in ref. [70].

**Table 6.** Summary of the specifications of the SALT ASIC.

| Variable | Specification |
|---|---|
| Technology | TSMC CMOS 130 nm |
| Channels per ASIC | 128 |
| Input / Output pitch | 80 µm/140 µm |
| Total power dissipation | < 768 mW |
| Radiation hardness | 0.3 MGy |
| Sensor input capacitance | 1.6–12 pF |
| Noise | $\sim 1000\,e^-$ @10 pF $+50\,e^-$/pF |
| Maximum cross-talk | Less than 5% between channels |
| Signal polarity | Both electron and hole collection |
| Dynamic range | Input charge up to $\sim 30\,000\,e^-$ |
| Linearity | Within 5% over dynamic range |
| Pulse shape and tail | $T_{\text{peak}} \sim 25$ ns, amplitude after $2 \times T_{\text{peak}} < 5\%$ of peak |
| Gain uniformity | Uniformity across channels within $\sim 5\%$ |
| ADC bits | 6 bits (5 bits for each polarity) |
| ADC sampling rate | 40 MHz |
| DSP functions | Pedestal and MCM subtraction, zero suppression |
| Output formats | Non-zero suppressed, zero suppressed |
| Calibration modes | Analogue test pulses, digital data loading |
| Output serialiser | Three to five serial e-links, at 320 Mbit/s |
| Slow controls interface | I2C |
| Fast digital signals interface | Differential, SLVS |

## 5.6 Readout chain

### 5.6.1 Hybrids

The UT hybrids are the FE boards for the microstrip silicon sensors. Being mounted in the detector acceptance, they were designed with mass minimisation in mind. They are low-mass flex PCBs composed of a stack of conductive and dielectric layers[32] with a total thickness of about 72 µm Cu, 176 µm polyimide and 63 µm adhesive. Their function is to distribute TFC and I2C signals, and route e-links from the ASICs to the the dataflex cables, while maintaining high signal integrity. Hybrids are designed to distribute low and high voltage (LV and HV) coming from the dataflex cable to the ASICs and sensors with high filtering performance. The UT uses two types of hybrids: VERA hybrids accommodate 4 SALT chips and they are used in most of the detector, while SUSI hybrids host 8 chips each and instrument the centre of the detector, where the strip density is double. In total, 888 VERA and 80 SUSI are needed to populate the full UT. Hybrids have been produced in panels, as visible in figure 32 (left), with VERA (SUSI) panels containing 8 (6) circuits each, fully instrumented with connectors and alignment holes. The connectors allowed a full electronic test to be performed without the need of cutting the hybrid away from the panel. Precut slits were added to ease the mounting on modules. In figure 32 (right) a picture of a SUSI hybrid in a final UT module is shown.

---

[32]DuPont$^{\text{TM}}$ Pyralux: two FR 7013 and two AP 8535R sandwiching an adhesive LF 1500.
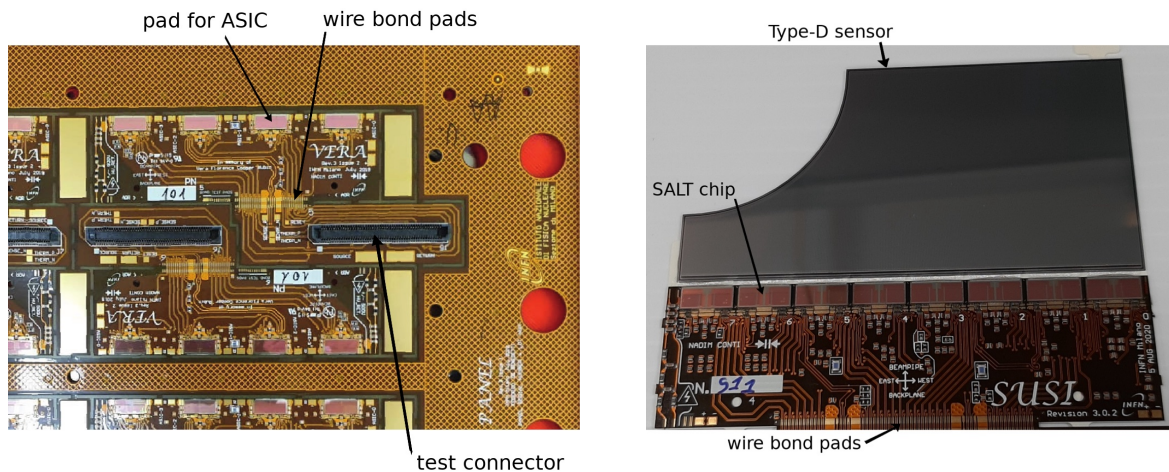
**Figure 32.** Left: two fully visible VERA hybrids, not yet equipped with ASICs, embedded in a carrier 8-hybrid panel before cutting. Right: SUSI hybrid mounted in a final UT Type-D module.

### 5.6.2 Passive PCB flex cables

Two families of passive flex cables are used to transfer signals and power from the peripheral electronics, located outside the acceptance, to the hybrids: pigtails provide the connection between outside and inside the detector box, dataflex cables provide the connection between pigtails and hybrids on the staves.

The dataflex cables are flexible PCB passive boards glued directly on the staves. They are designed to provide high fidelity signal integrity, be lightweight and maximise heat transfer from modules to the underlying support, while minimising the voltage drops across the board. The sensor bias voltage lines are isolated from the low voltage and data lines. Each dataflex cable connects to one pigtail via a high-density 400-contact connector[33] and to the modules via wire bonds. HV connectors on the outer end connect directly to HV cables that are fed through the detector box wall. The dataflex cables were designed and produced in three different sises, Short, Medium and Long, to match the various locations and types of modules along the stave. Their respective lengths are about 604, 700 and 748 mm. In figure 33 a picture of two sample dataflex cables is displayed. The stack[34] includes about 225 µm of polyimide, 70 µm of Cu and 100 µm of adhesive. A pictorial representation of the dataflex cable usage in the UT is found in figure 36 (left).

Pigtails are flexible PCB boards that provide connection between dataflex cables (inside the detector box) and the peripheral electronics (outside the box) discussed in section 5.6.3. The latter connection is made via a high-speed high-density connector.[35] A complex design was developed to cope with all the mechanical constraints and satisfy the electrical requirements related to the low voltage power and the SALT differential pair signals. The 272 installed pigtails are identical from the electrical point of view, despite the fact that dataflex cables can host varying numbers of chips. To accommodate the required lengths and angles due to the positions of the staves and the peripheral electronics, 11 different pigtail shape variants have been designed. Figure 34 shows three pictures of the pigtails. In the leftmost picture two different variants are shown for illustration. The two 1 mm thick

---

[33]MEG-Array™ 400 BGA connector from Amphenol™.
[34]DuPont™ Pyralux: LF 0110, AP 8535, LF 0100, AP 9121, LF 0100, AP 8525 and LF 0110.
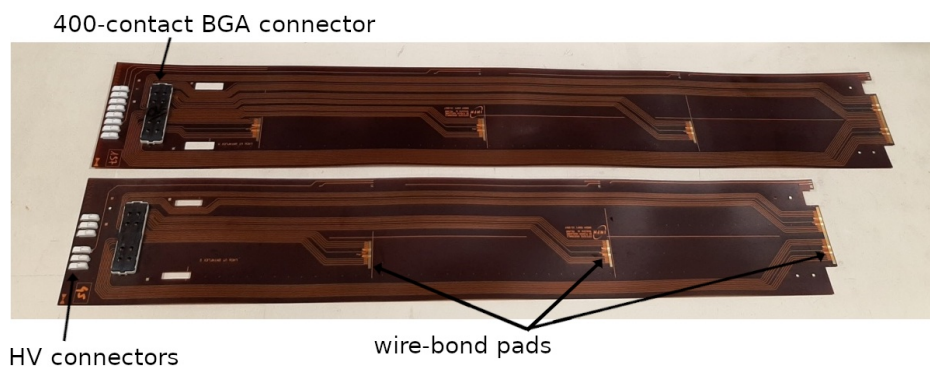[35]SEAF8-40-1-S-10-2-RA connector from SAMTEC™.

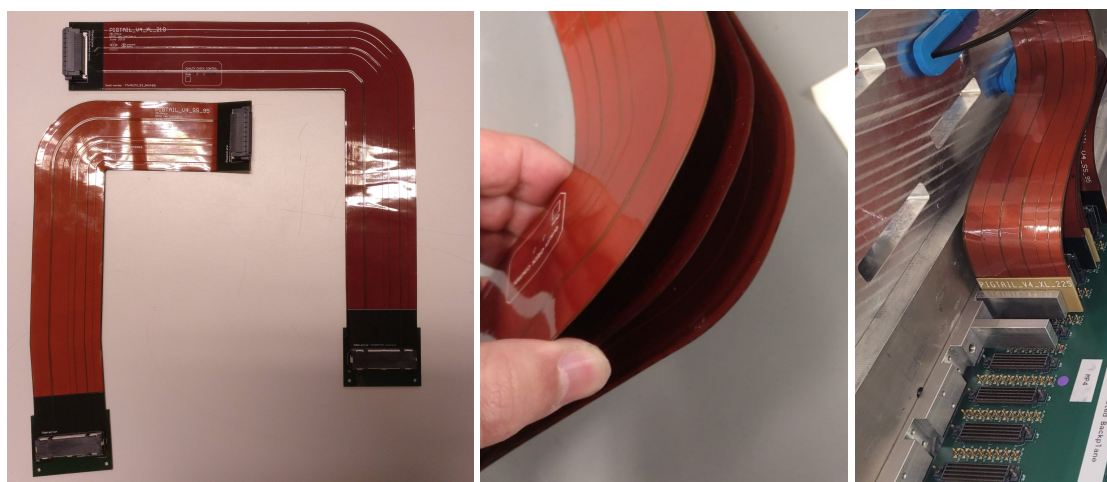**Figure 33.** Picture of one Short and one Medium dataflex cable.



**Figure 34.** Left: pigtails in two different shapes; middle: pigtail section showing the 3 subcables; right: installed pigtails.

FR4 stiffeners used to protect the fragile solder bumps under the connectors are also visible at both ends. The middle picture shows that each pigtail is composed of three flex subcables to confer sufficient flexibility to the pigtail for routing. The rightmost picture shows an example of bending when installed. Each subcable contains three copper layers of 35 µm thickness each, sandwiched between polyimide films and adhesive (adding up to about 0.25 µm of polyimide and 75 µm of adhesive per subcable). The inner layer is used for the differential pair signals and the outer layers for power lines. This provides good signal integrity and low resistance for the power traces. The pigtails cross the detector box walls through dedicated holes, which are sealed with foam to provide sufficient gas and light tightness.

### 5.6.3 Peripheral electronics

The periphery electronics processing interface (PEPI) units are responsible for readout and control of the detector. They connect to the staves via the pigtails. The full PEPI system comprises 24 backplanes, 24 pigtail power breakout boards (P2B2s), and 248 data and control boards (DCBs). The distribution of data, clock, and control signals is shown in figure 35. The GBTx, mounted on the DCBs, implements bidirectional links between the detector and the counting room with components
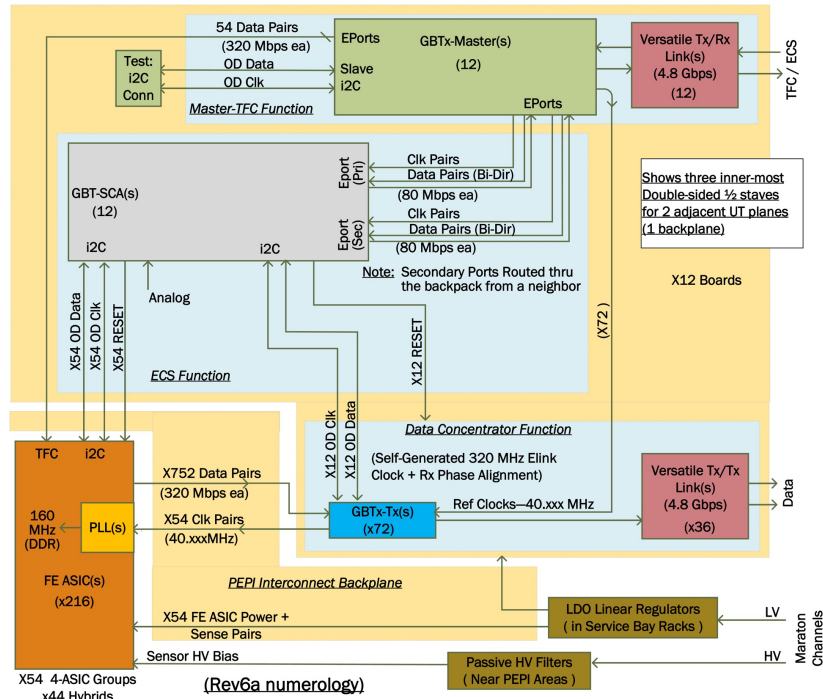
**Figure 35.** Schematic of LV, HV, data, fast and slow control signal distribution to UT FE electronics.

that are radiation hard up to 1 MGy [52]. Each PEPI unit houses 3 backplanes that route power, data, and control signals into and out of the detector volume and three P2B2s that route additional power. The backplanes support up to 12 DCBs that connect to the counting room via optical links.

### 5.6.4 Backplanes and power breakout boards

The routing of the backplanes balances the load carried by each DCB and organises the data to minimise resources needed for event reconstruction. A fully utilised backplane supports 6 half-staves read out by 12 DCBs. Due to space constraints in the area around the staves the size of the backplane is limited. Matching the twelve 400-contact pigtail connectors leads to a very high density of traces on the board. The signal routing is achieved by a 28-layer board with a PCB aspect ratio of 1:10 which is pushing the limits of manufacturability. There are two types of backplanes, referred to as true and mirror, with different layouts to accommodate the geometry of pigtails from the detector. Each type also has three variants: full, partial and depopulated. The full backplane is used for the high-occupancy region of the detector, while the depopulated one is used for the outer region of the detector planes, and partial backplanes are used everywhere else. While the full backplane transmits all data and control signals, on the partial and depopulated backplanes some traces are unused and grounded. The backplane routes all of the 1.5 V and 2.5 V power to the DCBs and some 1.2 V power to the SALTs. However, the majority of SALT 1.2 V power lines are routed through the P2B2 in order to leave room for data and control signals while maintaining the manufacturability of the backplane.

### 5.6.5 Data and control board

Each DCB supports one master GBTx, one GBT-SCA, six data GBTx chips, one VTRx [71], and two or three VTTx chips. The GBTx chip receives 320 Mbit/s sensor data from the e-ports connected

to the FE ASICs and repackages it into 4.8 Gbit/s data frames. In addition it controls data frames sent to and from the counting room via the GBT-SCA [16]. The GBT-SCA distributes slow control and monitoring signals to the detector. The VTRx incorporates a laser driver and optical receiver to convert between optical signals to and from the counting room and electrical signals to and from the GBTx, while the VTTx only transmits from the GBTx to the counting room. Communication between these chips is routed through the 16 signal and power layers of the DCB. In order to preserve high-frequency characteristics of the 4.8 Gbit/s communication between GBTx and VTTx/VTRx chips, a special laminate material[36] is used in place of standard FR4.

### 5.6.6  Data and control signal distribution

As shown in figure 35, between 6 and 12 e-ports of each data GBTx are connected to the FE ASICs. Each data GBTx receives data from either 2 or 4 ASICs belonging to the same 4-ASIC group, depending on the number of SALT e-ports and the location in the detector. The expected use of e-ports per SALT ASIC is shown by the numbers (3 to 5) in the boxes of figure 36 (left). Each pair of data GBTx transmits data to the counting room via a VTTx optical link. The master GBTx receives the LHC clock from the counting room via a VTRx optical link, and uses a PLL to generate a local clock with the same frequency (40 MHz). The master GBTx clock serves as a reference for the data GBTxs and GBT-SCA. One data GBTx per DCB transmits that clock directly to all FE read out by this DCB. The clock phase can be adjusted individually for each group of 4 SALT ASICs. TFC signals are received from the counting room via a VTRx optical link by the master GBTx and transmitted directly to all FE read out by the DCB. ECS signals are transmitted from the master GBTx to the GBT-SCA via a dedicated 80 Mbit/s e-port. The GBT-SCA distributes slow control signals to both the data GBTx chips and FE chips by the DCB. The GBT-SCA also distributes a GPIO signal that can reset the data GBTx and FE chips, and uses an ADC to monitor the voltage of the DCB and thermistors on the DCB and FE electronics.

### 5.7  Low voltage power

Low-voltage power for the UT hybrids and on-detector readout chain electronics is sourced from a bank of 12 power supplies[37] providing 144 8V/50A channels to four service bays situated around the UT detector. The service bays house 268 8-channel low-voltage regulator boards (LVRs) designed around the LHC4913 linear regulator chip. The LVRs segment the LV power into individual 3.6 A max derated channels with full differential remote regulation to the detector electronics at 1.26, 1.5 or 2.5 V nominal potential at load, depending on plug-in Current Control Mezzanine (CCM) boards. The individual channels also have the capability to be ganged in groups of two to boost the derated output to 7.2A for high-draw loads such as the GBTx chips and the 8-ASIC hybrid positions, with one channel providing the voltage regulation and the other sharing the output current with a specialised current-following CCM.

Groups of LVR channels served by a single power supply channel are referred to as *power groups* and these power groups have been optimised at the load side to maintain power independence of different backplanes. For the SALT ASICs, the non-negligible variations in voltage drop at different positions along the stave flex cable requires a *horizontal* grouping of hybrids at similar voltage

---

[36]Isola Terragreen$^{\text{TM}}$.

[37]Wiener$^{\text{TM}}$ MARATON HE LV.

| | Backplanes | | | | | | | |
| Inner | | | Middle | | | Outer | | |
|---|---|---|---|---|---|---|---|---|
| 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 |
| 4 | 4 | 4 | 3 | 4 | 4 | 3 | 4 | 4 |
| 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 |
| 3 | 3 | 3 | 2 | 3 | 3 | 2 | 3 | 3 |
| 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 |
| 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 |
| 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |

A A A A A A B C C B A A A A A

**Figure 36.** Left: illustration of the usage of dataflex cables on the different stave variants. Short cable shaded in blue, Medium in orange and Long in red. The modules hosted by the cable are highlighted with a coloured contour. The numbers 3, 4, 5 indicated the number of wire-bonded e-ports per chip on the given module. Right: arrangement of a single plane of the UT hybrids into power groups sharing a common LV output channel. Only one quadrant of the plane is shown. Boxes represent hybrids and the numbers (1 to 4) the power group.

drops rather than a *vertical* grouping along the stave. The grouping implemented in the detector is illustrated in figure 36 (right).

The LVR channels provide individual outputs referenced to shared LVR ground which is isolated from the local mechanical interface so that the floating power supply outputs remain floating through the LVR boards. The grounding of electronics is done at the backplane. Between the service bays and UT electronics, through approximately 8–10 m of cable, the LVRs must provide regulation across round-trip voltage drops ranging from 150 mV to 500 mV while maintaining smooth regulation performance, particularly at start-up. Extensive studies have been done to optimise the LVR regulation performance. The overall bandwidth is limited to a few hundred kHz. Faster voltage transients or oscillations are dealt with using passive decoupling networks at the hybrids or DCB input.

The UT LV power demands may be logically divided into SALT ASIC and DCB loads, with SALT ASICs requiring one 1.26 V LVR channel for VERA hybrids or a pair of ganged channels for the SUSI hybrids. DCB boards require a pair of ganged channels at 1.5 V for powering of the GBTx and GBT-SCA chips, and a single 2.5 V channel for VTTx and VTRx optical modules. Due to the low power demand for the latter, a single 2.5 V channel serves a pair of DCBs in the final UT system.

Monitoring and control for the LVRs is provided through a serial peripheral interface (SPI) to an on-board FPGA[38] which controls channel state and sequencing and provides state information. A dedicated mezzanine board provides a GBT-SCA interface to off-board electronics as well as monitoring analog outputs for real-time diagnostic information about voltage and current levels in the UT system. The mezzanines are in turn interfaced to the counting room via a control board in each LVR crate.

---

[38]Actel$^{TM}$ ProASIC3 FPGA.

## 5.8 Cooling

Evaporative $CO_2$ cooling is provided using a titanium tube embedded in the stave and running below the ASICs in a serpentine shape. The two-phase accumulator controlled loop is in common with the LHCb VELO detector, see section 2.4.3. The detector cooling system has to extract the power dissipated by the read-out chips, and keep the sensors at the target working-point temperature ($-5\,°C$) to prevent thermal runaway in presence of radiation damage. The total detector power to be extracted is expected to be about 4 kW, including 4192 ASICs (about 0.8 W/each), cables, sensors and environment. A safety margin of +25% is also considered in designing the system.

The core material of the box walls is a rigid polyetherimide-based polymeric foam.[39]  The foam is sandwiched in two carbon fibre reinforced polymer (CFRP) skins. The walls parallel (perpendicular) to the beam pipe are made of 24 mm (20 mm) thick core and 1 mm (0.5 mm) thick skins. A copper net (51 µm thick, 74% open area) covers the interior surfaces. The Be beam pipe is wrapped in a thermal blanket composed of 4 successive layers, a 0.2 mm PI PCB heater jacket directly on the beam pipe, ~ 0.1 mm PI tape, a 5 mm thick thermal insulation[40] and again ~ 0.2 mm PI tape. The interface between the box and beam pipe blanket is made of rigid polymeric foam,[39] EPDM foam and CFRP.

The detector box and beam pipe blanket were designed such that the temperatures on the surfaces never decreases below 13 °C, i.e. stays above the cavern dew point. Dry nitrogen, with a nominal dew point of $-50\,°C$, is circulated through the detector box via a supply line and a bubbler. Hot air from the electronics is forced to circulate near the pigtails outside the box to avoid condensation.

Four manifolds, two per side, one above and one below the detector box, distribute the cooling fluid to the detector, with one cooling loop per stave. Each loop is equipped with a restriction (0.2 mm orifice) at the entrance in order to avoid cross-stave effects arising from different heat loads. The nominal flow per stave is about 0.76 g/s. The pressure drop is driven by the inlet flow restrictions. Nominally, the total pressure drop in the detector is about 5 bar. On the stave, the temperature is stable to about $-0.5\,°C$. The $CO_2$ inlet temperature can be set from ambient temperature to $-30\,°C$.

The cooling requirements and properties of the detector have been studied using thermal and mechanical finite element models [72]. A result of the simulated thermal profile obtained on a central stave (variant C) is shown in figure 37 for an inlet $CO_2$ temperature of $-25\,°C$. The ASIC power is assumed to be 0.8 W per chip. This shows that the temperature of the innermost sensor (the most critically irradiated) can be kept well below the required temperature of $-5\,°C$.

## 5.9 Test beam results

To validate various aspects of the UT sensor design and readout, a series of beam tests were carried out. The test beam runs were focused on R&D and validation of three novel features of the sensor design (see section 5.4): (1) the double-metal region where the 190 µm sensor strip pitch is reduced to 80 µm to match the pitch of the custom UT ASIC input pads, (2) the top-side biasing scheme, and (3) the circular cut-outs of the Type-D sensors. Both unirradiated and irradiated sensors were tested, of full nominal size as well as miniature sensors. Sensor performance was measured, in particular signal-to-noise ratio (S/N), for exposures up to twice the maximum expected irradiation value over the lifetime of the UT detector.

---

[39]Airex R82.60™, from 3A Composites.

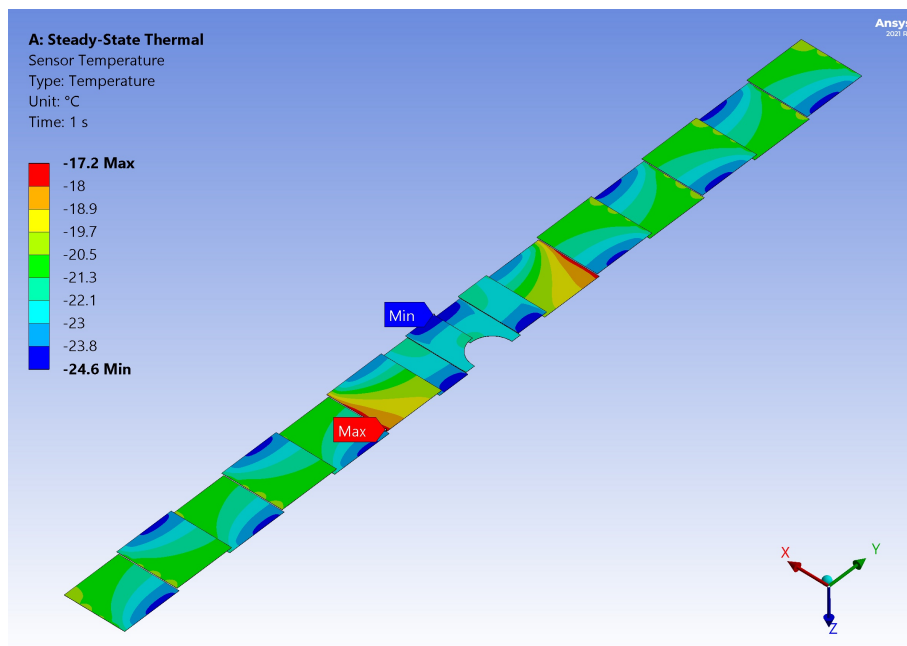[40]Pyrogel™ XTF from Aspen Aerogels.

**Figure 37.** Simulation of the thermal behaviour of a central stave (variant C) at an inlet $CO_2$ temperature of $-25\,^{\circ}\mathrm{C}$.

The key R&D areas listed above were studied in detail [73, 74]. Early test beam results, using the Beetle readout chip [75], showed that a S/N of about 12 can be expected for the Type-A sensors. The S/N of the Type-D sensors was measured to be higher, about 17, owing to the lower input capacitance of the shorter strips. The studies also demonstrated that the sensors with top-side biasing perform equally well as those that are biased directly from the back of the sensor. Studies of the Type-D sensors showed that they maintained excellent signal efficiency all the way up to the edge of the circular cut-out.

Test beam runs with early prototypes of the Type-A sensors showed a significant loss of signal efficiency in the PA region. Although this region covers less than 1% of the sensor's active area, the loss in efficiency was deemed to be unsatisfactory. Additional studies of the test beam data, as well as CAD device simulations,[41] were performed to conclude that the loss of efficiency was due to the capacitive coupling between the double-metal layers, inducing charge pick-up on other strips [76]. A new design was developed, with a thicker silicon dioxide insulating layer below the PA and with most of the PA metal moved outside the active area, and successfully tested. The results proved that the PA region exhibited no significant loss of hit efficiency.

A final test beam run [77] was conducted which successfully demonstrated the performance of a Type-A sensor with SALT v3.0 128-channel readout chips. The most relevant results are shown in figure 38. The left panel shows the distribution of collected charge, in ADC counts, for tracks with normal incidence. Fitting the distribution with a Landau function convoluted with a Gaussian function, results in a most probable value of 11.1 ADC counts, while the measured common-mode subtracted noise was about 0.9 ADC counts. Thus, a S/N of about 12 is obtained. An irradiated sensor was also tested and found to have a S/N about 10% lower, which is still large enough to meet the UT requirements on signal efficiency and noise hit rejection at the end of life of the detector.

---

[41]Simulatios were performed with Synopsys$^{\mathrm{TM}}$ Sentaurus TCAD.
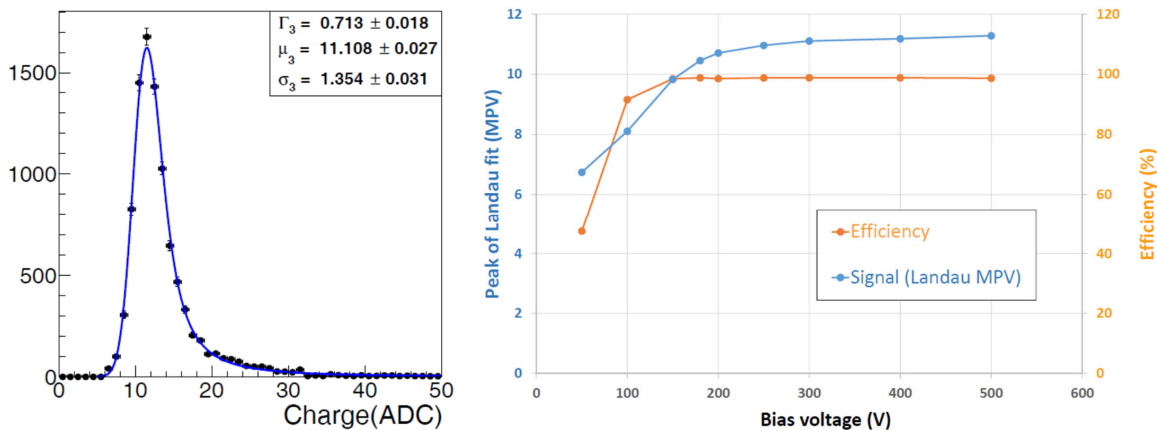
**Figure 38.** Test beam results for a Type-A unirradiated sensor for tracks with normal incidence. Left: distribution of collected charge (in ADC counts) at 300 V bias. The data are fitted with a Landau distribution convoluted with a Gaussian function. Right: most probable value of the Landau fit result and hit efficiency versus the applied bias voltage. Reproduced from [77]. CC BY 4.0.

## 5.10 Slice test

A full electronics read-out test on a prototype stave, dubbed *slice test*, was set up to validate the detector design with realistic power distribution and grounding. To read out a complete stave, two partially routed backplanes were used together with four DCBs and associated power distribution break-out boards. The LHCb MiniDAQ system was used to control the system, distribute the timing signals and send trigger commands to the FE electronics, and to process the data from the FE electronics. The stave was cooled using a bi-phase $CO_2$ cooling system (a reduced size version of the one installed in the LHCb cavern). The stave was operated in a light-tight box at temperatures between $-30$ and $15°C$.

Data were collected with different numbers of hybrids powered and configured to nominal settings, in order to study possible correlated noise effects. ASICs on each hybrid were configured individually, or all at the same time. Finally, tests with all hybrids powered and configured simultaneously were carried out. An example result of the noise before and after mean common mode suppression is shown in figure 39 for the innermost module on a Medium dataflex cable. The raw and common mode subtracted noise were around 1.1 and 0.9 ADC counts, respectively, when only that particular single module was configured. When all modules were operating, the raw noise increased by around 10% while the difference in the noise after common mode subtraction was negligible.

## 5.11 Simulation and reconstruction software

Unlike for the TT detector in the original LHCb experiment, the UT DAQ electronics does not cluster adjacent strip hits. All strip hits are saved with their ADC pulse heights. In spite of somewhat complicated UT raw bank organisation, driven by the readout granularity and limited computing resources available in the TELL40 boards, fast raw data processing is achieved via a smart iterator with the knowledge of the raw data structure, pointing directly to the raw bank in computer memory. In the simplest implementation, a strip hit in the raw bank is promoted directly to a hit with a space location used by the LHCb tracking software. Since only, 10% of tracks passing a UT sensor are expected to light up more than one strip, it is possible that this will be an optimal decoding scheme in the first level software trigger. A
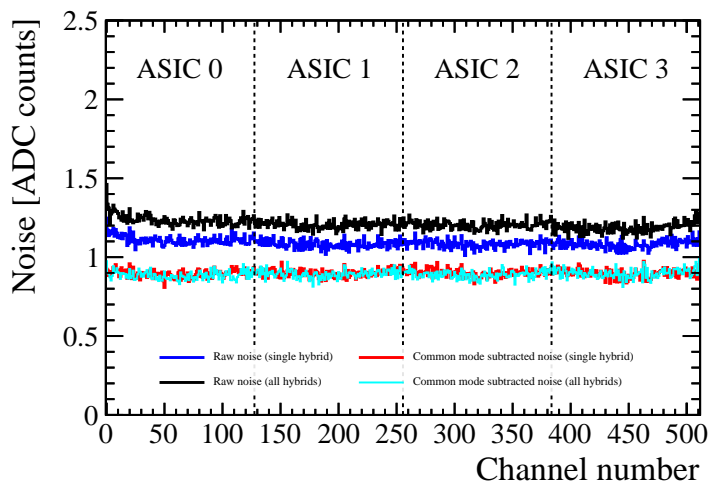
**Figure 39.** The raw and common mode subtracted noise measured in a single hybrid with only that hybrid powered and configured (blue and red curves, respectively) compared with the noise measured in the same hybrid when all hybrids in the stave were operating (black and cyan).
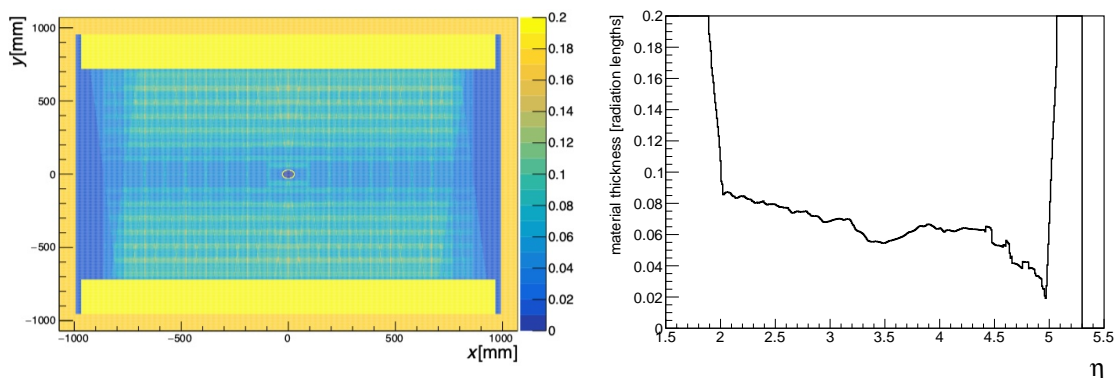


**Figure 40.** Material scan through the UT detector, with thickness given as a fraction of a radiation length. Left: thickness map in $xy$ plane for normal track incidence. Right: thickness as a function of pseudorapidity ($\eta$) as seen from the interaction point.

second version of the iterator over the UT raw data has been developed, which clusters on the fly while advancing over the raw bank. It checks for adjacent strip hits, by looking up the next element in the raw data, and returns the highest pulse-height strip, while suppressing the others. Detailed timing studies will be needed, once exact performance of the UT ASICs is known under running conditions, to determine if such algorithm can run in the first level or must be deferred to the second level of the software trigger.

A GEANT4-level simulation software was developed for the UT system. The detector representation includes technical details of the final design. The hierarchy of geometrical volumes has been structured to represent segmentation of the mechanical support and of the final assembly with detector software alignment in mind. For example, layers were split into left and right half-layers to follow the actual retractable C-frame design. A map of the material in the acceptance was obtained from simulation. The UT material thickness, expressed as a fraction of radiation length, is shown in figure 40. Compared to its predecessor tracker (the TT detector), the UT material in the region
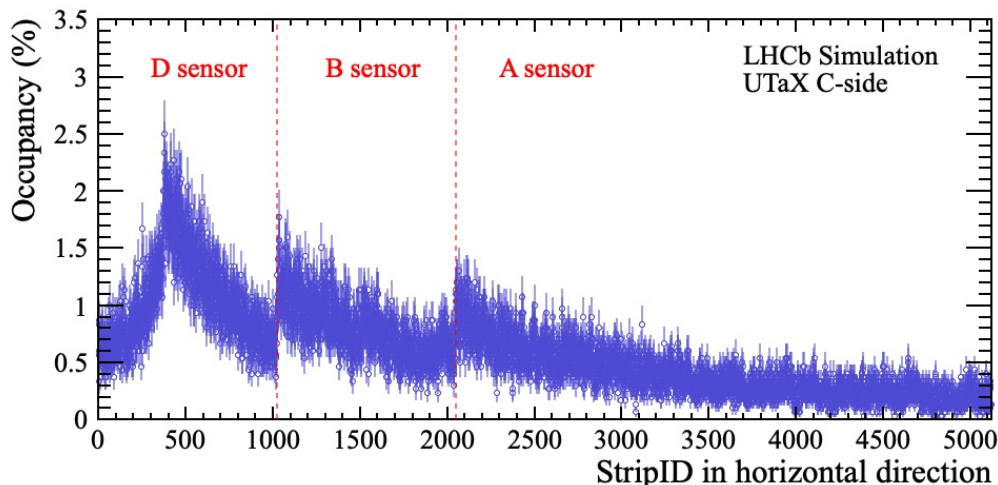
**Figure 41.** Expected UT strip occupancies for minimum bias events in the sensors near the detector midline ($y = 0$).

closer to the beam pipe ($4.4 < \eta < 5$, where the track density is maximal) has been reduced from about 15% to 4% of a radiation length.

The expected strip occupancies in the UT near the detector $y = 0$ mid line are illustrated in figure 41, where the average fraction of minimum bias events leaving a hit in the strip is shown as a function of strip number. The plot focuses on the Side C of the first UT layer (UTaX). The occupancies are kept below 3% even in the busiest region and are below 1% in most of the detector.

# 6 Scintillating fibre tracker

## 6.1 Overview

The tracker, located downstream of the LHCb dipole magnet, is responsible for charged particle tracking and momentum measurement. A momentum resolution and track efficiency for *b*- and *c*-hadrons comparable to the ones obtained in Run 1 and Run 2 must be achieved, while working in an environment with higher particle density. To cover the nominal LHCb acceptance it must have an area of about $6\,\mathrm{m} \times 5\,\mathrm{m}$ in the *xy* plane.

### 6.1.1 Detector requirements

The design of the tracker must take into account the following requirements:

- Performance: the tracker should provide a single hit position resolution of better than 100 μm in the magnet bending plane and a single hit reconstruction efficiency better than 99%.

- Rigidity: the mechanical stability of the detector must guarantee that the positions of the detector elements are stable within a precision of 50 (300) μm in *x* (*z*); the detector elements should also be straight along their length within 50 μm.

- Material budget: to limit further multiple scattering and secondary particle production, each of the 12 layers should not introduce more than 1% of a radiation length.
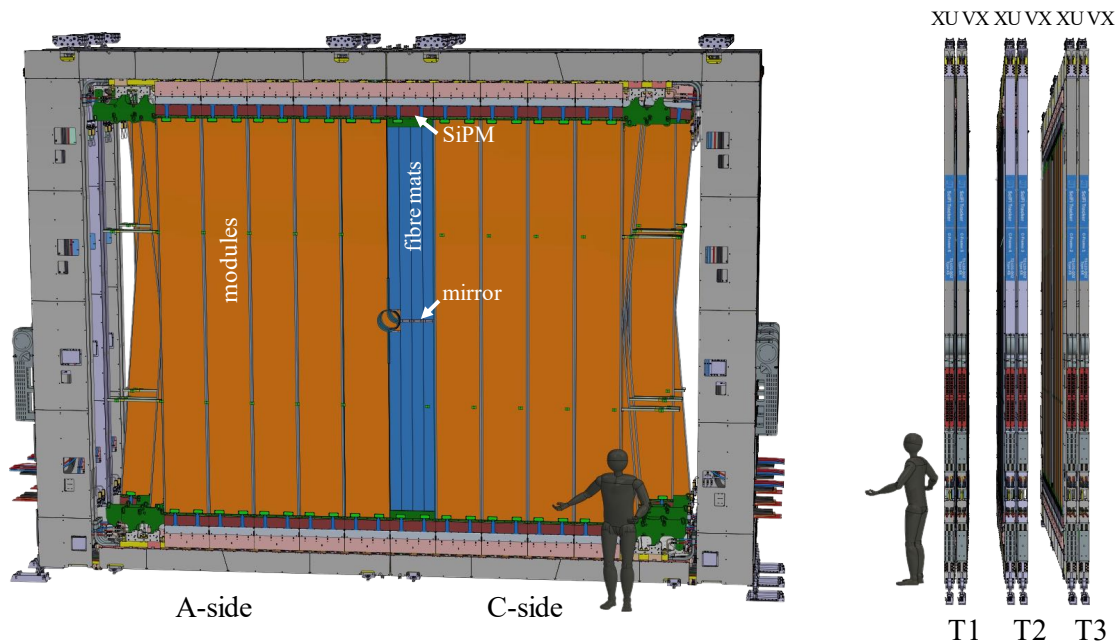
**Figure 42.** Front and side views of the 3D model of the SciFi Tracker detector.

- Radiation hardness: the tracker should operate at the desired performance over the lifetime of the experiment, where $50\,\text{fb}^{-1}$ of integrated luminosity is expected to be collected.

- Granularity: the tracker must have an occupancy low enough so that the hit efficiency is not impacted with an instantaneous luminosity of $2 \times 10^{33}\,\text{cm}^{-2}\,\text{s}^{-1}$ [78, 79].

A tracker design based on scintillating fibre (SciFi) technology with SiPM readout was chosen to fulfil these requirements.

### 6.1.2 Detector layout

The SciFi Tracker acceptance ranges from approximately 20 mm from the edge of the beam pipe to distances of $\pm 3186$ mm and $\pm 2425$ mm in the horizontal and vertical directions, respectively, with a single detector technology based on 250 μm diameter plastic scintillating fibres arranged in multilayered fibre mats. In total there are 12 detection planes arranged in 3 stations (T1, T2, T3) with 4 layers each in an $X - U - V - X$ configuration, as shown in figure 42. The $X$ layers have their fibres oriented vertically and are used for determining the deflection of the charged particle tracks caused by the magnetic field [79]. The inner two stereo layers, $U$ and $V$, have their fibres rotated by $\pm 5°$ in the plane of the layer for reconstructing the vertical position of the track hit.

Each station is constructed from four independently movable structures referred here as C-Frames, with two C-Frames on each side of the beam pipe. The carriages of the C-Frames move along rails fixed to a stainless steel bridge structure above the detector, supported by stainless steel pillars. To simplify production, each station is built from identical SciFi modules about 52 cm wide and spanning the full height, except for a few modules near the beam pipe. The T3 station is instrumented with six modules on each C-Frame. T1 and T2 stations have one less module on each side for instrumenting the smaller acceptance at those locations due to the opening angle of LHCb.

### 6.1.3 Detector technology

The detector modules are constructed as a honeycomb and carbon-fibre sandwich containing eight ~ 2.4 m long and ~ 13 cm wide SciFi mats made from six staggered layers of fibres. A thin mirror is glued to the fibre end to reflect additional light back to the readout side [80]. In the experiment, these mirrors are located near the $y = 0$ plane. From there, four fibre mats point upward and four downward, spanning a total height of almost 5 m. Near $x = 0$, a few special modules are used in order to take into account the presence of the beam pipe. Those modules have one mat shortened to accommodate the beam pipe radius, as seen in the centre cutaway module in figure 42. A more detailed description of the modules is found in section 6.3.

The optical signal from the scintillating fibres are detected by 128-channel arrays of SiPMs with a channel pitch of 250 μm. The SiPMs are discussed in section 6.4. At the readout end of each module, 16 SiPM arrays are bonded to a 3D printed titanium alloy cooling bar, which is aligned to the four fibre mats and housed in a cold-box, described in section 6.7.

### 6.1.4 Detector irradiation

Initially, the SciFi Tracker has a mean light yield about 18–20 photoelectrons for particles passing perpendicularly through the detector plane near the mirrors.[42] Irradiation will reduce the light yield. Fluka simulations have shown that up to 35 kGy of ionising radiation dose can be expected in the fibres in this region by the end of the lifetime of the detector. The total received dose drops off sharply to about 50 Gy at the readout end of the fibres, as seen in figure 43. Several irradiation campaigns were performed to study the effect on light output and a signal loss of about 40% is expected for hits in the most irradiated regions towards the end of the tracker's lifetime. This will still exceed 10 photoelectrons allowing for a single hit detection efficiency of 99% in the sensitive regions, assuming low selection thresholds (4 photoelectrons or greater), and a single hit position resolution measured to be 70 μm which is needed for efficient tracking of charged particles behind the magnet of LHCb.

At the position of the SiPMs, the expected total ionising dose over the lifetime of the experiment is relatively low, 50 Gy, such that the main concern will be the non-ionising energy loss (NIEL) which damages the silicon crystal lattice. A wall of borated polyethylene (5% w/w) between the RICH2 and the ECAL has been installed to reduce the back-scatter of low energy neutrons from the calorimeters by a factor of more than two in the region of the SiPMs (see also section 2.4.2). The inner region of shielding installed ±1 m around the beam pipe is 30 cm thick with the outer region having a thickness of 10 cm. Fluka simulations of the LHCb Upgrade show that a collected fluence of up to $6 \times 10^{11}$ 1 MeV$n_{eq}$/cm$^2$ is expected. The NIEL damage will increase the rate of pixel avalanches per channel initiated by thermal electrons, occurring without incident light, and is commonly referred to as dark noise or dark count rate (DCR). The SiPMs are to be cooled below $-40\,°C$ as the damage increases in order to suppress the rate of accidental clusters and maintain the required detector performance.

A significant portion of the detector infrastructure is dedicated to managing the radiation effects on the fibres and the SiPMs, as well as handling the large data rates generated by the increased instantaneous luminosity of the LHCb upgrade where, after zero suppression, up to 20 Tbit/s of data are sent to the DAQ from the SciFi Tracker FE electronics.

---

[42]The distribution is Poisson-like in its shape, due to variations in the total path of the charged particle through the scintillating fibre matrix, fluorescence processes, and saturation from large ionisation energy deposits.
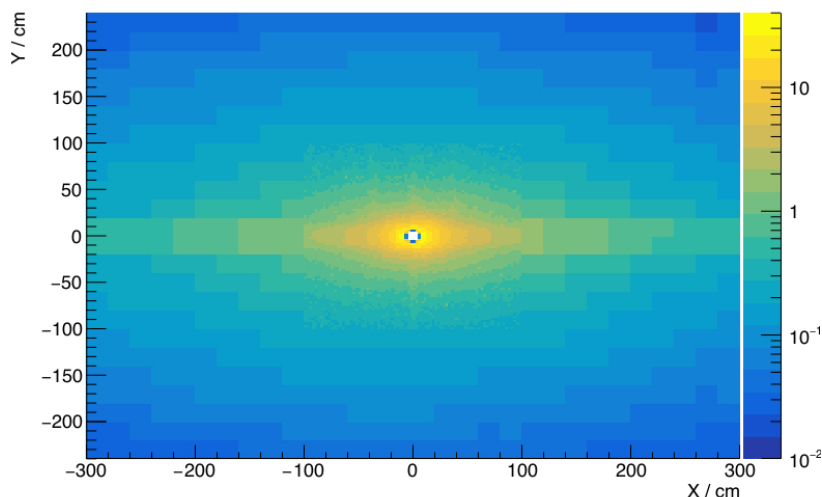
**Figure 43.** Map of the total expected ionising dose in kGy for an integrated luminosity of $50\,\text{fb}^{-1}$ at the T1 station of the SciFi Tracker from FLUKA simulations of the LHCb detector.

## 6.2 Scintillating fibres

Blue-green emitting double-clad plastic scintillating fibres with a 250 µm diameter were chosen for the LHCb SciFi Tracker.[43] The production of the over 11 000 km of fibre needed for the tracker was delivered on time and of a very high and constant quality, which led to negligible rejection rates.

The manufacturer states an intrinsic light yield of 8000 photons per MeV of ionisation energy deposited. The decay time constant was measured to be 2.4 ns, slightly faster than the 2.8 ns declared by the manufacturer [81, 82]. The double-clad fibre achieves a minimum trapping efficiency of 5.3% per hemisphere through total internal refection with the two cladding layers. The fibres have a mean nominal attenuation length of approximately 3.5 m. However, the transparency of the fibres will degrade due to the received ionising dose. This effect has a strong wavelength dependence and may result in losses that peak in the UV/blue regions. The emission spectrum of the fibre has a maximum at 450 nm. It has also been observed that the attenuation length degrades by 1–2% per year due to radical production resulting from the interaction of oxygen with the polymer [83–85].

### 6.2.1 Fibre quality control

The fibre supply was delivered in 48 individual shipments (batches) in weekly (100 km) or bi-weekly (300 km) intervals on 12.5 km spools. After the arrival of each shipment, fibre samples for quality control were prepared and the individual fibre spools were processed on a fibre scanner developed in-house. A detailed description of the quality control process can be found in ref. [85].

The attenuation length of the fibres was measured for each batch both by the manufacturer and upon reception by LHCb, using light from a photodiode and fitting the yield in a range between 1 and 3 m from the light source to avoid the contribution of a second exponential with shorter attenuation length.

The attenuation length was monitored over the full production time of almost two years. Averages over all spools (typically 24) of one shipment were calculated. Apart from an initial increase of the

---

[43]Fibres SCSF-78MJ by Kuraray™.

attenuation length during the preseries and early main-series production, the attenuation length was very stable throughout the full production, with a mean of $3.5 \pm 0.2$ m, well above the acceptance limit of 3 m.

The light yield was measured in a dedicated quality control setup (described in ref. [86]) at 2.4 m from the position of a $^{90}$Sr radioactive source (with magnetically selected 1.1 MeV electrons). An improvement of the attenuation length observed in the preseries production was also seen in the light yield. Otherwise, the single-fibre light yield, extrapolated to zero distance from the photon detector, was very stable around a mean of about $6.8 \pm 0.5$ photoelectrons, and was always above the acceptance limit of 5 photoelectrons.

### 6.3 Fibre mats and modules

#### 6.3.1 Fibre mats

As found in ref. [87] fibre mats are produced by winding six layers of fibres on a threaded winding-wheel with a diameter of approximately 82 cm. The winding puts the fibres in a regular hexagonal matrix with a horizontal fibre pitch of 275 µm. Titanium-dioxide loaded epoxy (20% w/w) is used to bond the fibres to each other and to provide some shielding for cross-talk and for the diffusion of light escaping the cladding. Thin black polyimide foils are bonded to the fibre mat on both sides to provide mechanical stability as well as some amount of light shielding.

Additionally, the fibre mats are cut precisely to the desired length (2424 mm) and width (130.6 mm) after having a polycarbonate end-piece added with precision holes for aligning the SiPMs to the fibre mat at the readout end, as seen in figure 44. The mirror end has two 2 mm-thick polycarbonate end-pieces for alignment and optical milling. The foil mirror[44] is bonded with epoxy to them and the fibre mat.
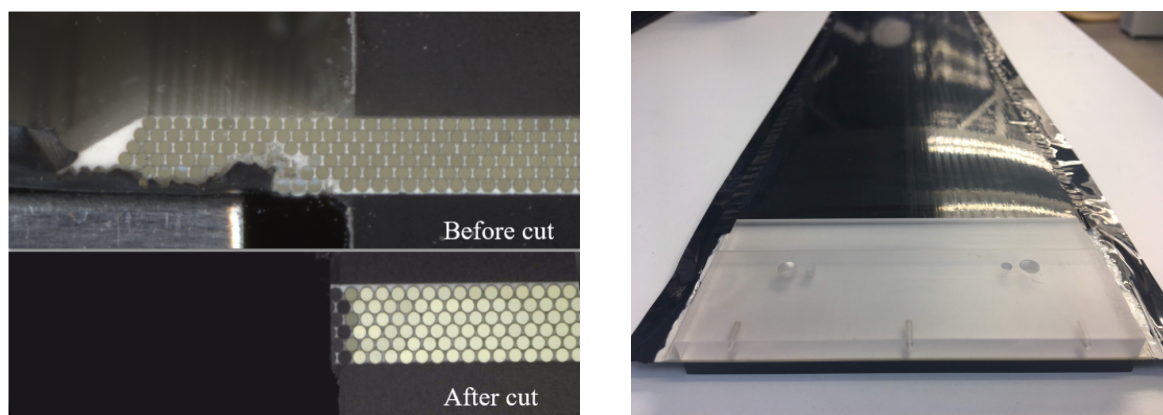


**Figure 44.** Left: view of a fibre mat with a microscope before and after the side fibres are cut away. Right: a photo of a fibre mat with the polycarbonate end-pieces and SiPM alignment holes. Reproduced from [87]. CC BY 4.0.

#### 6.3.2 Modules and C-Frames

Modules are built from eight fibre mats and two half-panels. Each half-panel is made of 19.7 mm thick polyaramid honeycomb cores laminated on the outer side with a single 0.2 mm carbon-fibre reinforced polymer (CFRP) skin. In order to minimise internal stresses, which could deform the module, a symmetric design was chosen, with the two half-panels sandwiching the fibre mats. The

---

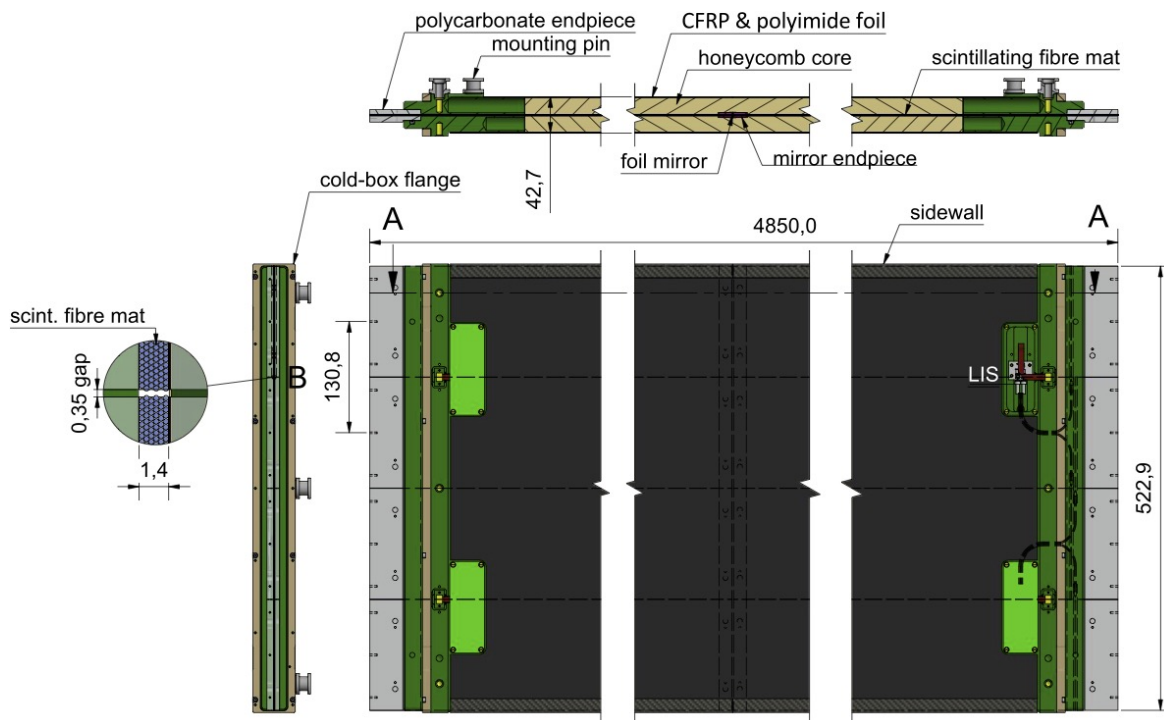[44]Enhanced Specular Reflector from 3M$^{\text{TM}}$ corp.

**Figure 45.** A cross-section and projections with dimensions (in mm) of a scintillating fibre module and its components. The outlines of the fibre mats and light injection system (LIS) are shown with dashed lines. The nominal gap between fibre mats is shown in a zoomed inset on the left. Reproduced from [87]. CC BY 4.0.

cross-section and design of a module is shown in figure 45. Each C-Frame contains two layers of modules with 10—12 modules in total, depending on the station. A full description of the module production can be found in ref. [87].

Special modules contain a circular cut-out in the half-panels to accommodate the radius of the beam pipe with a 20 mm safety radial gap. The top (bottom) fibre mat at this location has been shortened by 116 (116) mm for $X$-modules, and 139 (93) mm for stereo modules to accommodate the cutout resulting in a rectangular hole in the detector acceptance around the beam pipe. This results in three types of modules: nominal, $X$-beam-pipe, and Stereo-beam-pipe.

The modules are fixed in position on the C-Frames by the top-centre mounting pin. Five other mounting pins on the module constrain the rotations. The modules in the $X$-layer are positioned with a pitch of 531 mm, such that there is a gap of 8 mm between the edge fibres of neighbouring modules. For the stereo modules the horizontal pitch is 532 mm. There is also a 0.35 mm gap between neighbouring fibre mats within a module, as shown in the inset of figure 45, though the inefficiency gap is effectively slightly larger due to the likelihood of damaging the edge fibre during the cutting of the fibre mats, as shown in figure 44 (left bottom). A 2 mm gap between top and bottom fibre mats is also present. The total geometric inefficiency of one layer is approximately 1.7% (within the nominal acceptance).

### 6.3.3 Material budget

A module has a minimum estimated material thickness of 1.03% of $X_0$ for perpendicular tracks. The individual contributions are shown in table 7. The carbon-fibre U-shaped foils that enclose the sidewalls add an additional 0.145% of $X_0$ where they overlap. The region of the mirror end-piece

**Table 7.** Material budget contributions from the scintillating fibre module components. Densities and radiation lengths are taken from ref. [88].

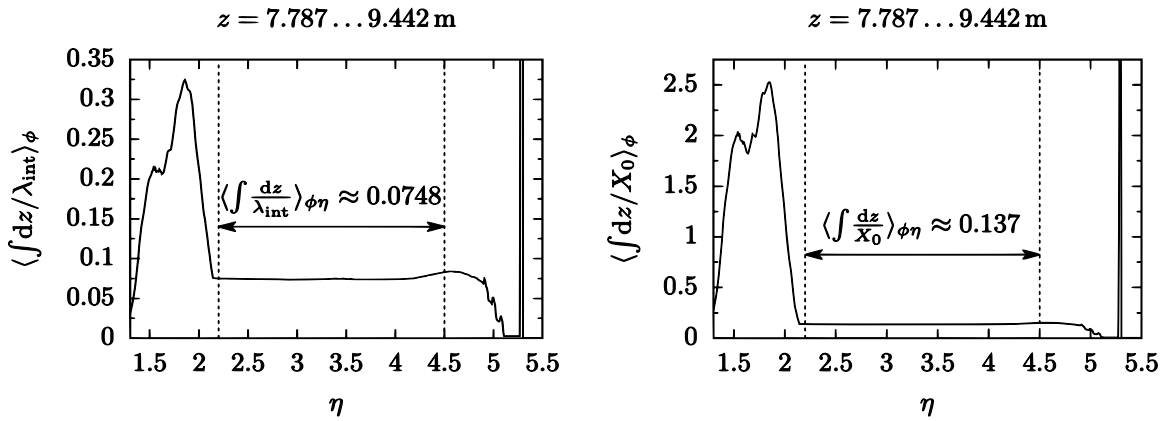| Material | Thickness ( μm) | Density (kg/m$^3$) | Layers per module | $X_0$ ( cm) | Fraction of $X_0$ (%) |
|---|---|---|---|---|---|
| *widespread* | | | | | |
| Fibre mat | 1350 | 1180 | 1 | 33.2 | 0.407 |
| Honeycomb core | 19700 | 32 | 2 | 1300 | 0.303 |
| CFRP skin | 200 | 1540 | 2 | 27.6 | 0.145 |
| Glue | 260 | 1160 | 2 | 36.1 | 0.144 |
| Polyimide foil | 25 | 1410 | 4 | 35 | 0.029 |
| *local* | | | | | |
| Sidewalls | 200 | 1540 | 2 | 27.6 | 0.145 |
| Mirror polycarbonate | 2000 | 1200 | 2 | 34.6 | 1.15 |



**Figure 46.** The traversed amount of material, averaged over azimuthal angle $\phi$, in units of (left) hadronic interaction length and (right) radiation length as a function of $\eta$ for a sample of simulated $B_s^0 \to \phi\phi$ decays. The total average is also shown for $\eta$ between 2.2 and 4.5 (range limited to the acceptance of the SciFi Tracker shown with dashed lines).

of the fibre mat adds 4 mm of polycarbonate, corresponding to an additional 1.15% of $X_0$. In the acceptance, a particle passing through the 12 SciFi Tracker layers traverses a minimum of 12.4% of an $X_0$, or a bit more depending on angle and exact location.

Results from particle tracking simulations using a sample of $B_s^0 \to \phi\phi$ decays, shown in figure 46, estimate the material budget for the three stations, averaged over pseudorapidities in the range $2.2 < \eta < 4.5$, to be 7.48% of a hadronic interaction length, and 13.7% of a radiation length.

## 6.4 Silicon photomultiplier assemblies

The tracker will use a total of 524 288 SiPM channels implemented as 4096 128-channel arrays to detect the light from the scintillating fibres. The production version of the array, manufactured by Hamamatsu™, is here referred to as *H2017*. Each channel comprises $4 \times 26$ pixels connected in parallel with a pixel size of 57.5 μm $\times$ 62.5 μm resulting in a single channel with dimensions of
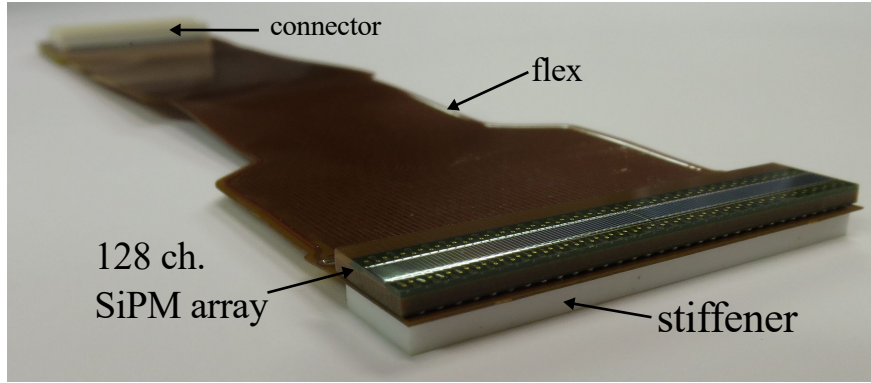
**Figure 47.** An H2017 SiPM array bonded to a flex cable. The white stiffener is visible on the lower side of the flex cable.

230 μm × 1625 μm. The channel pitch is 250 μm. The array is composed of two 64-channel dies wire-bonded to a PCB with a 220 μm gap between the dies.

A 100 μm thin transparent epoxy window protects the silicon and bonding wires, such that the array can be pressed against the scintillating fibres. The PCB is also instrumented with a Pt1000 temperature sensor on the back to monitor the temperature. The 32.6 mm wide PCB is bump bonded to a PCB flex cable which has been equipped with a connector and a stiffener, glued to the back side, as seen in figure 47.

### 6.4.1 Bias and photon detection efficiency

The pixel of the SiPM is a reverse-biased photodiode operated in Geiger mode. The breakdown voltage, $V_{BD}$, at which the avalanche process can occur has been measured to have a mean value of 51.75 V and varies by up to ±300 mV across the array [89]. The breakdown voltage for every channel in every array has been measured and stored in a database for use during operation. The nominal operating voltage in the SciFi Tracker is 3.5 V above $V_{BD}$; details of the SiPM bias distribution system are reported in ref. [90] The $V_{BD}$ is linearly dependent on the temperature with a coefficient of $(60 \pm 2)$ mV/K requiring a stable cooling system, as several parameters depend on the value of the bias over $V_{BD}$. A signal multiplication (gain) value of $(1.01 \pm 0.01) \cdot 10^6$ V$^{-1}$ was measured for these devices.

The photon detection efficiency (PDE) for the H2017 SiPM was measured with both current and pulse frequency methods described in ref. [89] and found to have a peak value of $(43.5 \pm 3.5)$ % for nominal bias.

The avalanche is quenched as the bias voltage drops below $V_{BD}$ due to the increased current over a so-called quench resistor. The resistors have a range of 470 kΩ−570 kΩ in the H2017 devices with a 25 Ω difference between odd and even channels. No simple correlation was observed between the quench resistor value and $V_{BD}$.

### 6.4.2 Cross-talk and correlated noise

Infrared photons produced in the primary pixel avalanche can be absorbed in neighbouring pixels causing additional pixels to fire, either in time or with a short delay due to the depth and location of the photon absorption. This is referred to as direct and delayed cross-talk. A reduction of the cross-talk over older SiPM designs has been achieved by the addition of trenches in the silicon between pixels.

Trapped charge carriers in the silicon of a pixel that has previously fired can result in a delayed pixel avalanche once the bias in the pixel exceeds the breakdown voltage and the trapped charge is freed. This is called after-pulsing. Cross-talk and after-pulsing are collectively referred to here as correlated noise.

The correlated noise probabilities were measured as a function of the overvoltage [89]. The results are shown in figure 48 (left). Direct cross-talk has a probability of 3.3% of occurring at nominal overvoltage. Delayed cross-talk has a probability of 3.7% and an exponential decay time constant of $(17.7 \pm 0.4)$ ns for the H2017 sensors.[45] After-pulsing is negligible in these devices.
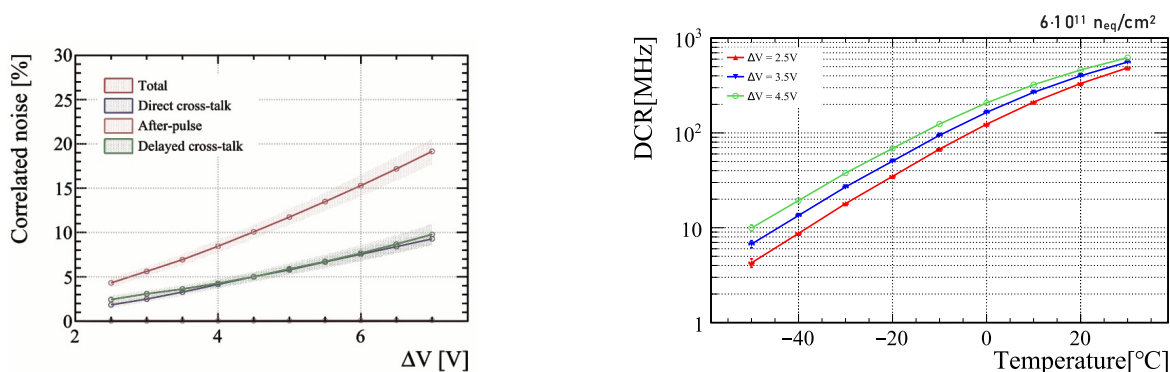


**Figure 48.** Left: correlated noise probabilities for an H2017 detector as a function of $\Delta V$. Right: dark-count rate for an irradiated SiPM as a function of temperature for three overvoltages. Reprinted from [91], Copyright (2020), with permission from Elsevier.

### 6.4.3 Dark noise

The photodetectors are expected to accumulate a total fluence of $6 \times 10^{11} 1\,\mathrm{MeV}\,n_\mathrm{eq}/\mathrm{cm}^2$ and 50 Gy of ionising radiation over the lifetime of the experiment. The ionising dose at this level has been shown to not have significant impact on the SiPM performance. The accumulated displacement damage, however, increases the rate of single pixel avalanches caused by thermal excitation of electrons, which increases the DCR by several orders of magnitude over the lifetime of the detector. A DCR of 14 MHz per channel is expected towards the end of life of the detector, as seen in figure 48 (right).

Coupled together with cross-talk, single-pixel dark counts have some significant probability to create signal amplitudes of two or more pixel avalanches which appear similar to low amplitude signals caused by particle tracks. Channels with dark-noise amplitudes above a set threshold will be mistaken for real signal clusters unless they are rejected by other means, such as the clustering algorithm described in section 6.5.5. In addition to screening neutrons with the PE shielding, a further reduction of the DCR by a factor of 100 is achieved by cooling the irradiated SiPMs from room temperature to $-40\,^\circ\mathrm{C}$, as it is described in section 6.7.

### 6.5 Front-end electronics

The FE electronics consists of three types of boards: the PACIFIC board, the Clusterisation board, and the Master board, shown in figure 49 (left). The PACIFIC board performs the digitisation of the analog signals received from the SiPM. The Clusterisation board performs zero-suppression and clustering

---

[45]The first batch that was delivered (about 10% of the total), is slightly noisier, with direct and delayed cross-talk probabilities of 3.2% and 5.6%, respectively.
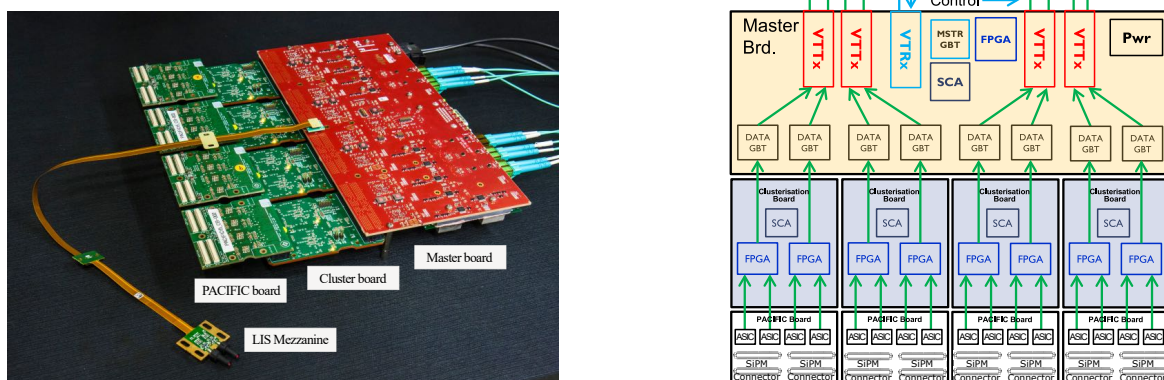
**Figure 49.** Left: picture of assembled Master, Clusterisation and PACIFIC boards. Reprinted from [91], Copyright (2020), with permission from Elsevier. Right: corresponding schematics of signal data routing. Reproduced from [87]. CC BY 4.0.

of the signals. The Master board serves multiple roles, such as the distribution of control and clock signals, monitoring, as well as distributing the low voltage and SiPM bias to the other boards.

Figure 49 (right) shows the layout and data path of one set of FE electronics. Each end of a SciFi module requires two sets of electronics, each one consisting of 4 PACIFIC boards, 4 Clusterisation boards and one Master board, to digitise and transmit data from 16 SiPMs. The boards are fixed to a thick aluminium chassis which is cooled by chilled demineralised water.

### 6.5.1 Data flow summary

SiPM avalanche pulses are processed and digitised in the PACIFIC ASIC (described in more detail in section 6.5.2) with a system of three hierarchical threshold comparators with four possible results, which are encoded in a 2-bit output word. Four channels are serialised together and transferred from the PACIFIC to a Microsemi Igloo2™ FPGA for clusterisation at a rate of 320 Mbit/s. The data output after serialisation of the PACIFIC has a total rate of 10.24 Gbit/s delivered for each SiPM array. The FPGAs are programmed to perform a cluster search algorithm per SiPM array, as discussed in section 6.5.5, and to calculate the position for each found cluster.

Clustered data from each bunch crossing are labelled with an event header and transferred to the data serialiser where the SciFi Tracker FE electronics follows the architecture described in section 2.3. Eight GBTx chips serialise the eight data streams on each Master board (4.8 Gbit/s per link) and the VTTx transmitters [15] push the data to the DAQ. In total, the SciFi Tracker transmits approximately 20 Tbit/s to the BE over 4096 data links. A bi-directional VTRx is used for the ECS and the TFC commands for each set of electronics. The GBTx, VTTx, and VTRx are described further in detail in section 10.

### 6.5.2 PACIFIC ASIC

The PACIFIC is a 64-channel ASIC with current mode input and digital output developed to read out SiPMs. It is implemented in a 130 nm CMOS process.[46] Each of the 64 channels contains an analog processing, digitisation, slow control, and digital output synchronised with the 40 MHz bunch clock

---

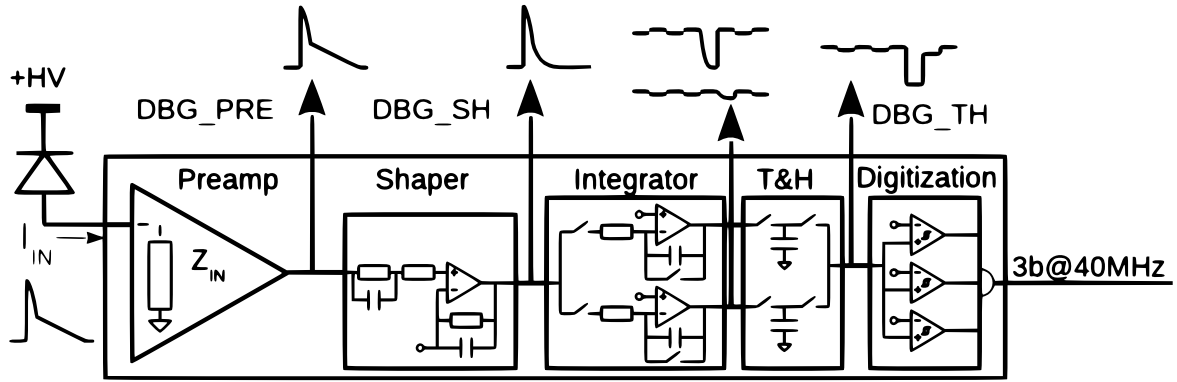[46]By TSMC™ Taiwan Semiconductor Manufacturing Company.

**Figure 50.** PACIFICr5q Channel block diagram. Reproduced with permission from [92].

of the LHC. Several versions of the PACIFIC ASIC were developed. The version used in the SciFi Tracker is the PACIFICr5q. A detailed description of the ASIC can be found in ref. [93].

A simplified representation of the analog processing is shown in figure 50. It consists of a preamplifier, a shaper, and an integrator. The interleaved, double-gated integrator operates at 20 MHz to avoid dead time as one integrator is in reset while the other collects the signal. The two integrator outputs are merged by a track-and-hold to provide a continuous measurement.

The output voltage of the track-and-hold is digitised using three comparators acting like a nonlinear flash ADC. The result of the three comparators is encoded into a two bit datum. Four channels are serialised together. The design of the PACIFICr5q is complemented by auxiliary blocks such as voltage and current references, charge injection, control DACs, and power-on-reset circuitry.

An additional feature allows for the voltage present on the SiPM anode to be fine tuned using an internal configuration over a range from 100 to 700 mV, to account for the variations in $V_{BD}$ between SiPMs, sharing a common external bias. Four SiPMs ($8 \times 64$ channels) share a single HV bias channel. SiPMs have been selected and grouped such that the variation between the SiPMs is within the tuneable range of the PACIFIC.

### 6.5.3 PACIFIC analog circuit simulation

A 10 photoelectron SiPM signal with an arrival time of 0 ns at the input has been simulated at 4 V above the breakdown voltage. The signal amplitude over time after the shaper circuit is shown in figure 51 (left) for two separate *pole-zero* shaper settings, pz5 and pz6. The signal has a positive component for 10 ns with a small undershoot afterwards. The pz5 setting is intended to create a larger undershoot compared to pz6. The integrated charge of the shaper signal that falls in one bunch crossing period from 0 to 25 ns is measured in one integrator. The data points in figure 51 (right) are the track-and-hold values (sampled output of the integrator) for separate signals for a range of arrival times. Data points with a negative arrival time have been partially integrated in the previous bunch crossing by the other integrator. The threshold setting relative to the maximum value is indicated by dashed lines in figure 51 (right). This will ensure a relatively flat efficiency for separate signals with the same number of photoelectrons occurring at different times with respect to the clock phase of the integrator. Dark-noise signals will arrive randomly in time and be added on top of any signal pulse and charge spillover across bunch crossing windows.
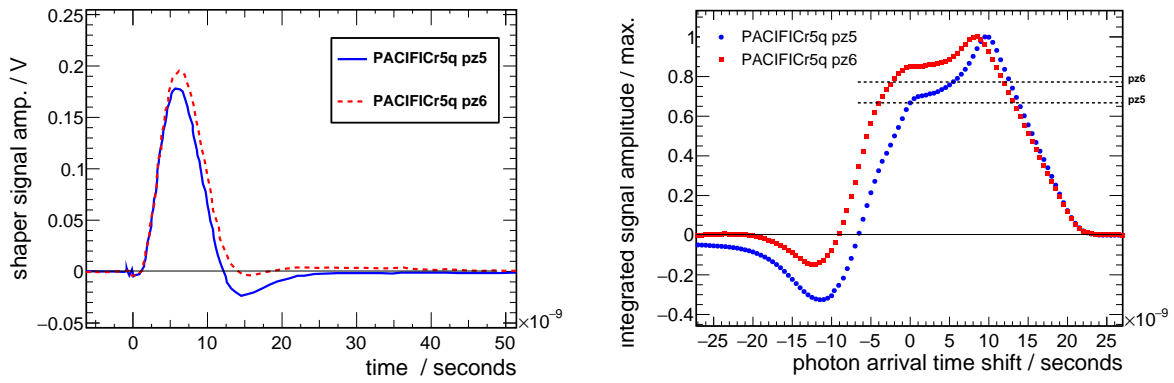
**Figure 51.** Left: the simulated shaper amplitude as a function of time for a single 10 photoelectron signal at 4 V above breakdown. Right: the track-and-hold output values as a function of signal arrival time. The dashed lines indicate the nominal threshold value with respect to the maximum value for two settings (pz5 and pz6, see text).
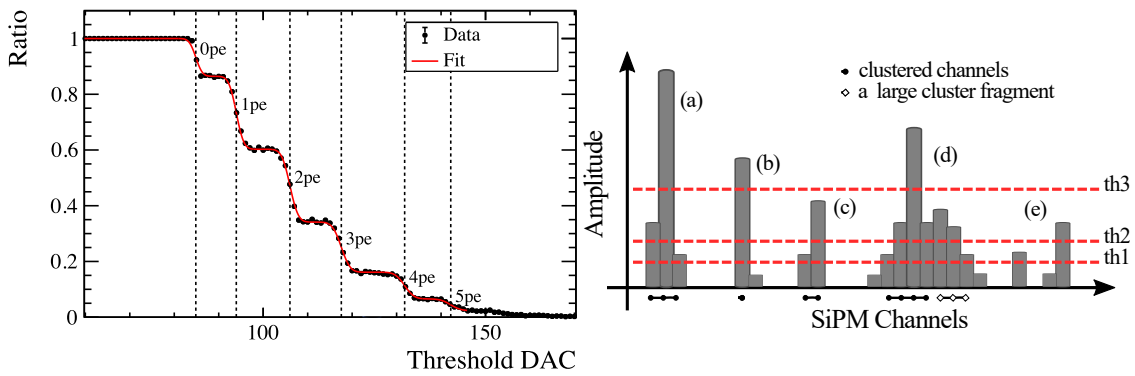


**Figure 52.** Left: an example threshold calibration for one comparator of one channel, showing the ratio (number of data above threshold to the total number of data) as a function of the threshold value. The red curve is the result of a fit. The vertical dashed lines through the steps highlight the discrete photoelectron amplitudes. Reproduced with permission from [95]. Right: a diagram of the clustering algorithm.

### 6.5.4 Threshold calibration

The setting of the thresholds has a large impact on the single hit efficiency and dark-count rate. The calibration of the comparators which are used to set the signal thresholds consists of two parts: (a) determining the ratio of events over threshold for each DAC setting for each of the three comparators of each channel ($3 \times 524$ thousand) and (b) an offline fit to the data to extract the calibration constants. The threshold DAC scan is performed under pulsed illumination provided by the light injection system (LIS). An example of the results for one threshold is shown in figure 52 (left). The fit to the data is based on an analytical description of the Poisson-like SiPM spectrum described in ref. [94], which includes contributions from cross-talk. It not only allows one to determine the threshold DAC values corresponding to the discrete photoelectron amplitudes, nominally 1.5, 2.5, and 4.5 photoelectrons, but also gives access to other important parameters such as the mean light intensity of the LIS. A more detailed discussion of the calibration procedure can be found in ref. [95].

### 6.5.5 Clusterisation FPGA

The PACIFIC 2-bit output data are processed by the Clusterisation board. The FPGA clustering algorithm groups neighbouring channels from the same SiPM array into clusters, calculating an 8-bit cluster position which reduces the data volume from 10.24 Gbit/s per SiPM to less than 4.8 Gbit/s. The combination of PACIFIC thresholds and channel selection in the FPGA suppresses the rate of accidental clusters from dark noise while maintaining a high track hit efficiency.

The clustering algorithm is best explained by the example shown in figure 52 (right), where the three hierarchical thresholds th1, th2 and th3 of the comparator are visualised. A cluster is formed when the sum of the weights of two or more neighbouring channels exceeds the weight of th2, such as clusters (a) and (c) in the figure. A weight of 1, 2, and 6 is attributed to channels which exceed th1, th2, or th3 respectively. Exceptionally, a cluster is also formed when a single channel has an amplitude greater than th3, such as cluster (b) in the same figure. The channels around (e) will not form clusters.

The 8-bit barycentre of a cluster is calculated from a weighted average of all participating channels in the cluster, rounding to a half channel position. This digital half channel precision is enough to provide a position reconstruction resolution better than 100 μm.

A maximum of four channels can be included in a single cluster before it is flagged as *large* in a ninth bit and combined with subsequent cluster fragments, such as (d) in figure 52 (right). An unweighted geometric barycentre is determined for the large clusters. The maximum number of clusters per event that can be sent by a clusterisation FPGA is limited. In the high occupancy region, a flexible data format is used, which sends a maximum of 16 clusters per SiPM per across additional bunch crossings when needed. In the rest of the detector a fixed data format is used which has a limit of 10 clusters, sent synchronously with each bunch crossing.

An FPGA may fail due to the passage of ionising particles over time. Irradiation tests performed at the CHARM facility at CERN indicated that the chosen FPGAs remain re-programmable up to 23 Gy corresponding to approximately $10\,\mathrm{fb^{-1}}$ of integrated luminosity at LHCb. In total, the boards received up to 300 Gy of ionising dose and $3 \times 10^{12}\,1\,\mathrm{MeV}n_{eq}/\mathrm{cm^2}$. The observed speed degradation indicates that the loss will be lower than 5% during the detector lifetime and should not affect the operation of the detector. Increased power consumption was not observed to any significant level. During the irradiation tests three FPGAs (out of 26) ceased to respond and had to be power-cycled to make it function again. No FPGA was permanently damaged.

### 6.5.6 Master boards

The cluster data from the FPGA is serialised by the GBTx ASICs [13] on the Master board and shipped over optical fibres. The FE architecture is such that the cluster data of each SiPM are sent over a single fibre to the TELL40 DAQ. There are four VTTx (8 links) on each Master board while the boards are controlled by using the ECS through the VTRx connector [96].

A housekeeping FPGA on the Master board provides slow and fast control to the light-injection system, the voltage-level shifting necessary to drive the monitor LEDs connected to the GBTx status outputs, and the power-up reset signals to the data GBTx chips.

The Master boards are powered by an external 8 V power supply. The radiation tolerant DC-DC converter (FEASTMP) are used to power the various other components on the FE boards. Each FE box contains 2 Master, 8 PACIFIC and 8 Clusterisation boards. It requires approximately 200 W of power and is cooled by circulated chilled water.

### 6.5.7 Light injection system

For calibration and commissioning purposes a light injection system (LIS) is implemented at the ends of the SciFi modules, as it can be seen in figures 45 and 53. The system consists of an external GBLD-based light driver and a scratched optical fibre which injects light into the scintillating fibres. Fast and slow control signals, as well as power, are transferred through a PCB flex cable connected to the Master board. The pulsed injection is controlled by two GBLD laser driver chips on the LIS mezzanine via I2C such that a distribution from one to five photoelectrons is observed by each SiPM channel. This range was chosen such that the three DAC threshold values of the PACIFIC can be set to correspond to the desired 1.5, 2.5, and 4.5 photoelectrons.

### 6.6 Mechanics and alignment

The SciFi Tracker C-Frames are hung from the rails of the former LHCb Outer Tracker bridge, downstream of the magnet. The C-Frames are made from extruded aluminium profiles. The majority of the service cables and cooling pipes are distributed along these profiles. Additional aluminium covers are fixed to the outside of the profiles to provide electromagnetic shielding and additional stiffness to the profiles. The total mechanical structure when completely assembled with cables, pipes, detector modules and electronics weighs slightly less than 1.5 tonnes.

A system of adjustment screws on the table, bridge, and carriages at the top and bottom of each C-Frame allow for adjustment in three spatial dimensions plus rotations with a precision of about 100 µm. A threaded drive mechanism which is fixed at the bottom of the C-Frame allows for a controlled movement to the final run-position when the C-Frames are closing around the beryllium beam pipe.

#### 6.6.1 Survey

Each SciFi module has four laser tracker survey points, two at each end. There are four more on the beams of the C-Frames which are used during installation and alignment of the detector. The location of these targets are known to an accuracy better than 0.2 mm with respect to their nominal design positions. Additionally, the curvature of all modules has been measured by photogrammetry with reflective targets on the surface of all the modules after installation on the C-Frames. The data from each half of a module was fit to a single plane, where the standard deviation from the plane is of the order of 0.1 mm, approximately the resolution of the measurement. A single plane fit to the entire C-Frame layer has maximum deviations of 1.5 mm, typically localised at the edges or corners of the plane [97].

#### 6.6.2 Real-time 3D position monitoring system

An online 3D metrology system has been developed to permanently monitor the evolution of the position of one detection plane in each of the three SciFi Tracker stations while they are exposed to the magnetic field and other slowly varying experimental conditions. The system relies on triangulation measurements by 24 BCAM cameras [98] of passive reflective targets placed on the detector surface. The camera positions (orientations) are known at the level of a few µm ( µrad). A sequential image acquisition cycle of all cameras provides one determination of the detector position approximately every minute. The intrinsic accuracy of relative movement measurements is better than 100 µm, and averaging over cycles, for up to one hour, can lead to an improvement of the precision to better than 20 µm. These measurement are used to study the geometric evolution of the SciFi Tracker detector, and to validate the alignment constants obtained from the offline tracking alignment algorithm.

## 6.7 Sensor cooling

The SiPMs are coupled to the SciFi modules inside a cold-box. The details can be seen in figure 53. There is a gap of 0.48 mm between each of the 16 neighbouring SiPM sensors on a module. They are glued to a 3D-printed titanium alloy cooling bar and pressed against the fibre ends with springs. The whole assembly is housed in a multipart 3D-printed polyamide shell containing an expanded polyurethane foam layer. A 0.12 mm thick tin-plated copper foil is glued to the shell to improve thermal uniformity on the outer surface. The cooling liquid is supplied and returned via vacuum-insulated stainless steel pipes.
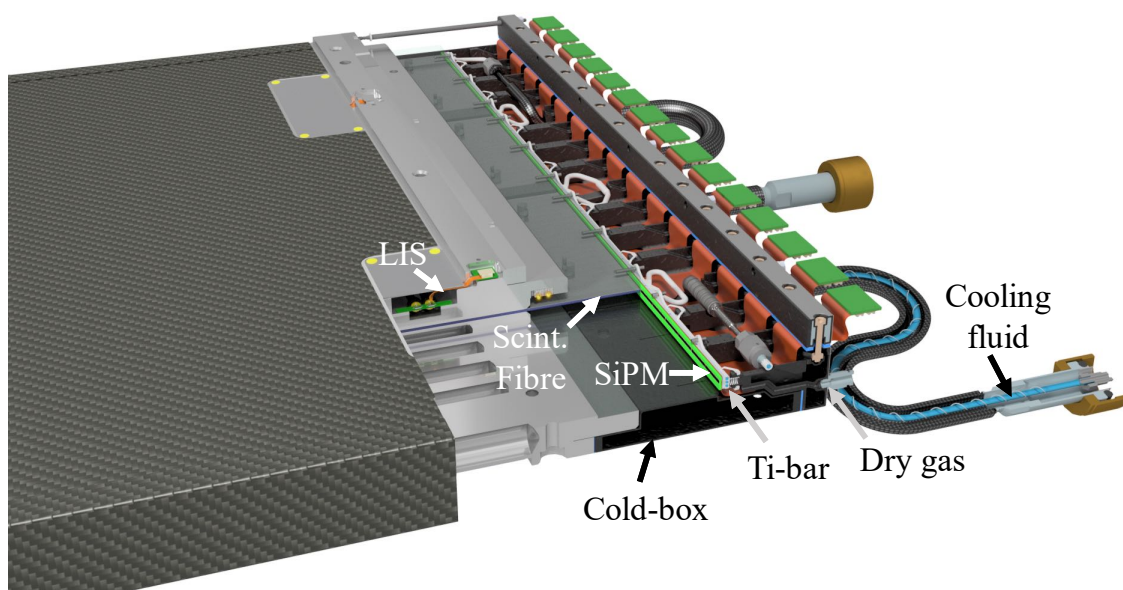


**Figure 53.** A cutaway view of the cold-box fixed to the fibre module.

The SiPM cooling circuit uses a single phase thermal transfer fluid.[47] The delivered coolant temperature can be adjusted between +30 and −50 °C, depending on the desired cooling performance needed. The fluid is circulated through vacuum insulated lines to the detector by a dedicated plant located in the shielded underground area of LHC Point 8. To remove any risk of frost and condensation building up on the cold-bar or SiPMs, every cold box is individually supplied with dry air at a dew point of −70 °C. Thin heater wires are wrapped around the cold-box of each module, as well as transfer bellows on the modules to maintain a stable outer temperature above the cavern dew point. The dry air flow rate out of each box is monitored individually. Four power supplies[48] are needed for the detector which provide up to 50 W of heating power per cold-box.

## 6.8 ECS and DAQ

The SciFi Tracker control system is implemented in the LHCb ECS platform described in section 10.4. The ECS mainly includes the controls of the high-voltage power supplies for the SiPMs, the low-voltage power supplies (electronics and heating wires), the FE and BE electronics (including DAQ).

---

[47]Both fluoroketone and $C_6F_{14}$ are possible.
[48]Wiener™ MARATON.

Additionally, the ECS provides monitoring of the voltages, temperatures and heating powers, of the water cooling system, of the SiPM cooling system, including the integrated vacuum system, of the dry gas system for controlling humidity inside the cold-box, and of the BCAM position monitoring system.

The SciFi Tracker DAQ uses multiple data formats. Two formats are used for clustered data in order to cope with the varying occupancy across the detector, as described earlier. In addition to the two basic data formats and some TFC related outputs, a special data format provides a nonzero-suppressed data mode which allows for the output of the raw 2-bit data along with the clustered data. This requires sending only a fraction of the channels per each bunch crossing, due to the larger amount of information. Detailed documentation can be found in refs. [99] and [100].

## 6.9 Simulation and reconstruction software

To study the tracking performance a detailed SciFi Tracker simulation was performed. It consists of three main parts.

First, the energy deposition for each particle that traverses a fibre mat is computed in the Gauss application [101]. In Boole, where the signal created from the energy deposited is simulated and digitised, the number of photons produced in each channel at the track hit location is calculated assuming 8000 photons produced per MeV of energy deposition. The survival probability for these photons to reach the SiPM is dependent on the properties of the mirrors and the total integrated ionising dose along the fibre. The fibre properties can be adjusted to take into account ageing and exposure to radiation.

A second standalone Geant4-based simulation was developed to propagate optical photons in a single irradiated fibre in order to produce the survival probability map, shown in figure 54, that can be used in the Boole simulation of the detector. The total photon signal is divided amongst neighbouring channels based on the crossing angle of the track and a slight smearing to account for the cluster widths observed in test beam data. The photon propagation simulation and the full detector simulation have been tuned to data obtained from experimental measurements of irradiated fibres and test beam data. A detailed description is given in ref. [102].
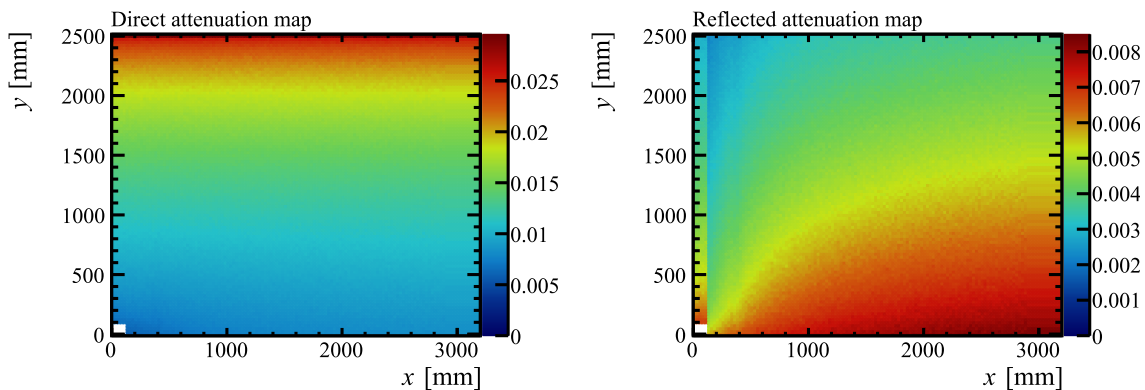


**Figure 54.** The survival probability (attenuation) map of direct and reflected photons in the SciFi Tracker after $50\,\mathrm{fb}^{-1}$ from simulation. Reproduced with permission from [102].

In the third part, the SiPM signal and PACIFIC chip digitisation are simulated. A signal in each channel is obtained based on the quantum efficiency and other properties of the SiPM, the thresholds of the PACIFIC comparators, and the electronics response function of the ASIC, as it can be seen

in figure 51 (right). The dark-noise avalanches of the SiPM are also (optionally) simulated at this stage and added to the analog signal of each channel before digitisation. The simulated digitisation of the thresholds outputs are then passed to the clusterisation output and the position of a cluster can be calculated and encoded. At this stage the output of the simulation corresponds to what is obtained on detector electronics output. The simulations are used for the development of the analysis and monitoring chain, as well as testing various scenarios such as the impact of threshold settings on the production of dark-noise clusters, single hit efficiency, and tracking performance.

### 6.9.1 Spillover clusters

A spillover signal in a given bunch crossing is a signal due to a particle associated to an interaction from a previous or following bunch crossing. There are two main sources of spillover signals. The first is associated with real hits from particles from previous or following bunch crossings that end up inside the current integration window due to flight times and timing offsets. The second is associated with the extended shape of the analog electronics pulse, which will result in a (positive or negative) charge contribution to the PACIFIC integrators in multiple bunch crossings. To mitigate the first type

**Table 8.** The average number of clusters per event occurring in the current bunch crossing for different PACIFICr5q models based on a sample of $B_s \rightarrow \phi\phi$ events generated at the given delay from the current crossing.

|                  | −50 ns | −25 ns | 0 ns | +25 ns |
|------------------|--------|--------|------|--------|
| PACIFICr5q pz5   | 76     | 350    | 4082 | 34     |
| PACIFICr5q pz6   | 91     | 503    | 4021 | 25     |

of spillover, an undershoot at the end of the pulse was added and tuned to reduce the effect arising from −50 and −25 ns crossings. The average number of clusters observed in the current (0 ns) bunch crossing from events generated in the bunch crossings (−50 to +25 ns) are displayed in table 8 for two different models of the PACIFICr5q settings (see section 6.5.3).

### 6.9.2 Dark-noise cluster rates

Initial estimates indicate that the tracking algorithms performance is degraded if the rate of dark clusters per SiPM increases above 2 MHz, or approximately 200 dark-noise clusters across the entire tracker per bunch crossing. However, this rate is dependent upon the single channel DCR and the thresholds set in the PACIFIC. Ideally, the thresholds are set as low as possible to maximise the efficiency, while accepting a tolerable amount of dark-noise clusters that can be removed in the track finding algorithms. A comparison of the number of thermal-noise clusters per bunch crossing for two different electronics response configurations is shown in figure 55. The data points are generated from a detailed model of the DCR, PACIFIC response function, SiPM cross-talk, and the clustering algorithm. A simple power law function, overlaid in the figure, describes well the data.

### 6.10 Test beam results

A slice test of two SciFi Tracker modules coupled to a complete set of nearly final FE readout electronics was performed in 2018, with a preliminary standalone version of the LHCb 40 MHz PCIe40 readout [103]. The system was tested for its single particle hit reconstruction efficiency and position reconstruction resolution using a beam from the Super Proton Synchrotron (SPS) at CERN's
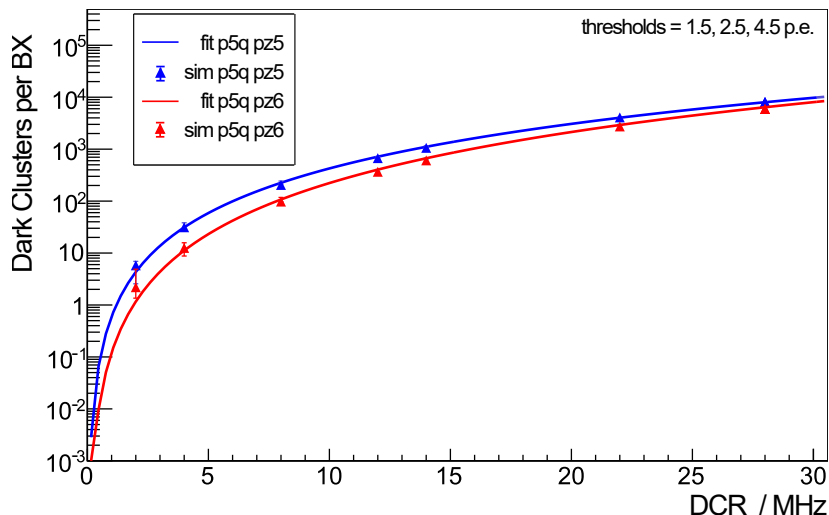
**Figure 55.** Simulated number of thermal-noise clusters per 25 ns clock cycle as a function of the DCR. The curves stem from a power law fit to the data points. The comparison is made for two different models of the PACIFIC settings (in blue the pz5 settings, in red the pz6). The three threshold values used are 1.5, 2.5 and 4.5 photoelectrons.

North Area consisting of 180 GeV/$c$ pions, protons and/or muons. The SiPMs were not irradiated and kept at room temperature with chilled water.

To ensure that the reconstructed particle tracks are of good quality and to provide a precise hit location, the TimePix3 telescope in the H8A area was operated synchronously with the PCIe40 readout. The telescope provides a fine time stamp that divides the 25 ns long clock cycles into 96 intervals of about 0.25 ns.[49] This fine time stamp is required as the arrival of the particles from the SPS are not synchronised in time with the DAQ system, unlike the particles that will be generated at the LHC.

For every trigger and high-quality track from the telescope a search is made in the SciFi Tracker data for a corresponding signal cluster (one or more neighbouring channels with a signal). The relative position of the found cluster is compared to the telescope track position in the module to determine the position resolution of the fibre modules, as well as the single hit efficiency.

### 6.10.1 Single-hit efficiency

A log-normal distribution describes well the measured inefficiency distribution observed in the test beam as can be seen in figure 56 (left). From these data a mean hit efficiency for the sensitive regions of the SciFi Tracker (excluding gaps) is estimated to be $0.993 \pm 0.002$ (quoting the standard deviation of the log-normal as an estimator for the expected spread).

### 6.10.2 Hit position resolution

The hit position residual is defined as the difference between the cluster position and the position of the telescope track extrapolated to the fibre plane, after having applied alignment corrections using Millepede [104]. An example hit position residual distribution, merging the data from three beam spot positions, at the centre and outer edges of one 32.6 mm wide SiPM arrays, is shown in figure 56 (right). As the telescope resolution contributes insignificantly, the width of the residual

---

[49]Due to the multiple PLLs used to generate this fine time binning, the intervals are not all equally long.
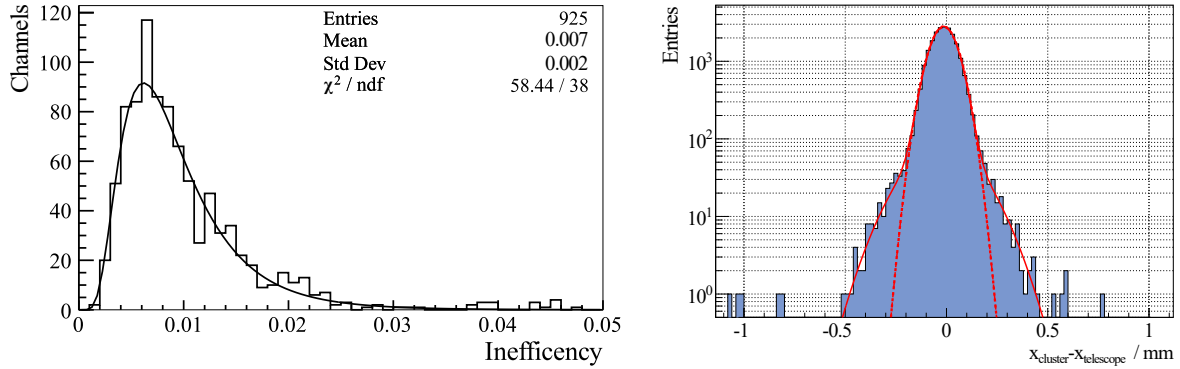
**Figure 56.** Left: the single hit inefficiency of the five most efficient fine time bins of all channels measured at several positions across the module. The distribution is fit to a log-normal distribution. The mean and standard deviation of the log-normal distribution is also shown. Right: an example hit position residual distribution fitted with a single (dashed curve) and double (solid curve) Gaussian function. Gap regions have been excluded.

distribution can be taken as the resolution of the SciFi Tracker module. The distributions of the hit position residuals measured across the modules in the test beam campaign indicate a single hit position resolution of 64±16 μm for perpendicular tracks.

# 7 RICH

Charged hadron discrimination, namely the separation between pions, kaons and protons, is a crucial aspect of the LHCb physics programme. Hadron PID is provided in LHCb by the RICH system in the 2.6–100 GeV/$c$ momentum range, and plays a central role in the measurements performed in LHCb with Run 1-2 data [105, 106]. The RICH system allows to: distinguish between final states of otherwise identical topologies, e.g. $B^0_{(s)} \rightarrow \pi^+\pi^-, K^+\pi^-, K^+K^-$ decay modes; heavily reduce the combinatorial background in decay modes involving hadrons in the final state, such as $B^0_s \rightarrow \phi\phi$, where $\phi \rightarrow K^+K^-$, that would be prohibitively large without PID requirements; perform the flavour tagging of a $B^0_{(s)}$ meson at the production vertex, relying on charged kaon identification from the $b \rightarrow c \rightarrow s$ decay chain. The information provided by the RICH system is also used to suppress the combinatorial background at the HLT2 level (see section 10.2).

The overall layout and concept of the RICH system remains unchanged with respect to Run 1-2 LHCb [1], although critical modifications were needed to allow the system to operate at the higher design luminosity while maintaining a performance comparable to that of Run 1 and Run 2 [107, 108]. It consists of two detectors, RICH1 and RICH2, as shown in figure 57. RICH1 covers an angular acceptance from 25 to 300 mrad in the magnet bending plane and from 25 to 250 mrad in the vertical direction. The photon detector planes are located above and below the beam pipe, where the residual magnetic field is minimal. RICH2 is located downstream the dipole magnet, covering an angular acceptance from 15 to 120 mrad in the magnet bending plane and 15 to 100 mrad in the vertical direction. The photon detector planes are located on the sides. In both detectors the Cherenkov photons produced inside fluorocarbon gaseous radiators are reflected outside the LHCb acceptance by means of a system of spherical and planar mirrors, focusing the ring images on the photon detector planes.

RICH1 is located upstream of the dipole magnet and employs a $C_4F_{10}$ gas radiator with a refractive index $n = 1.0014$ for Cherenkov radiation of wavelength $\lambda = 400$ nm at standard temperature and
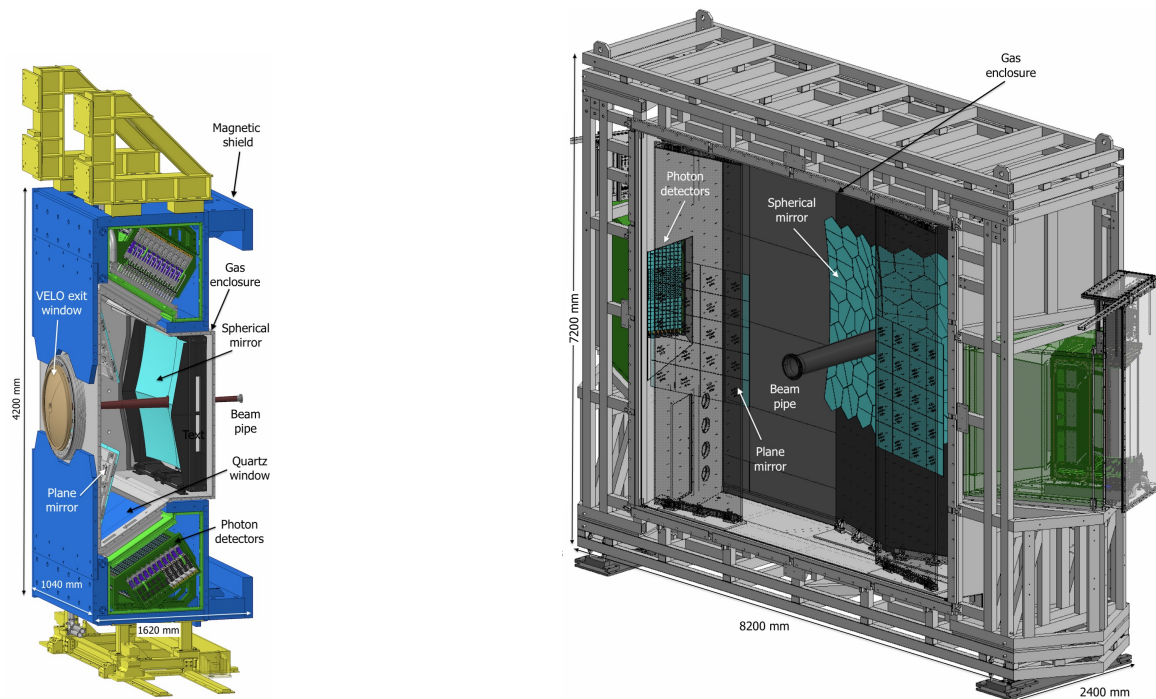
**Figure 57.** Schematic view of the (left) RICH1 and (right) RICH2 detectors. Reproduced from [109]. © 2022 IOP Publishing Ltd and Sissa Medialab. All rights reserved.

pressure (STP), allowing to provide PID in the momentum range between 2.6 and 60 GeV/$c$. The average path length of particles inside the radiator is approximately 110 cm. RICH2 is designed to provide PID for higher momentum particles, between 15 and 100 GeV/$c$, with a CF$_4$ gas radiator with $n = 1.0005$ for Cherenkov radiation of $\lambda = 400$ nm at STP, and an average track path of 167 cm.

In order to read the detectors out at 40 MHz rate, the full photon detection chain was replaced in both RICH1 and RICH2 detectors, since the former hybrid photon detector (HPD) [110] had embedded FE electronics limited to a 1 MHz output rate. The HPDs have been replaced with multi-anode photomultiplier tubes (MaPMTs) equipped with new FE electronics. The upgraded photon detection modules are described in section 7.1.

One of the key parameters driving the performance of the RICH system is the efficiency of the pattern recognition algorithm, optimal for detector occupancies not exceeding 30% as determined from experience in Run 1 and Run 2 operations.[50] With the five-fold increase in the instantaneous luminosity, a redesign of the RICH1 optics was necessary to reduce the peak occupancy, as described in section 7.2. The optical system and mechanical envelope of RICH2 was left unchanged, but redesigned support structures to house the new photon detectors were required, as reported in section 7.3.

Monitoring and control systems and the calibration procedure have also been updated, to cope with the changes in the photon detection chain and the readout infrastructure, as described in sections 7.4 and 7.5. The upgraded RICH system has been designed to improve the single photon resolution and to keep the excellent PID performance provided in Run 1 and Run 2 in the more challenging conditions of Run 3. Its expected performance is discussed in section 7.6.

---

[50]The occupancy is defined as the number of fired channels over the total number of channels in a given region at 40 MHz readout rate.

## 7.1 The upgraded photon detection chain

The design of the upgraded photon detection chain has been optimised in order to cope with the highly nonuniform occupancy expected in the RICH system, ranging from about 30% in the central region of RICH1 down to 5% in the peripheral region of RICH2, if the previous photon detectors were kept unchanged. The largest hit rates correspond to Cherenkov photons associated to the large number of tracks produced at high pseudorapidity. Given the observed occupancy distribution, the detector geometry and channel granularity have been optimised taking into account existing overall mechanical constraints and the number of readout channels, which have significant impact on the cost. As described in the following, the photon detection planes are subdivided into two regions having different granularity, with the aim of keeping the optimal performance while maximising the cost savings. In addition, in order to ensure stable operations of the upgraded RICH detectors, an evaluation of the photon detection chain performance under high radiation fields has been performed as described in section 7.1.5.

### 7.1.1 Photon detectors

The main parameters driving the choice of photon detectors have been good spatial resolution on a large active area, high detection efficiency in the wavelength range 200–600 nm, and very low background noise to allow single photon detection despite the high occupancy foreseen. The MaPMTs had already been considered during the first construction of the RICH detectors, but were rejected mostly due to the limited fill factor ($\sim 40\%$) of such devices at the time.[51] The latest models available on the market are instead characterised by a fill factor exceeding 80% and could be adopted as the upgraded RICH photon detectors. The selected MaPMT models both consist of a matrix of $8 \times 8$ anodes. RICH1 and the central region of RICH2 are equipped with 1-inch MaPMT modules[52] with a pixel size of $2.88 \times 2.88\,\text{mm}^2$, ideal for the high occupancy areas of the RICH system. The outer region of RICH2 has been equipped with a 2-inch device[52], with a pixel size of $6 \times 6\,\text{mm}^2$. The decision to install detectors with a coarser granularity in the peripheral regions of RICH2 allowed a significant reduction in the number of MaPMT units and readout channels, while having a negligible impact on the overall RICH performance as demonstrated by simulation studies. The MaPMTs installed in the upgraded RICH detectors, together with the schematic view of their internal structure, are shown in figure 58.

A total of 1888 (768) 1-inch MaPMTs are installed in RICH1 (RICH2) and 384 2-inch MaPMTs are installed in RICH2. Over 3500 units, including spares, have been purchased by the RICH collaboration and quality-assured to verify the requested technical specifications, among which are a gain larger than $10^6$ and a dark-count rate less than $< 2.5\,\text{kHz/cm}^2$. A full set of quality assurance (QA) tests was implemented to qualify the whole MaPMT production. Two of the typical QA parameter scans, the signal amplitude as a function of the HV and the quantum efficiency (QE) as a function of the wavelength, are shown in figure 59.

---

[51]The fill factor is the fraction of active area with respect to the total detector area.

[52]Hamamatsu$^{\text{TM}}$ R13742 (1-inch) and R13743 (2-inch). The R13742 and R13743 models are custom variants of the commercial models R11265 and R12699, respectively. The difference between the custom and commercial units stands on the requested technical specifications described in the text.

**Figure 58.** Left: the MaPMTs selected for the upgraded RICH detectors with the 2-inch model on the left and the 1-inch model on the right. Right: scheme of the internal structure of the MaPMT. Reproduced from [122]. CC BY 4.0.
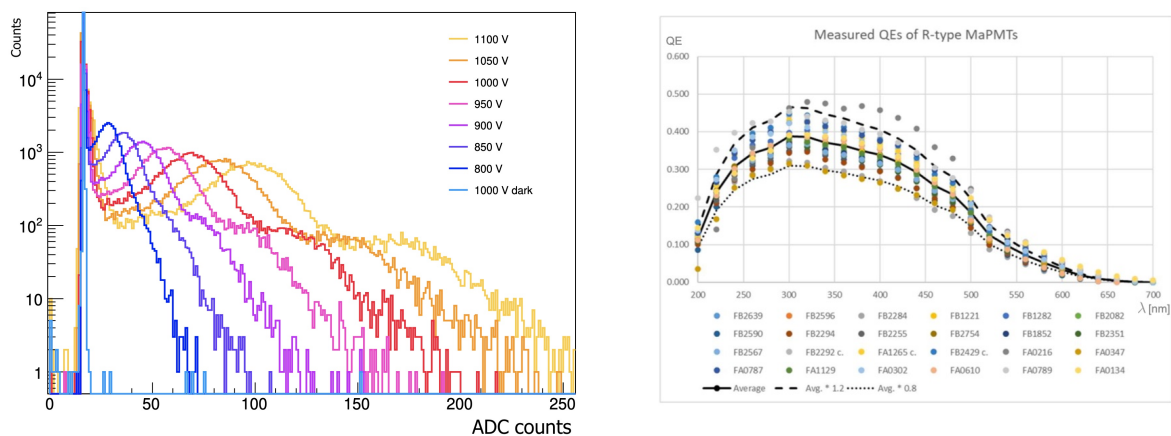


**Figure 59.** Left: typical signal amplitude spectra for a pixel as a function of the HV value. Reprinted from [111], Copyright (2023), with permission from Elsevier. Right: QE curves for a batch of 1-inch MaPMTs from the production: the ultra bi-alkali photocathode allows to reach excellent QE values.

### 7.1.2 Front-end electronics and elementary cell

The average hit rate in the high occupancy regions of RICH1 can exceed $10^7$ hits/s per pixel on average. Furthermore, the estimated total integrated dose over the detector lifetime, in the regions closer to the beam pipe is estimated to be about 2 kGy for RICH1 [112]. Radiation-hard fast readout electronics is therefore needed, with low power consumption to minimise heating. This motivated the design of a custom 8-channel front-end ASIC named CLARO [113]. CLARO, shown in figure 60, is designed in 350 nm CMOS Austria Micro Systems (AMS) technology, with the exception of configuration registers as will be discussed in section 7.1.5.

Each CLARO channel is composed of an analogue transimpedance amplifier followed by a discriminator. Converted and discriminated input current signals trigger asynchronous digital pulses at the output. The output signals have a voltage swing of 2.5 V, and a variable length allowing time-over-threshold measurements.

A 128-bit register allows CLARO single channel configuration. In particular, to allow for channel-by-channel gain differences, input signals can be attenuated by factors 1, 1/2, 1/4, and 1/8.
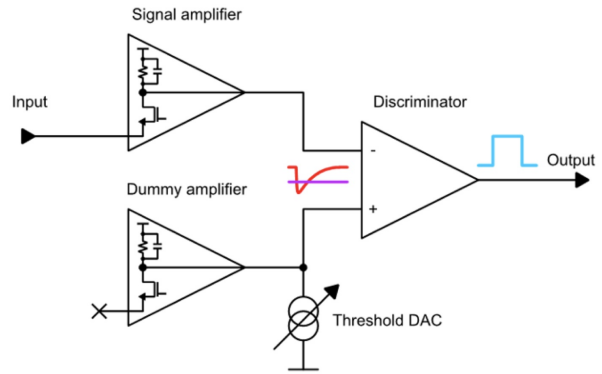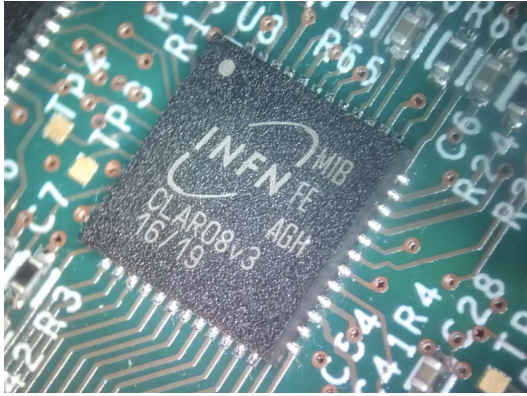
**Figure 60.** Left: CLARO ASIC with its packaging. Reproduced from [113]. © 2017 IOP Publishing Ltd and Sissa Medialab. All rights reserved. Right: block schematic of a CLARO channel. The purpose of the dummy amplifier is to give each channel a differential structure, improving the power supply rejection ratio and allowing DC-coupled input to the discriminator. Reproduced from [114]. © 2022 IOP Publishing Ltd and Sissa Medialab. All rights reserved.
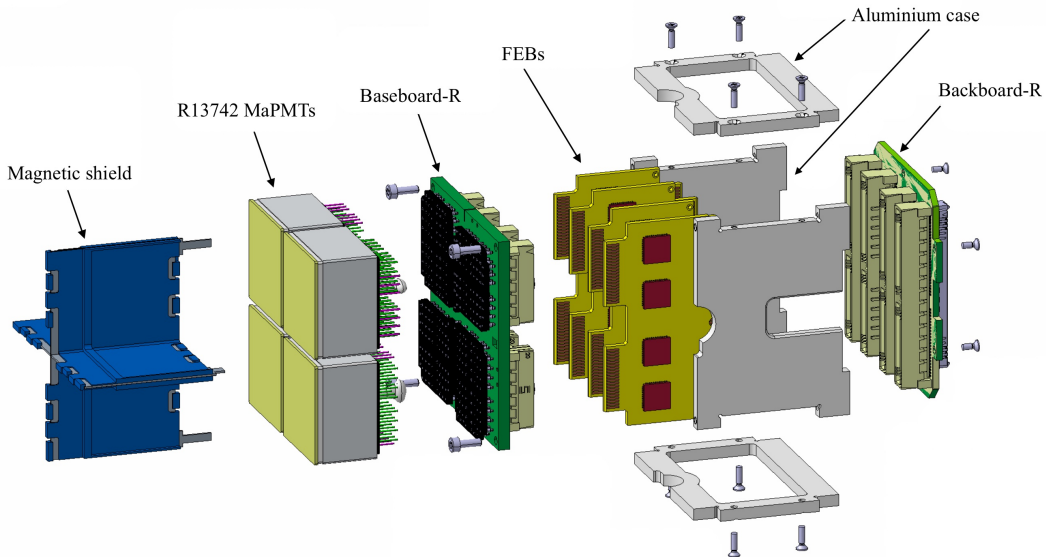


**Figure 61.** Exploded view of the R-type elementary cell. Reproduced from [122]. CC BY 4.0.

Individual thresholds can also be set, with the possibility to cancel the discriminator offset. Thresholds are calibrated with test signals injected at the input through a dedicated test capacitor. A more detailed description of the CLARO design and its functionalities can be found in ref. [113].

The CLARO number of channels matches the $8 \times 8$ pixel modularity of the MaPMTs and allows placing the ASIC as close as possible to the MaPMT anodes, minimising the parasitic capacitance at the input and the susceptibility to electromagnetic interference noise.

The readout system was arranged in compact units named elementary cells (ECs). Two types of ECs, adapted to the different MaPMT models, are used: the R-type elementary cell (EC-R) and the H-type elementary cell (EC-H). A view of the EC-R is shown in figure 61. It reads out four 1-inch MaPMTs, for a total of 256 pixels in approximately $2 \times 2$ square inches. The MaPMTs are
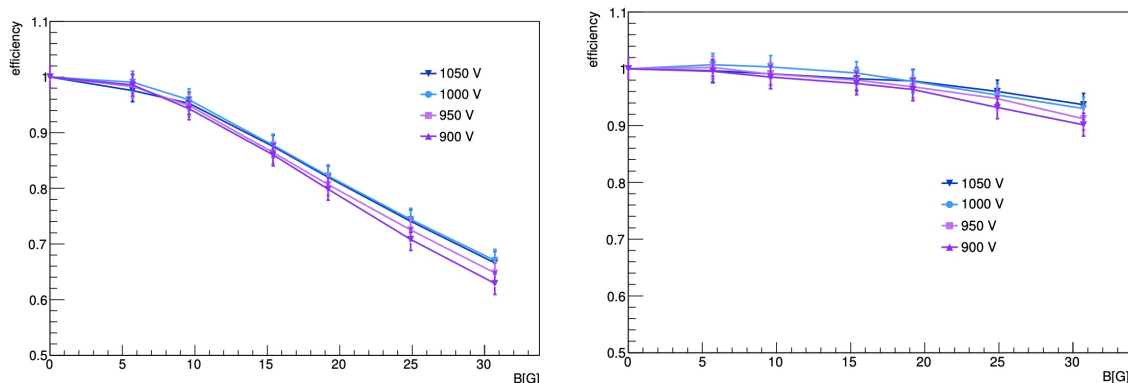
**Figure 62.** Counting efficiency as a function of the longitudinal magnetic field for an edge pixel, at different values of HV, for an EC-R (left) without and (right) with the magnetic shield.

plugged into a baseboard, which hosts four $3\,M\Omega$ resistive dividers in parallel, to bias the dynodes of each MaPMT. The last two dynodes of each chain can be powered by dedicated supply lines in high-occupancy regions, where the drawn current is higher and can induce nonlinear effects in MaPMT gain. A magnetic shield is placed in front of the MaPMTs in the RICH1 EC-Rs where, even inside the magnetic shield, the stray magnetic field from the LHCb magnet is up to about 2 mT. The shield is cross-shaped and made of a $500\,\mu m$ thick mu-metal. It deflects the field lines, attenuating the magnetic field that reaches the MaPMT by a factor of approximately 20, down to a value where its effect on the performance of the MaPMTs becomes negligible, as shown in figure 62.

The baseboard propagates the anode signals to four front-end boards (FEBs), hosting eight CLARO ASICs each (four on each face of the board). The FEBs are in turn connected to a backboard routing the output signals to the photon detector module digital boards (PDMDBs), described in section 7.1.3, through two high-density connectors. The CLARO power supply and control signals are generated on the PDMDBs and are routed through the backboard as well. A 3.0 mm thick and 40.5 mm long aluminium case serves as a mechanical support structure for the electronic components and allows thermal transfer by conduction, with the heat dissipation from the voltage dividers enhanced by copper layers inside the baseboard. Temperature monitoring is also ensured by temperature probes. There are 472 EC-Rs in RICH1 and 192 in the central region of RICH2.

The EC-H, shown in figure 63, reads out a single 2-inch MaPMT. Accordingly, it consists of a single $2.5\,M\Omega$ voltage divider and two FEBs with half the CLARO channels disabled. There are 384 EC-Hs in the peripheral region of RICH2.

### 7.1.3 Photon detector module digital boards

The PDMDB is required to transport the digitised photon detector signals away from the high-radiation region of the detector without introducing dead time and ensuring the interface with the LHCb ECS.

An FPGA-based approach is adopted as a flexible way to capture and format the data and to interface between the different electrical signalling standards of the front-end ASICs and GBT chipset. A comprehensive set of measurements at a number of irradiation facilities, reported in section 7.1.5,
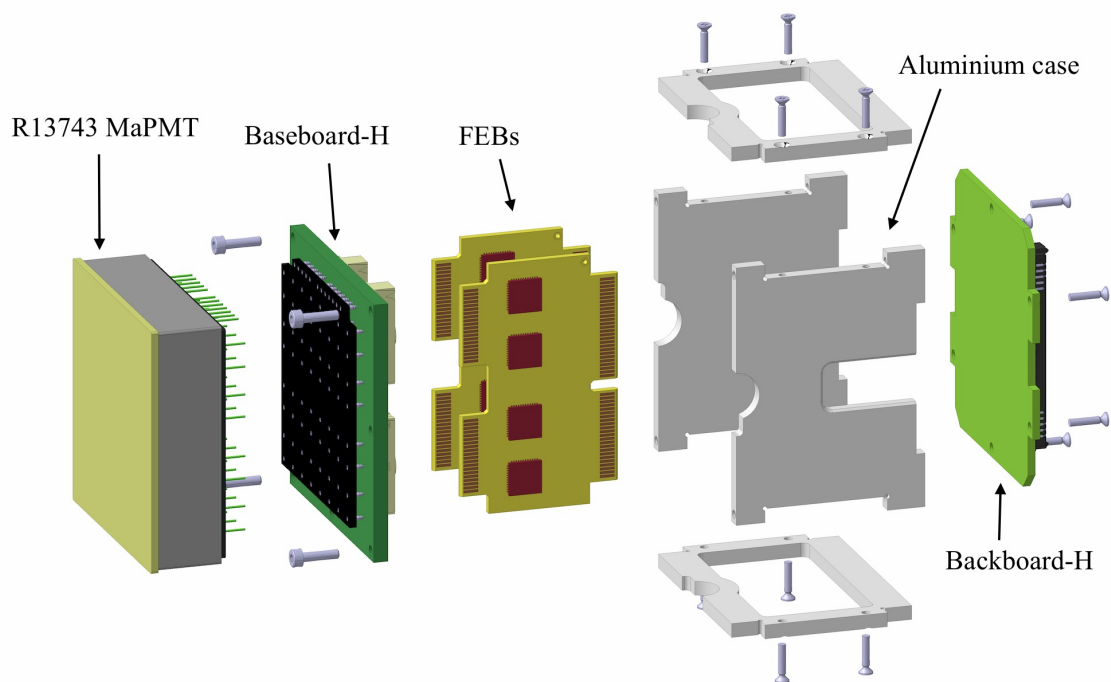
**Figure 63.** Schematic view of the H-type elementary cell. Reproduced from [122]. CC BY 4.0.

has demonstrated that the chosen FPGA[53] is sufficiently tolerant to the effects of radiation in the RICH environment, provided certain mitigating design features are incorporated. Nevertheless, a modular design, with the radiation-hard components on pluggable modules, allows these parts to be reused in case it becomes necessary in the future to replace the FPGAs.

Two variants of the PDMDB are used, corresponding to the different granularity of the photon detector planes. A pair of back-to-back PDMDB-Rs is coupled to a group of four EC-Rs and a single PDMDB-H is coupled to a group of four EC-Hs. The assembly of four ECs and one or two PDMDBs is called a photon detector module (PDM). Each PDMDB hosts one timing and control module (TCM) and up to three data transmission modules (DTMs), implemented as pluggable mezzanine boards, following the concept outlined above.

The TCM is a $3 \times 6\,\mathrm{cm}^2$ module that provides an interface for the fast- and slow-control data exchanged between a PDM and the LHCb ECS. The physical link is implemented using a VTRx and GBTxs operating in bidirectional forward-error-correction mode. The initial configuration of the TCM is programmed into its e-fuses to ensure proper operation at power-on. Configuration protocols provided by the TCM include I2C to program the DTM GBTx, JTAG to program the FPGAs, SPI to configure the CLARO ASICs, ADCs for temperature and voltage monitoring for the PDM, DACs to generate the voltage level for CLARO test pulse generation and general purpose input/outputs (GPIOs) for local resets and digital control.

The DTM is a $3 \times 6\,\mathrm{cm}^2$ plug-in module that provides the high-speed data transmission interface for the PDM. There are three (two) DTMs on each PDMDB-R (PDMDB-H). The physical uplink is
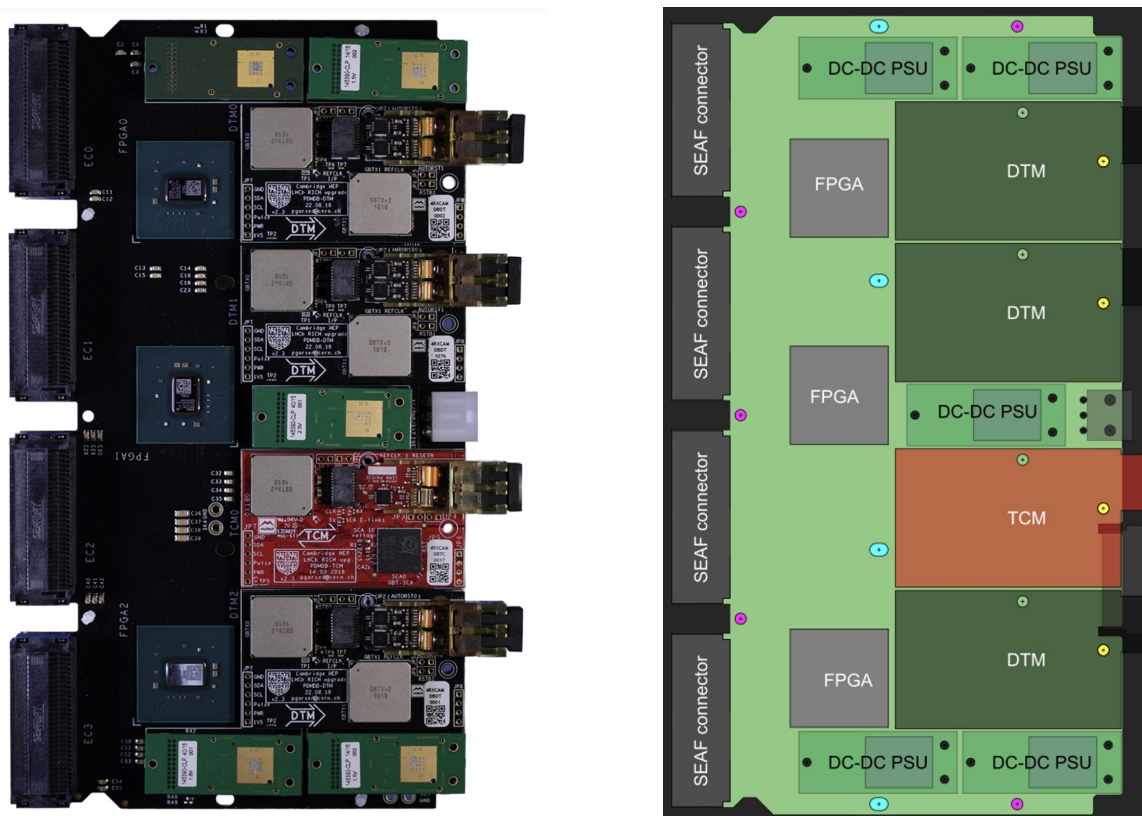
---

[53]Xilinx Kintex-7.

**Figure 64.** Left: picture of a fully populated PDMDB-R board. Right: schematic view of the PDMDB-R board main components. The PDMDB-H differs by having one less FPGA and DTM with respect to the PDMDB-R.

implemented by the VTTx dual optical transmitter with each channel connected to a GBTx ASIC, each operating in wide-bus transmission mode. The GBTxs are configured through their I2C configuration port. The two GBTxs and the FPGA are connected to a dedicated TCM I2C bus master.

The PDMDB motherboard acts as a bridge for the signals between the ECs and the TCM and DTMs. The board also incorporates local power regulation for the FPGAs as well as for the active components on the ECs, TCM and DTMs using CERN FEASTMP DC-DC converters. The only active components on the motherboard apart from the DC-DC converters are the FPGAs. These receive the 2.5 V LVCMOS digital outputs of the CLARO ASICs. No zero-suppression is applied, thus the FPGAs effectively sample the CLARO data at 40 MHz and transport the sampled data transparently to the up-links with constant latency.

### 7.1.4 Photon detector columns

The ECs and PDMDBs are arranged into two types of PDM: the PDM-R, composed of four EC-Rs and two back-to-back PDMDB-Rs, installed in the whole RICH1 and in the central region of RICH2; the PDM-H, composed of four EC-Hs and one PDMDB-H, installed in the peripheral regions of RICH2.

For both RICH1 and RICH2 detectors, six PDMs are assembled on a T-shaped aluminium structural element, referred hereafter as T-bar, to build one RICH column, including the distribution of services and the cooling circuit. In order to minimise the production of specific mechanical components for RICH1 and RICH2, the T-bar is kept identical between the two detectors. The T-bar
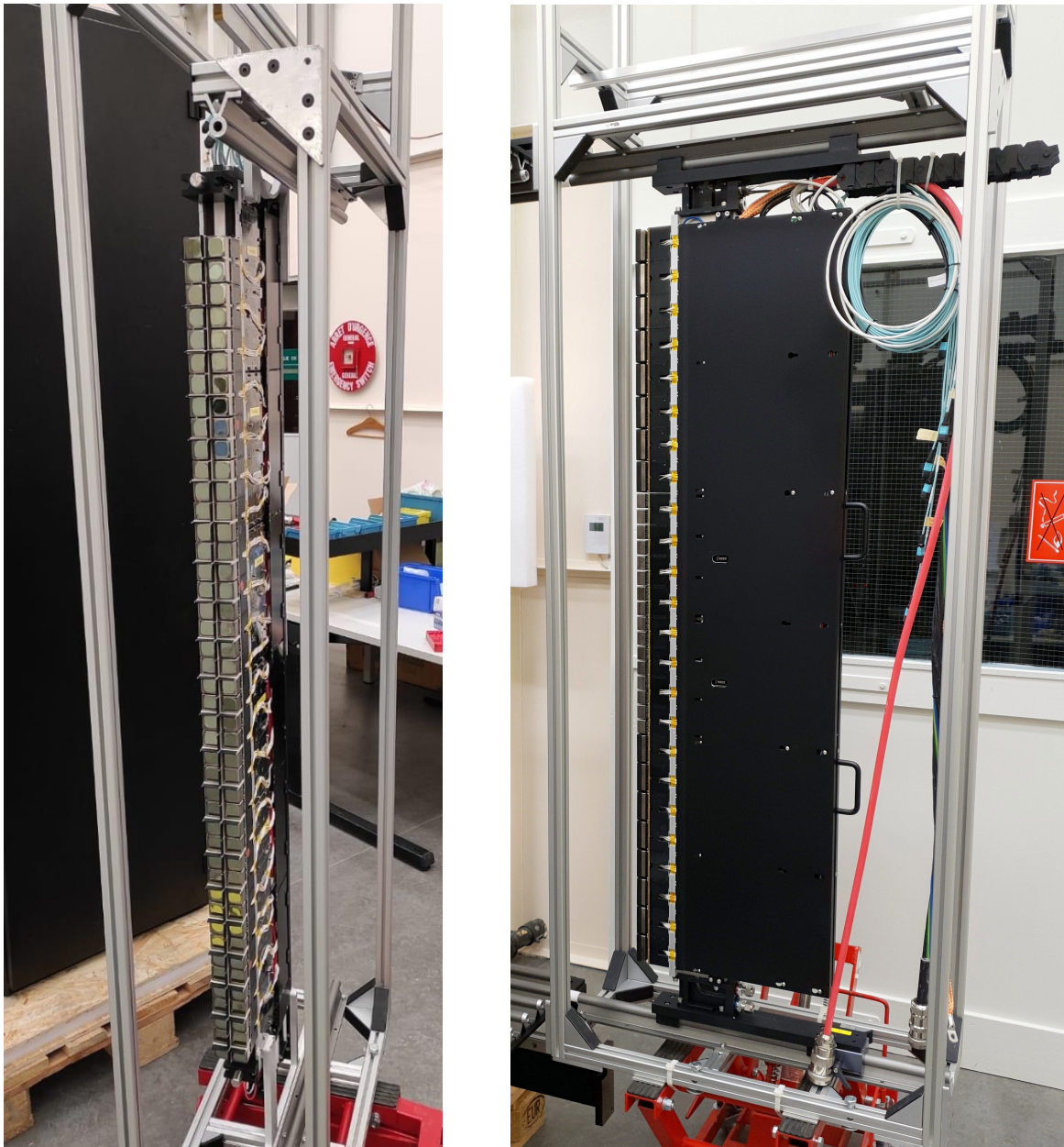
**Figure 65.** Pictures of (left) a completed RICH1 column (front view) and (right) a RICH2 column (side view, ECs on the left), with the photon detector chain and the complete set of services.

provides a precise reference for the positioning of the ECs and PDMDBs, at the level of 0.2 mm. The overall length of a column is approximately 1.6 m with a width of 55 mm and a depth of 40 cm. Fully populated RICH1 and RICH2 columns are displayed in figure 65.

The ECs are fixed at the front of the T-bar base and plugged into the corresponding connectors of PDMDB connectors which are mounted on the T-bar sides. The active components on the PDMDB are a significant source of heat and are therefore cooled by means of an aluminium plate that provides a thermal coupling with the T-bar. The thermal exchange is favoured by using commercial thermal pads placed on the active elements of the PDMDBs.

RICH1 columns contain 22 ECs, instead of the 24 corresponding to six PDMs, since the two ECs at the upper and lower end of each array are not mounted to facilitate the installation, handling and maintenance, while maintaining the complete acceptance. For the same reason, only four PDMs (16 ECs) are installed in the outer column of each RICH1 photon detector plane. RICH2 columns are fully populated with 24 ECs and have four PDM-Hs interleaved by two PDM-Rs.

The cooling of the photon detection chain is achieved by circulating a fluoroketone coolant[54] in two 6 mm diameter ducts deep-drilled into the spine of the T-bar. The cooling keeps the temperature at the MaPMT surface well below 30°C.

The instrumented columns require the distribution of services such as power supply cables, data, TFC and ECS optical fibres, and monitoring devices for the slow controls and DSS. These services run on both sides of the T-bar. The LV supply is provided by two power supply channels[55] for RICH1 columns, while one channel is used for the RICH2 columns. Dedicated distribution boards, that are located at one end of the column, provide the 2.5 V supply by means of DC-DC regulators. Each MaPMT is supplied by a high voltage (900 V) and an intermediate voltage (90 V) to power the last dynode to mitigate possible nonlinear effects on the gain within the MaPMTs. The HV supplies are provided by common floating ground A1538DN CAEN boards. The HV supply has a common floating ground for the complete column. Control signals of the PDMDBs and data are transmitted through long-distance optical fibres with twelve-fibre ribbon with MPO-to-LC connectors[56] used to fan-out the optical links to the individual connectors on the PDMDBs. To monitor the overall temperature, a total of 168 (112) temperature probes are installed on RICH1 (RICH2) columns. These temperature sensors are monitored through the ECS as described in section 7.4.1. For safety purposes and in order to cope with eventual network disruptions, a thermo-switch (normally closed), is mounted directly on the T-bar and will issue a DSS alarm if the temperature exceeds 35°C.

At each end of the column, the T-bar is fixed to a trolley composed of an interface plate and two open cylindrical bearings, made of low friction and electrically insulating polymer, that slide on cylindrical rails. In this way it is possible to easily extract the columns from their operational position for maintenance. The alignment of the columns inside the corresponding rack can be affected by mechanical tolerances and differential thermal dilation, that are compensated by a small degree of freedom of the bearing at the top-end of each column. This clearance is recovered by preloaded washer springs acting between the floating trolley and the T-bar end.

### 7.1.5 Irradiation campaigns and mitigations

According to FLUKA simulations, at the RICH photon detection system location during the whole upgrade phase (corresponding to about 50 fb$^{-1}$ of integrated luminosity) a total ionising dose of 200 kRad and a fluence of $3 \times 10^{12}$ 1 MeV $n_{eq}$/cm$^2$ and $1 \times 10^{12}$ high-energy hadrons (HEH) per square cm are expected, where the estimations include a safety factor of two. Several irradiation campaigns have been carried out in order to assess the impact of radiation on photon detection chain components, in particular on MaPMTs, CLARO ASICs and PDMDB FPGAs, taking into account additional safety factors with respect to the expected dose. Irradiation tests on other components, such as passive elements, cables and mechanical components, were also performed, showing no radiation-induced degradation.

---

[54]Novec 649$^{\text{TM}}$.

[55]Wiener$^{\text{TM}}$ MARATON.

[56]MPO, multifibre push on connector and LC, Lucent connector are standard connectors for optical fibres.

Radiation effects on MaPMTs elements have been carefully studied. Radiation-induced effects on photocathode sensitivity and secondary emission ratios were found to be negligible. Radiation damage on optical entrance windows was also investigated. Multiple samples of MaPMT windows made of borosilicate and UV-transmitting glass were irradiated at different particle fluences. UV glass windows were found to suffer substantially smaller degradation than borosilicate ones and were therefore chosen for the installed MaPMTs.

Radiation hardness tests of a CLARO prototype have been performed with neutrons, X-rays and protons, as described in refs. [115, 116]. Further tests with ion, proton and mixed-field high-energy beams where performed on the first two full versions of CLARO, which showed soft single event latchup (SEL) events at values of linear energy transfer (LET) of about $20\,MeV/mg/cm^2$. As a result of these tests a third version of CLARO was produced, where configuration registers were resynthesised using cells radiation-hard by design [117, 118], which exhibited a higher LET threshold for SEE with respect to AMS standard cells. Tests on this CLARO version [112] have confirmed that the design modification was effective, with threshold for SEU and SEL having increased by a factor of three with respect to the previous versions.

FPGAs are sensitive to SEU which may flip configuration bits disrupting the correct operation of controlled devices and of the FPGA itself. Therefore, the PDMDB FPGAs have been tested with various species of ionising particles. In particular, to emulate as closely as possible the LHC environment, the FPGAs have been tested under a mixed neutron and HEH irradiation fields at the CHARM facility [119] at CERN. The SEU cross-section within the configuration memory has been measured by counting single- and multi-bit errors arising while emulating a fixed pattern through the FPGA logic. The estimated SEU cross-section over the full 19 Mbit FPGA configuration memory (CRAM) has been determined to be $(1.02 \pm 0.37) \times 10^{-7}$ $cm^2$/device [120]. No SEL events were observed in the tests at CHARM while a SEL threshold of approximately $15\,MeV\ cm^2/mg$ has been found when irradiating the FPGAs with Kr, Ni and Ar ions. This threshold is considered acceptable but risk mitigation actions have been put in place as described below. In order to further decrease the risk of logic failures, the PDMDB firmware was designed in order to keep minimum complexity reducing the CRAM usage to about 30 kbit. As much data processing as possible was shifted to the BE DAQ boards, resulting in a worst-case upper limit of approximately 28 logic failures per hour.[57] The PDMDB output data are presented to the TELL40 processing logic after a constant delay that simply adds to the optical fibre propagation delay. As a result, any synchronous processing required by the data transmission protocol can be safely performed in the TELL40, therefore saving substantial logic resources in the PDMDB FPGAs and reducing significantly the probability of radiation-induced upsets. Furthermore, the most critical parts of the FPGA logic are protected using an extended triple modular redundancy technique that also allows selective partial reconfiguration of the FPGA without disrupting the logic operation. Finally, a fast recovery procedure is implemented in the BE electronics and, for redundancy, in the slow control system.

## 7.2 RICH1 optical and mechanical systems

The upgraded RICH1 detector, located upstream of the LHCb dipole magnet between the VELO and the UT, underwent major design changes and a subsequent rebuilding. Nevertheless, the fluorocarbon $C_4F_{10}$ gas radiator was retained and the angular acceptance was kept unchanged.

---

[57]The number is relative to all the FPGAs used in RICH1 and RICH2 when considering the worst case scenario irradiation level of RICH1 everywhere, a safety factor of four, and averaged over an operational time corresponding to $50\,fb^{-1}$.
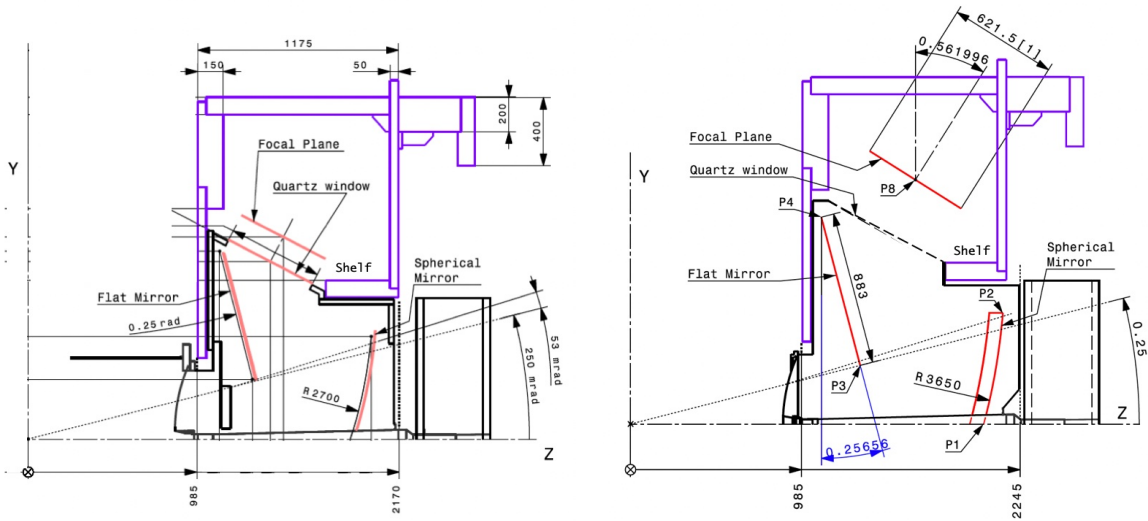
**Figure 66.** The optical geometries of (left) th<e original and (right) the upgraded RICH1.

A schematic (CAD model) of the RICH1 detector is shown in figure 57. RICH1 is aligned to the LHCb coordinate axes and occupies the region beyond the VELO exit window $947.5\,\mathrm{mm}$ $\leq z \leq 2245\,\mathrm{mm}$ and $\pm920\,\mathrm{mm}$ in $x$. The $z$ axis follows the beam line which is inclined at $3.6\,\mathrm{mrad}$ to the horizontal.

The optical layout of the upgraded RICH1 has been modified to reduce the larger hit occupancy expected in the central region of the detector. The occupancy has been halved by increasing the focal length of the spherical mirrors by a factor of approximately $\sqrt{2}$, which also improves the Cherenkov angle resolution by reducing mirror aberrations. A comparison of the previous and upgraded optical layout is shown in figure 66. As a consequence of the increased focal length, the photon detector planes have been moved parallel to and outwards from the beam line by approximately $270\,\mathrm{mm}$ ($225\,\mathrm{mm}$ in $y$, $150\,\mathrm{mm}$ in $z$). To minimise the material budget within the acceptance, lightweight carbon-fibre spherical mirrors are used and all other components of the optical system are located outside the acceptance. The average RICH1 material budget inside the LHCb acceptance is $\sim 4.8\%\,X_0$. Planar (flat) mirrors reflect the image from the tilted spherical mirrors onto the photon detector planes.

The arrays of MaPMTs, described in section 7.1.4, are located at the upper and lower focal planes, and each array occupies an active area of $605 \times 1199\,\mathrm{mm}^2$. The MaPMTs are shielded from the $60\,\mathrm{mT}$ fringe field of the LHCb dipole by magnetic shielding boxes made of ARMCO® iron, placed above and below the beam line outside the LHCb acceptance. These shields are retained from the original RICH1 detector, however with the so-called *shelves* cut off by $70\,\mathrm{mm}$ to ensure photon acceptance. To allow MaPMT column extraction and insertion for maintenance and installation, additional apertures of $666 \times 462\,\mathrm{mm}^2$ and $810 \times 56\,\mathrm{mm}^2$ were machined on the sides of each shielding box and covered with removable $10\,\mathrm{mm}$ thick plates. Inside the shielding boxes, the MaPMTs have additional local mu-metal shielding as described in section 7.1.2 and are able to work efficiently in fields of $3\,\mathrm{mT}$. To guide the new design, the magnetic field was simulated with the OPERA/TOSCA software and later measured at a position displaced by about $10\,\mathrm{cm}$ from the nominal MaPMT plane. These studies confirmed that the MaPMTs will operate in a magnetic field in the range $0.6$-$2.2\,\mathrm{mT}$, with the axial field below $1\,\mathrm{mT}$ for all MaPMTs.
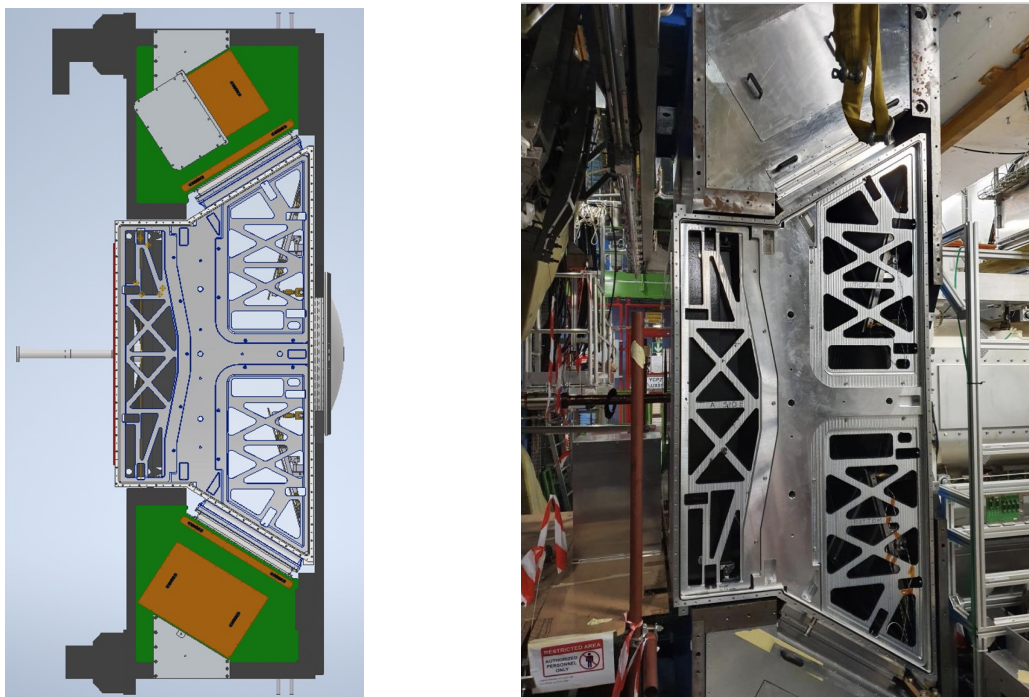
**Figure 67.** Left: side view CAD layout of the RICH1 gas enclosure. Right: photo of the RICH1 gas enclosure after its installation in the LHCb cavern.

### 7.2.1 Gas enclosure

The purpose of the gas enclosure is to contain the $C_4F_{10}$ gas radiator, to provide an optical bench for all optical components, and to ensure gas and light tightness. A schematic and a picture of the gas enclosure are shown in figure 67. The $C_4F_{10}$ radiator gas pressure follows the atmospheric pressure within $\pm 3$ mbar. The total gas volume is approximately $3.8\,\text{m}^3$.

The enclosure is machined from 30 mm thick aluminium alloy tooling plates. The six sides are bolted and epoxy-sealed at their edges, and internally sealed with flexible silicone sealant[58] to ensure leak tightness. The side faces of the gas enclosure are open to allow access for installation of mirrors and to the beam pipe. The structure has removable stiffening hatch-plates at the side apertures to prevent deflections of the structure when the side panels are removed, following loading with the mirrors, or when under ambient operational pressure. During normal operation, the sides are sealed by 15 mm aluminium panels. The maximum deflection of the superstructure is limited everywhere to 150 µm.

The upstream and downstream faces of the gas enclosure have apertures to allow passage of particles with minimum scattering within the LHCb acceptance. The upstream face attaches to 300 µm thick stainless steel bellows, the so-called *VELO seal*, which provides a gas-tight, mechanically compliant (longitudinal $\pm 10$ mm and transverse $\pm 1$ mm) seal to the downstream face of the VELO vacuum tank

The downstream face is closed by a low-mass (16.2 mm thick, estimated 0.7% $X_0$) exit window manufactured from a sandwich of two 0.6 mm thick carbon fibre skins filled with 15 mm of foam.[59] The window is sealed to a flange (a fixed fin machined in the beryllium beam pipe) using a 1 mm thick

---

[58]Bluestar CAF4 Silicone Sealant.

[59]Airex R82.80.

**Figure 68.** Pictures of (left) the spherical and (right) the bottom flat mirror assemblies.

opaque moulded silicone diaphragm,[60] as used in the original RICH1 [1]. The deflection of the window due to the gas enclosure pressure differential is approximately ±0.8 mm. All removable aperture covers (VELO seal, exit window, quartz window and side doors) are sealed with 4.5–5.8 mm diameter EPDM O-rings.[61] The gas enclosure is supported by the lower magnetic shield through mounts that allow its alignment to the nominal beam line. Unloaded with optical components but including the quartz windows, described below, the overall weight of the gas enclosure is approximately 1130 kg.

Square apertures above and below the beamline allow Cherenkov light to reach the MaPMTs, located behind. The apertures are sealed with polished fused silica windows 8 mm thick, of dimensions 655×475 mm². The windows are each fabricated from three equal-size panes, glued together along one edge and then glued into an additional frame. The six quartz panes were individually coated with an antireflective coating with approximately a quarter wavelength of $MgF_2$, which provides a gain in transmitted light at 270 nm of an additional 3.5%. The transmission has been measured on test samples over the range of interest of 270 nm to 500 nm to be better than 95%.

### 7.2.2 Mirrors

Four tilted CFRP spherical mirrors and 16 glass planar mirrors are used to focus the Cherenkov photons onto the photon detector planes, positioned outside the LHCb detector acceptance. Pictures of the mirror assemblies are shown in figure 68.

Each spherical mirror has a width of 740 mm, a height of 650 mm, a thickness of 33 mm and a radius of curvature of 3650 mm. Each mirror has a weight of 2.9 kg. The mirrors are arranged into four quadrants centred around the beam line. The inner corner of each mirror has a quarter-circle cutout in order to accommodate the beam pipe with a clearance of approximately 12 mm. In order to allow a mirror alignment with a precision of the order of tenths of mrad, each mirror is supported at the outer three corners by means of spherical rod end adjusters, bolted to a CFRP frame made of a 2-inch square tubular structure. The CFRP frame is divided along the $y$ axis into two C-shaped

---

[60]Dow Corning Sylgard 186, with 5% black pigment added.
[61]Ethylene Propylene Diene Monomer (M-class) rubber.

halves, each one supporting two mirrors positioned in the vertical direction, with a mirror-to-mirror separation of 3 mm to allow for alignment and as a clearance for deformations. The top (bottom) pair of mirrors are aligned to point to the same top (bottom) centre-of-curvature. Each C-shaped half has a weight of approximately 10 kg, including the weight of the mirrors, and is bolted to V-shaped blocks sitting on a cylindrical load rail positioned at the floor of the gas enclosure.

The planar mirrors have a width of 370 mm, height of 440 mm, thickness of 8 mm and a radius of curvature larger than 60 m. Each mirror has a weight of 3.3 kg. The 16 mirrors are arranged in two sets of eight mirrors each, positioned outside of the detector acceptance above and below the beam line, with a mirror-to-mirror separation of 3 mm as in the case of spherical mirrors. The tilt of the plane is 257 mrad to the vertical. Each mirror is bonded at its centre to a polycarbonate mount, bolted into machined pockets on four rigid 1-inch thick aluminium support frames, each frame supporting four mirrors. In addition, a polycarbonate ring centred on each mirror is bonded over a larger area than the polycarbonate mounts and secured to the support frame to retain the mirrors in case of failure of the polycarbonate mounts.

Each support frame with its mirrors weighs approximately 43 kg and is bolted to V-shaped blocks sitting on a rail bolted to the front panel of the gas enclosure.

The mirror quality is characterised by the diameter $D_0$ of the circle which contains 95% of the light intensity from a point source placed at the mirror centre of curvature (CoC) imaged at the CoC of the mirror. The $D_0$ for all planar and spherical mirrors was found to be better than the specification of 2.5 mm. The spherical and flat mirrors were individually aluminised with 10 nm chromium (adherence layer) and 100 nm aluminium (reflective layer). Additional enhancement in reflectivity and protection against oxidisation were ensured by coating the mirrors with 70 nm $SiO_2$ and 60 nm $HfO_2$. The reflectivity of all mirrors is everywhere > 90% in the wavelength range 260 nm to 500 nm, peaking around 95%.

There are three stages to the mirror alignment process: prealignment on the optical rig before installation, survey in situ, and alignment with data. The prealignment on the optical rig is crucial to the process as the CoC of the mirrors are inaccessible when the mirrors are installed in RICH1. At this stage the upper and lower pairs of the two spherical mirrors are aligned to a common CoC. For the planar mirrors, the top and bottom set of mirrors are aligned to form a single plane parallel to both support frames which will point to the corresponding photon detector plane. Survey points on the spherical mirror frame and on the flat mirror backing plates then reference the CoCs and the flat mirror tilts respectively.

### 7.2.3 Photon detector region

The MaPMT columns make use of common cooling and electronics, but have custom support mechanics and services. Eleven RICH1 columns are arranged side by side to form an 11×22 array of elementary cells. One such array is placed below the beam pipe, with a second above. Both arrays are horizontal in the plane perpendicular to the beam pipe, and are tilted with an angle of 562 mrad with respect to the vertical towards the interaction point.

RICH1 columns are supported from the face opposite the MaPMTs and held in place on rails allowing for easy removal for maintenance. Pivoting around the rails is prevented by precision alignment pins at both ends of the columns on the same face as the MaPMTs. The rails are held at the correct angle and aligned on the MaPMT chassis, a mechanical support structure which acts
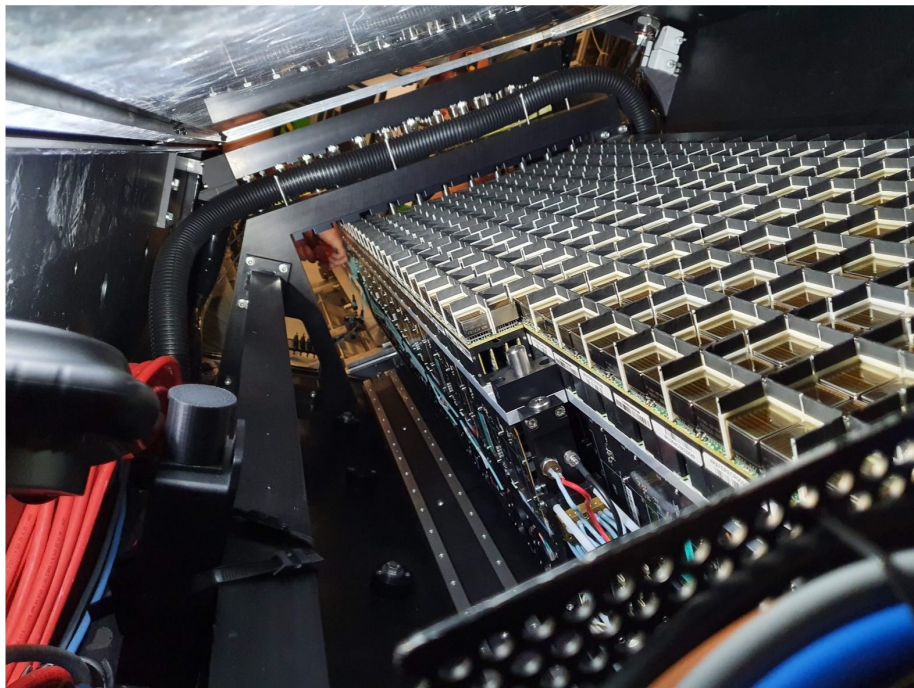
**Figure 69.** Picture of RICH1 columns while being inserted into their support structure, the lower MaPMT chassis. The rails and alignment structures are visible on the left side. The chassis is mounted to the soft-iron magnetic shielding that surrounds the MaPMT region. The push connector for the copper-carried services is visible as well in the left-most column. Reproduced from [109]. © 2022 IOP Publishing Ltd and Sissa Medialab.

as an optical bench for the columns. At each end, the chassis has precision slots to match the alignment pins mounted on the columns. A picture of the lower chassis hosting the corresponding columns is shown in figure 69.

Columns are inserted and removed along the rails from one side of the chassis. At the far end, services carried on copper cables (low-voltage, high-voltage and monitoring cables) are connected to the column through a push connector. Services that require manual disconnection (cooling and data fibres) are attached to the column at the extraction side. This allows the column to be removed for maintenance by accessing only one side of the detector.

### 7.3 RICH2 photon detector planes

RICH2 is located downstream of the dipole magnet, covering an angular acceptance of $\pm 15 - \pm 120$ mrad in the horizontal direction and $\pm 15 - \pm 100$ mrad in the vertical direction. It extends between 9500 and 11832 mm along the $z$-axis. The photon detector planes are located on the LHCb detector sides.

The superstructure including the large entry and exit windows, the two magnetic shields located on each side and the optical system (mirrors and their supports) inside the enclosure and the quartz windows are retained without any modification [1]. These components have been demonstrated to sustain the occupancy and radiation levels expected in LHCb upgrade running conditions.

The main upgrade of RICH2 concerns the installation of the new photon detector planes. The integration of the new photon detection chain and its services in RICH2 requires a new mechanical

structure to host the photon detector planes. The structure holding the columns is made of aluminium profiles, allowing to place up to fourteen RICH2 columns side by side with a clearance of 1 mm between them. This arrangement of the columns form the array of the photon detectors. Twelve columns were found to be enough to cover the detector acceptance and were installed, with two empty slots at the periphery, as shown in figure 70, for a total of 288 ECs.

Two such structures, or *racks*, are installed into the magnetic shields on the LHCb Side A and Side C. The arrays are placed vertically and parallel to the RICH2 structure which is not exactly perpendicular to the LHC beam. Both arrays are tilted by an angle on the horizontal plane of 1.065 rad with respect to the LHCb $x$ axis, with the first column towards the interaction point being the farthest from the beam pipe.

Each rack is installed on a trolley, mounted on rails, which allows the movement of the complete photon detector system perpendicular to the focal plane for installation, maintenance and dismounting of the rack. Furthermore, anchor points between the rack and the trolley, allow to adjust the position of the racks. Inside each rack, top and bottom rails allow to slide each column individually in the direction perpendicular to the focal plane. The rails are made of hard anodised aluminium and fastened on base plates. On the upper side of each rack, cable chains connected to all columns route the electric cables and optical fibres to the patch panel located above the rack. Thanks to this design, each column can be extracted fully independently, allowing to continue the operations on the other columns. This arrangement greatly simplifies installation and maintenance with respect to Run 1-2 setup.
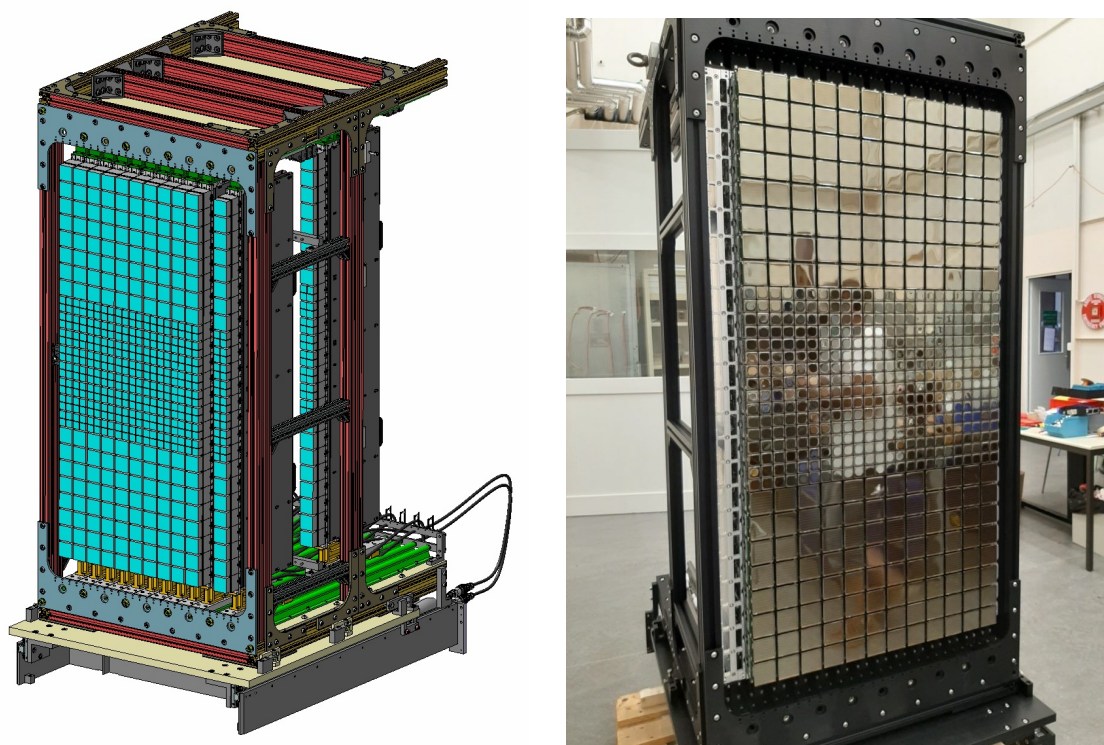


**Figure 70.** Fully assembled and commissioned RICH2 photon detector array. Left: CAD view; right: photograph taken in the assembly area. Reproduced from [109]. © 2022 IOP Publishing Ltd and Sissa Medialab. All rights reserved.
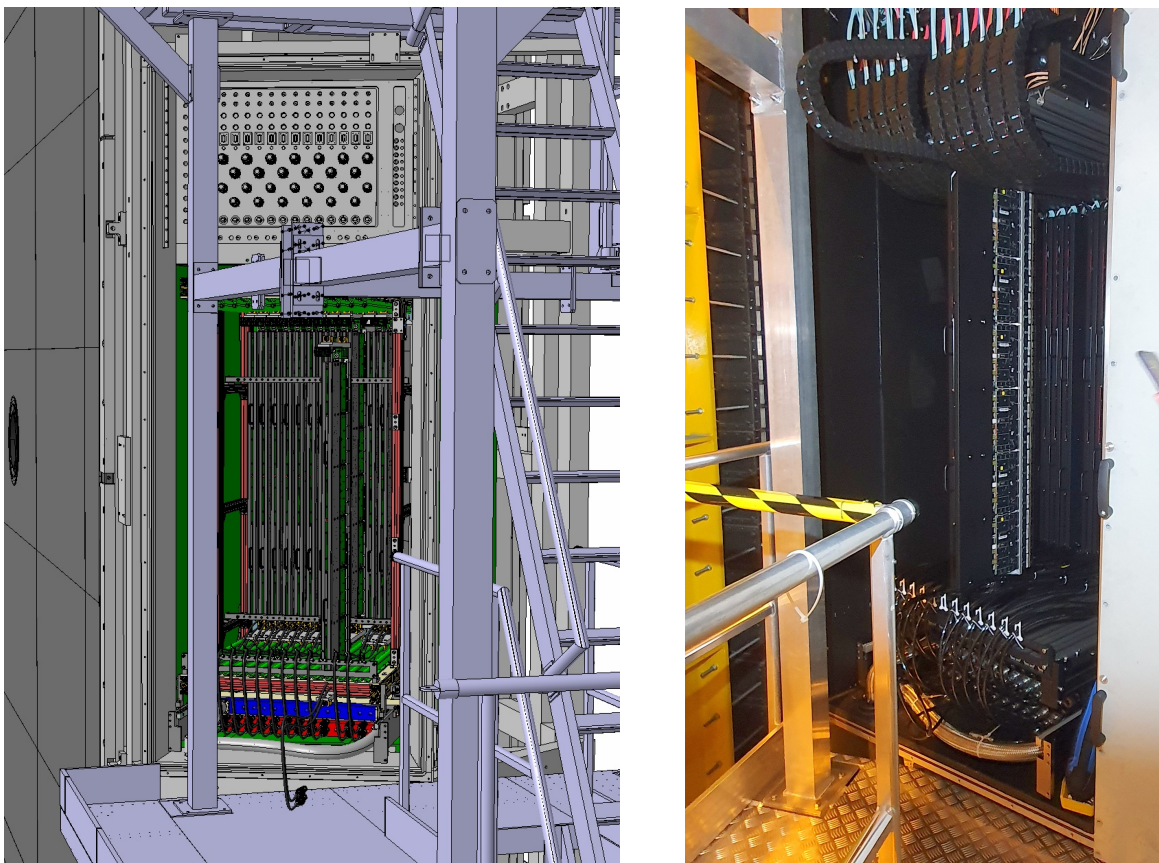
**Figure 71.** RICH2 photon detection system inside its enclosure. Left: CAD view; right: photograph.

Two cooling manifolds are located under the trolley. The first distributes the cooling fluid to each column in parallel while the second collects back the fluid. The manifolds are connected to the two transfer lines via custom-made feedthroughs. Each column is connected with two polyurethane hoses and double shut-off couplers to the manifolds. The couplers offer the possibility to disconnect a column avoiding purging the fluid or closing any valve.

A patch panel installed at the top of the photon detector enclosure provides the interface where all services are connected to the columns. It provides the connections to the safety ground, LV, HV, DSS, TFC, ECS and data transmission lines. The interface has been designed to provide an easy way to disconnect the columns and to ensure a light- and gas-tight enclosure. Nitrogen is flushed permanently to ensure a dry atmosphere in the enclosure, minimising the risk of condensation and therefore maintaining a good dielectric environment.

## 7.4 Monitoring, controls and data acquisition

Monitoring, controls and data flow of the RICH system are fully integrated in the LHCb online structure, with TFC and ECS controlling timing and detector control and configuration, and optical links transporting data from FE to BE electronics.

### 7.4.1 Detector control system

The RICH detector control system (DCS) is a subset of the ECS that controls LV power supplies and monitors detector safety and environmental parameters for both the photon detection system and the gas radiators. A variety of sensors are interfaced using an embedded local monitoring board (ELMB) [121], while a large number of temperature sensors mounted on baseboards, FEBp and backboards are read out using the ADC of the SCA readout chip with a resolution of approximately 0.5°C. As these temperature readings are only available when the detectors are operational, a separate Pt100 sensor mounted on every photon detector column, read out via the ELMB in a 4-wire configuration, provides additional information when the detectors are switched off. In addition, two temperature sensors are installed at the input and output of each cooling manifold. A further check on the circulation of the cooling fluid is performed by using two pressure sensors per cooling manifold. Temperature and humidity in the photon detector enclosure are monitored with dedicated sensors. A safe switch-on procedure is implemented, where automatic configuration of the SCA readout chip at power-up allows the read out of the temperature sensors without any action of the operator. If the automatic configuration fails, the detectors are switched back off.

When the detector is operational, parameters from sensors allow to constantly monitor the electronics, the detector environment and the condition of the cooling system to ensure safe operation. These parameters are used by the ECS to check for possible signs of abnormal conditions and take automatic actions. Possible actions include switching off the LV system and in extreme conditions also the MaPMT HV system. Finally a smaller number of sensors are connected to the DSS running on a PLC system with many redundancies, ensuring an additional layer of detector safety.

The DCS also collects information about the temperature and pressure of the Cherenkov gas radiators, as any change in the gas density affects directly the refractive index. Both RICH radiator gases are left free to follow the atmospheric pressure variations. Their temperature is the same as the LHCb cavern temperature, which is kept stable at a value of $20.0 \pm 0.5$°C. Temperature and pressure are recorded in a database and are extracted by the LHCb event reconstruction software to calculate the correct refractive index.

### 7.4.2 DAQ controls, monitoring and data flow

The RICH control system operates on a minimal set of devices, composed by one PDMDB, one SOL40 and one TELL40, and replicates the commands to up to thousands of FE and about a hundred BE devices. The SOL40 provides, via each of its 48 bidirectional optical links, the reference 40 MHz clock and TFC commands to the PDMDBs through their TCMs. Uploaded commands are decoded by the combination of a GBTx and an SCA readout chip while the FE data is sent by the PDMDB to the TELL40, via unidirectional optical links, through DTMs, depending on the variant of the PDMDB.

The relevant PDMDB registers are periodically monitored through the GBT server and differences between writing and readings in any of such registers will raise errors.

TFC commands are used at power on and configuration time to enable communication buses, initialise temperature sensors, load the firmware on the PDMDB FPGAs and set thresholds on the CLARO discriminators. A stateless implementation of the PDMDB firmware, where data are transported transparently towards the links to the BE, was chosen in order to minimise the impact of SEE due to radiation. This approach required to move to the TELL40 the association of event data with the corresponding BXID, needed for event building.

Since the pixel occupancy in the RICH detector varies by orders of magnitude between RICH1 and RICH2, their centre and periphery, and the highest occupancy region corresponds to a relatively small fraction of the acceptance, the TELL40 input bandwidth had to be carefully optimised to minimise the number of needed boards to save on costs.

Where possible, the bandwidth was kept under control by applying a simple lossless compression algorithm. Then, after a Monte Carlo simulation of the hit occupancy, two different configurations of TELL40 have been implemented: where the bandwidth could be balanced, 48 input links were used, while 24 where used elsewhere, allowing the number of needed boards to be reduced by about 30%. Each TELL40 merges and compresses BXID aligned data from all the connected input links into packets that are transferred to the host event builder (EB) server via PCIe, up to a maximum bandwidth of 102 Gbit/s. Although the EB network allows for a maximum average bandwidth of 90 Gbit/s, after optimisation none of the RICH EB servers exceeded 70 Gbit/s, being limited by the maximum instantaneous bandwidth rather than its average.

## 7.5 Calibration of photon detectors and front-end electronics

The single photon detection efficiency is a crucial parameter of the RICH system, and is mainly driven by intrinsic properties of the MaPMTs, such as the photocathode quantum efficiency, the collection efficiency at the first dynode and the single photon gain. Additional contributions to the detection efficiencies arise from the anode signal digitisation, provided by the corresponding CLARO channel by means of a programmable threshold.

Dedicated calibration procedures are used in order to minimise the inefficiency arising from the threshold setting, while ensuring the suppression of noise hits due to the MaPMT and FE electronics pedestals. In addition, the same procedures allow to monitor the single photon gain variation with time and ageing of each MaPMT channel, and the stability of each CLARO channel.

Each CLARO channel is calibrated by using DAC scans, where a known charge is injected at the input in 256 steps of $15.6 \times 10^3$ electrons worth of charge each, for different CLARO settings. It allows to determine the conversion between a threshold DAC code and the corresponding absolute charge. An example of the output of DAC scans performed on a RICH2 column is reported in figure 72.

Threshold scans are performed in order to find the set of front-end working points that maximise the single photon efficiency. The photon detection chain is illuminated at very low light intensity (single photon regime), and the threshold of the CLARO comparator is decremented in unit steps. It is found that the best operational points correspond to threshold settings that are five steps above the pedestal. The distribution of the threshold settings for all the RICH2 channels, converted to absolute charge as determined from DAC scans, is compared to the single photon peak distributions of the corresponding MaPMT anodes at different HV values in figure 73. Threshold scans, providing the integral pulse height spectrum, are also used to estimate the single photon peak position for each channel, allowing to implement the monitoring of gain variations with MaPMT ageing.

### 7.5.1 Time alignment

The prompt Cherenkov radiation and focusing mirror optics lead to the nearly simultaneous time-of-arrival to the detection plane of photons from a track in the RICH detector. This unique feature allows the application of a time gate at the FE electronics to exclude out-of-time background hits from the output data whilst accepting the photon signals within a narrow time interval. Figure 73 shows the distribution of
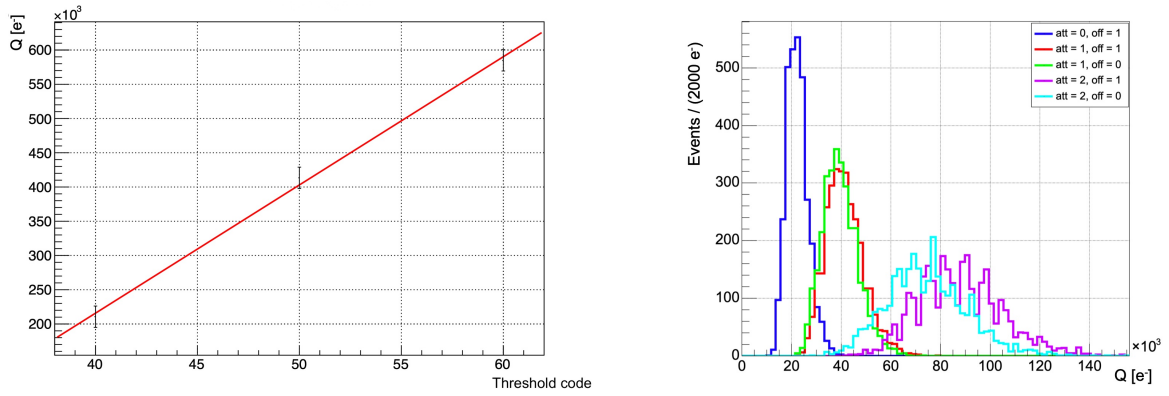
**Figure 72.** Typical output of the DAC scans procedure. On the left, the calibration of a single CLARO channel with offset bit enabled and no attenuation is shown. The charge (in units of electron charge) corresponding to a threshold DAC code (th) is determined by the linear relation $Q = Q_0 + Q_{th} \cdot$ th. On the right, the distribution of the charges (in units of electron charge) corresponding to one threshold step ($Q_{th}$) for a RICH2 column is shown. The linearity of the threshold setting as a function of the injected charge is found to be excellent for all attenuation and offset values.
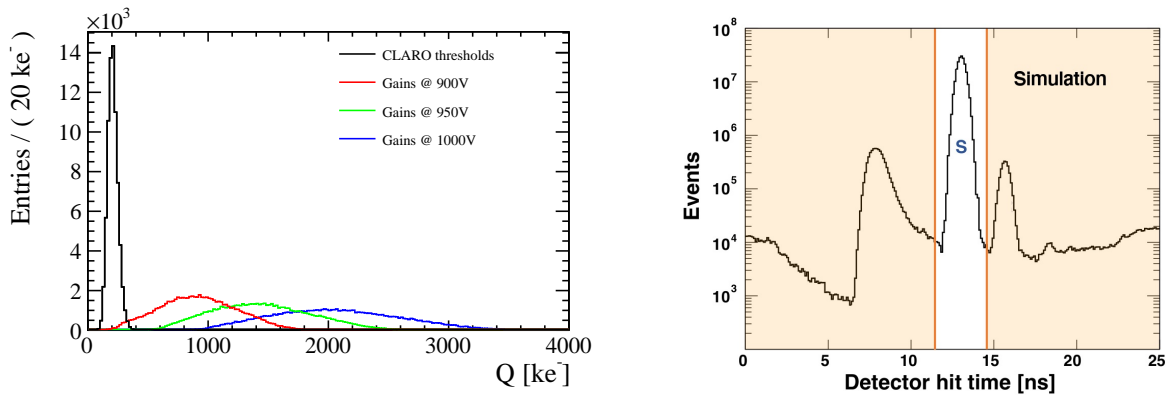


**Figure 73.** Left: distribution of RICH2 thresholds of the CLARO comparator converted into absolute charge (black). The mean and standard deviations of the distribution are $(207.58 \pm 0.16) \times 10^3$ electrons and $(39.64 \pm 0.10) \times 10^3$ electrons. The threshold settings can be compared to the pixel gain at 900 V (red), 950 V (green) and 1000 V (blue). Reprinted from [111], Copyright (2023), with permission from Elsevier. Right: RICH1 simulated photon detector hit time distribution showing the signal peak (S) and a possible time gate in the front-end electronics. Reproduced from [122]. CC BY 4.0.

photon hit times in RICH1 from a simulation. The signal peak spans approximately 2 ns due to the spread of primary interactions in LHCb, which dictates the minimal width for the FE time gate. In practice, the combination of CLARO time walk, channel-to-channel variations, MaPMT transit time spread and digital sampling rate at the FE electronics require a time gate whose width must be set to 3.125 ns or doubled to 6.250 ns if needed. In addition to the background from the beam interactions in figure 73, the time gate can exclude sensor noise, such as MaPMT cross-talk and afterpulses [122]. The achieved background reduction significantly improves the PID performance using the RICH pattern recognition algorithms.

Gate generation and the time alignment procedure are implemented in the PDMDB FPGA firmware. The time gating logic exploits the deserialiser embedded in every input-output logic block

of the FPGA which can operate at gigabit rates. The deserialiser samples the CLARO signals using both edges of the 160 MHz clock and shifts the sampled data at 320 Mbit/s into an 8-bit shift register. This byte can be checked against specific signal patterns using a lookup table, which is a readily available memory resource with a small logic footprint in the general purpose logic of the FPGA. If the CLARO signal pattern matches one of the configured lookup table patterns, a hit is registered on the 40 MHz system clock edge. The programmable lookup table allows flexibility between data-taking modes such as different time gate widths, edge detection and basic spillover checks.

The time gate is applied at a fixed latency with respect to the LHCb clock. The FPGA receives the 40 MHz system clock and 160 MHz sampling clock from the GBT, where the clock phases can be adjusted over the 25 ns range in fine steps of 49 ps. This allows the position of the time gate to be fine-tuned with respect to the signal time-of-arrival in the RICH detector.

## 7.6 Expected performance

The performance of the RICH detectors is evaluated using the LHCb simulation framework. The main parameters used to evaluate the performance are the Cherenkov angle resolution $\Delta\theta_C$ and the photoelectron yield. The Cherenkov angle resolution is estimated starting from the single-photon resolution, $\sigma_\theta$, which can be split into roughly independent contributions:

- *chromatic*, due to the chromatic dispersion of the radiators which leads to a dependence of the Cherenkov angle on the photon energy;

- *emission point*, due to the tilting of the spherical mirrors which leads to a smearing of the Cherenkov angle depending on the point of emission of the photons along the track;

- *pixel error*, due to the finite size of the MaPMT pixels;

The single photon resolution is investigated from simulated events, using a simplified Cherenkov angle $\theta_C$ reconstruction, assuming that the photon is emitted from the track at the middle of the radiator length and that it hits the centre of the relevant detector pixel. The reconstructed Cherenkov angle is then compared with the expected value from the simulation. The photoelectron yield, $N_{ph}$, i.e. the number of Cherenkov photons emitted by a track detected on a Cherenkov ring by the MaPMTs, is studied using high momentum tracks with $\beta \sim 1$ where the Cherenkov angle as well as the number of emitted photons are maximal. For these tracks, the total Cherenkov angle resolution is given by:

$$\Delta\theta_C = \frac{\sigma_\theta}{\sqrt{N_{ph}}} \oplus C_{tracking}, \tag{7.1}$$

where the constant factor $C_{tracking}$, added in quadrature, is the contribution from track reconstruction uncertainties. This includes the uncertainties associated to multiple scattering, to the track curvature inside the RICH radiator volumes, and to the intrinsic resolution of the tracking detectors. As such $C_{tracking}$ is a function of momentum, and it takes an average value of 0.35 mrad while asymptotically tending to approximately 0.15 mrad at high momentum as determined from simulation studies [108].

### 7.6.1 Simulation setup and typical output

The simulation used to obtain the results presented in this section includes the up to date information on the RICH geometry and on the properties of its optical system and photon detection chain. The

properties of the individual photon detector channels, such as the MaPMT gain and noise, are accounted for; in addition, background sources such as cross-talk, afterpulses and scintillation photons produced by charged particles in the CF$_4$ radiator are included as well [123]. In the simulation, the values related to the mentioned aspects of the individual channels are based on the data acquired during quality assurance procedures with a 900 V bias voltage for the MaPMTs. The study is performed using a sample of 10 000 events containing a $B_s^0 \to \phi\phi$ decay as typical signal events. The simulated data are obtained with the standard upgrade configuration, corresponding to an instantaneous luminosity of $\mathcal{L} = 2 \times 10^{33} \, \mathrm{cm^{-2} \, s^{-1}}$. The PID performance, after the application of the reconstruction algorithms, is reported for tracks in the momentum range of 2–100 GeV/$c$ and $p_T$ larger than 0.5 GeV/$c$. Only the tracks that traverse the full LHCb tracking system acceptance are used.

The average MaPMT quantum efficiency curve, which is used in the simulation, is shown in figure 74, together with a typical PID performance curve, representing the probability to misidentify a pion as a kaon versus the probability to correctly identify the particle as a kaon. The expected hit occupancy in the MaPMT as a function of the MaPMT identifier is reported in figure 75.

### 7.6.2 Performance studies

A single particle simulation is used to evaluate the best achievable photoelectron yield ($N_{\mathrm{ph}}^{\mathrm{optimal}}$) and the Cherenkov angle resolution. The simulation is configured to provide 80 GeV muons, which ensures that the tracks are saturated and are not significantly curved by the magnetic field to minimise the uncertainty arising from the tracking system. In addition, the acceptance region where the RICH performance is expected to be optimal is used: the polar angle of the tracks is required to be in the 90–180 mrad and 40–90 mrad for RICH1 and RICH2, respectively. The results are summarised in table 9. Consistent results are found when using $B_s^0 \to \phi\phi$ events provided that the selected tracks fulfil the constraints described above for the single particle simulation. The photoelectron yield decreases when the optimal track requirements are relaxed. This is shown in table 9 where a typical photoelectron yield, $N_{\mathrm{ph}}^{\mathrm{typical}}$, is reported. The typical photoelectron yield values are lower than the optimal ones, mainly due to the limitations in the acceptance due to the beam pipe region. As it can be extracted from table 9, the contributions to the total Cherenkov angle resolution from the RICH and tracking systems at high momentum are of similar magnitude, and slightly dominated by the $C_{\mathrm{tracking}}$ term. The $C_{\mathrm{tracking}}$ contribution used here is the asymptotic limit of approximately 0.15 mrad at high momentum, dominated by the intrinsic resolution of the tracking detectors.

## 8 Calorimeters

To cope with the new LHCb readout scheme, the FE and readout electronics of the electromagnetic and hadronic calorimeters have been entirely redesigned and replaced. Moreover, two subdetectors of the previous calorimeter system, namely the Scintillating Pad Detector (SPD) and the PreShower (PS) [125], have been removed, given their reduced role in the new LHCb all-software trigger.

The layout of both the electromagnetic calorimeter (ECAL) and hadronic calorimeter (HCAL) remains unchanged for the upgrade. A complete description of the geometry and technological aspects can be found in ref. [125]. To minimise the required modifications, the ECAL and HCAL calorimeter modules, their photomultiplier tubes, Cockroft Walton (CW) bases and coaxial cables were also maintained unmodified. However, to keep the same average anode current of the phototubes
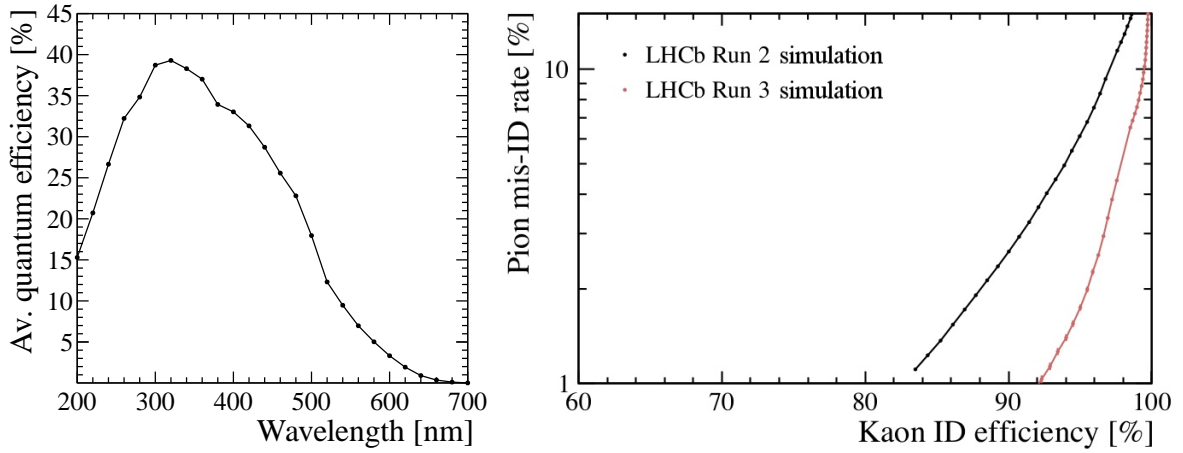
**Figure 74.** Left: average quantum efficiency of the MaPMTs used in the RICH detectors. Right: a typical PID performance of the kaon identification obtained from the LHCb software for the configuration described in the text (red). A corresponding curve for the Run 2 conditions (prepared using the simulation with LHCb Run 2 geometry and luminosity, as reported in ref. [124]) is shown for reference (black).
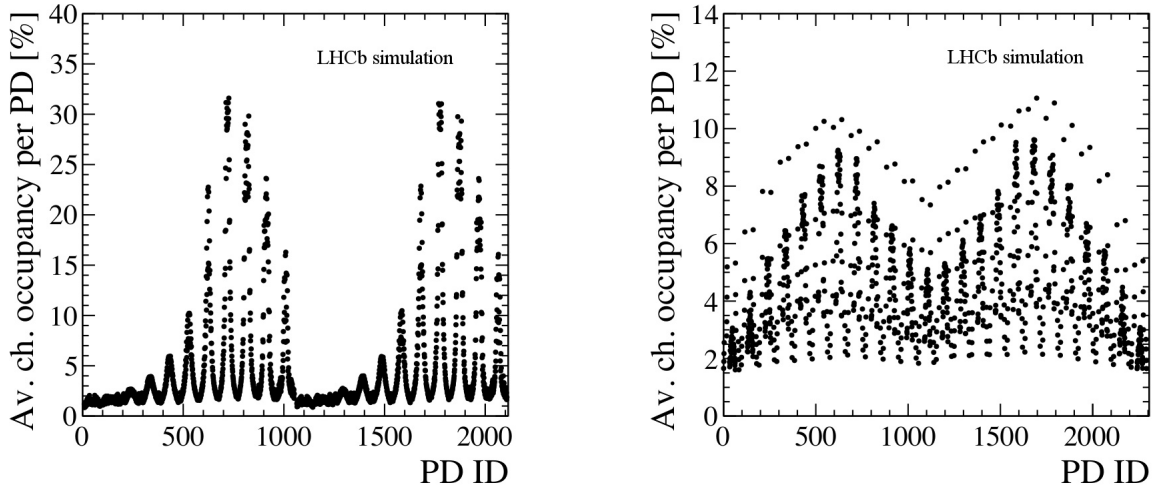


**Figure 75.** Average expected occupancy per channel for different MaPMTs in the (left) RICH1 and (right) RICH2 detector.

**Table 9.** Simulated performance of the upgraded RICH detectors. For RICH2, the values are given for the inner detector regions populated with the 1-inch MaPMTs.

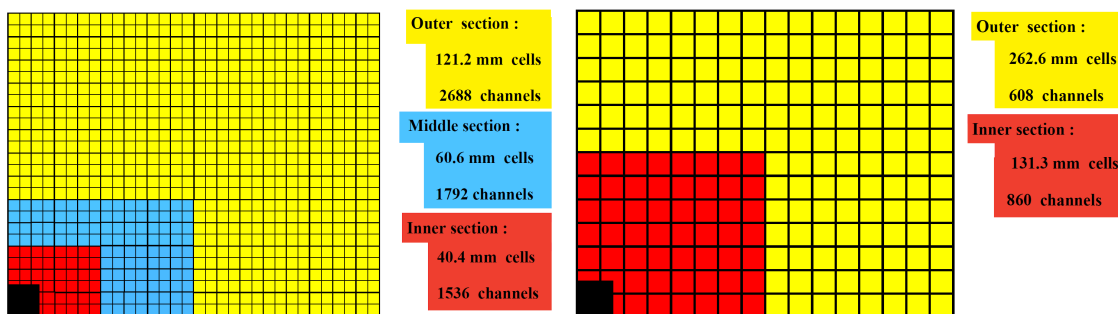| | Photoelectron yield | | Cherenkov angle resolution [mrad] | | | | |
|---|---|---|---|---|---|---|---|
| | $N_{ph}^{optimal}$ | $N_{ph}^{typical}$ | chromatic | emission point | pixel | $\sigma_\theta$ | $\Delta\theta_C$ |
| RICH1 | 63 | 59 | 0.52 | 0.36 | 0.50 | 0.81 | 0.18 |
| RICH2 | 34 | 30 | 0.34 | 0.32 | 0.22 | 0.52 | 0.17 |

**Figure 76.** Lateral segmentation of (left) the ECAL and (right) the HCAL. One quarter of the detector front face is shown. Reproduced from [125]. CC BY 3.0.

at the higher luminosity, their high voltage has been reduced implying an increased gain of the amplifier-integrator in the FE cards. This modification is described in section 8.2.1.

The FE electronics boards have been fully redesigned to comply with the 40 MHz readout frequency, but their number and format have been chosen such that they are compatible with the existing crates and racks, as described in section 8.2.

The decision to keep the calorimeter modules, their PMTs and CW bases, assumes that they can operate at radiation levels corresponding to the foreseen integrated luminosity. This is discussed in section 8.1.3.

## 8.1 General detector structure

The LHCb calorimeter system presents a classical structure of an electromagnetic calorimeter followed by a hadronic calorimeter. The most demanding performance constraint concerns the identification of electrons and photons with an optimal energy resolution requiring the full containment of the showers from high energy particles. For this reason, the thickness of the ECAL was chosen to be 25 radiation lengths [126]. On the other hand, the trigger requirements on the HCAL resolution do not impose a stringent hadronic shower containment condition, thus its thickness is limited to 5.6 interaction lengths [127], due to space limitations.

To account for different hit densities across the calorimeter surface, the ECAL is segmented laterally in three regions referred to as inner, middle and outer, with increasing dimensions going from the beam pipe outwards, as shown in figure 76. The HCAL is segmented in two regions with a larger granularity with respect to the ECAL, given the typical spread of hadronic showers. The regions are segmented in cells of transverse dimensions roughly projective with respect to the interaction point. Their dimensions are optimised to provide uniform measurements of the *transverse energy*, $E_T = E_c \sin\theta$, where $E_c$ is the energy measured by a cell and $\theta$ is the angle between the vector pointing to the centre of the cell from the interaction point. This quantity is particularly useful for hadron selection at trigger level.

The two calorimeters share the same basic detection principle: scintillation light from plastic scintillator modules is transmitted to the PMTs by wavelength-shifting fibres. Fibre bundles from calorimeter modules are then fed to the PMTs. In order to have a constant energy scale across the calorimeter surface the gain of the ECAL and HCAL PMTs is set proportionally to the distance from
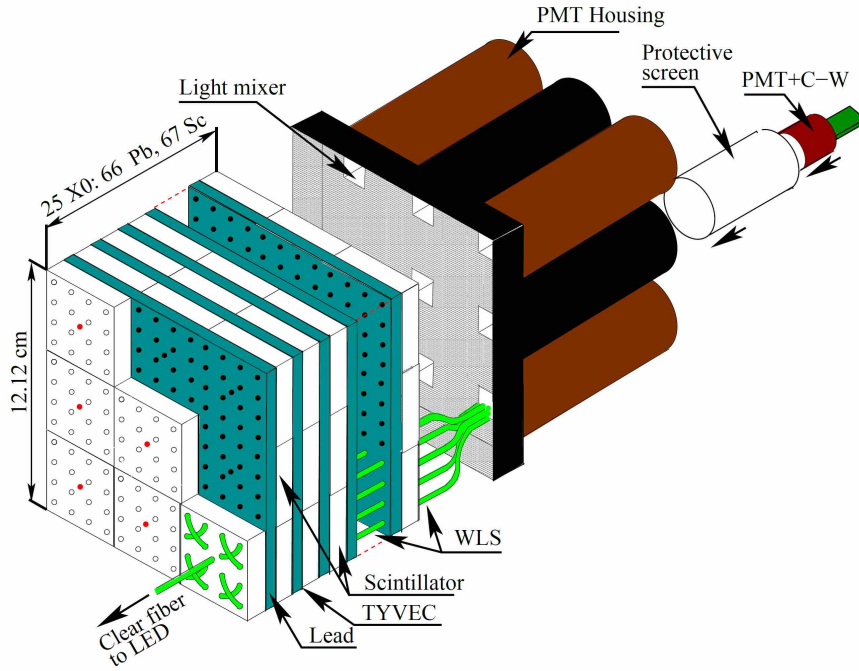
**Figure 77.** Schematic of an ECAL cell. Reproduced from [128]. CC BY 4.0.

the beam pipe of the corresponding modules. Since the light yield of an HCAL module is a factor 30 less than ECAL modules, the HCAL phototubes operate at higher gain.

### 8.1.1 The electromagnetic calorimeter

The ECAL front surface is located at about 12.5 m from the interaction point. The square cell sizes for the inner, middle and outer regions are 121.2 mm, 60.6 mm, 40.4 mm side for outer, middle and inner regions, respectively, and scale with the distance from the beam-pipe in order to make the particle rate per cell roughly uniform. The outer dimensions of the ECAL match projectively those of the tracking system, $\theta_x < 300$ mrad and $\theta_y < 250$ mrad, while the inner angular acceptance of ECAL is limited to $\theta_{x,y} > 25$ mrad around the beam pipe, where $\theta_x$ and $\theta_y$ are the polar angles in the $xz$ and $yz$ planes in the LHCb reference frame.

The ECAL cells have a *shashlik* structure, as shown in figure 77, with alternated scintillator (4 mm) and lead (2 mm) layers. The scintillation light readout is performed by dedicated phototubes PMTs.[62] The total number of cells is 6016. The energy resolution of a given cell, measured with a test electron beam, is parametrised as [125]:

$$\frac{\sigma(E)}{E} = \frac{(9.0 \pm 0.5)\,\%}{\sqrt{E}} \oplus (0.8 \pm 0.2)\,\% \oplus \frac{0.003}{E \sin\theta} \tag{8.1}$$

where $E$ is the particle energy in GeV, $\theta$ is the angle between the beam axis and the line from the LHCb interaction point to the centre of the ECAL cell. The second contribution is a constant term taking into account mis-calibrations, nonlinearities, energy leakage out of the cell and other effects, while the third term is due to the noise of the electronics which is evaluated on average to 1.2 ADC counts [1].
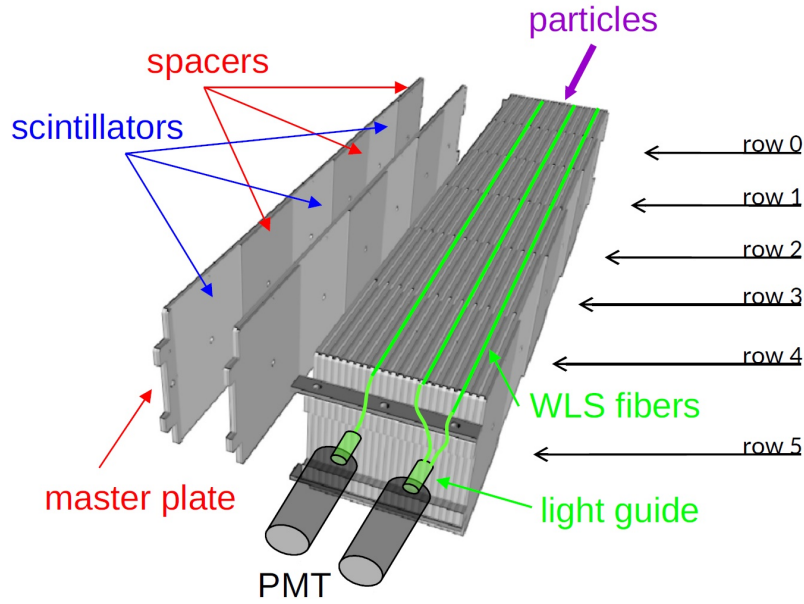
---

[62]Hamamatsu R7899-20.

**Figure 78.** Schematic of an HCAL cell. Reproduced from [128]. CC BY 4.0.

### 8.1.2 The hadronic calorimeter

The HCAL is a sampling tile calorimeter with a thickness of 5.6 interaction lengths. The sampling structure consists of staggered iron and plastic scintillator tiles mounted parallel to the beam axis (figure 78) to enhance the light collection. The same PMT type as in ECAL is used for the readout. The HCAL has a total of 1488 cells arranged in an inner and an outer region, segmented in square cells with sides of 131.3 mm and 262.6 mm, respectively. The energy resolution, as measured in beam tests with pions, is parametrised as follows:

$$\frac{\sigma(E)}{E} = \frac{(67 \pm 5)\,\%}{\sqrt{E}} \oplus (9 \pm 2)\,\%$$

(8.2)

where $E$ is the deposited energy in GeV [128].

### 8.1.3 Radiation effects and ageing

The effects of radiation on the calorimeter system components was assessed in a series of measurement campaigns at different irradiation facilities. The radiation resistance of calorimeter components such as scintillators, wavelength shifter fibres, PMTs and CW bases has been studied to extrapolate their lifetime over the foreseen integrated luminosity of the upgraded LHCb. First irradiation tests performed on an ECAL module prototype, indicated that the performance of the innermost cells remains satisfactory up to $\sim 25$ kGy at the position of the shower maximum, corresponding to about 20 fb$^{-1}$ of integrated luminosity [125]. Tests conducted at CERN PS IRRAD facility with 24 GeV protons and on ECAL inner modules, placed in the LHC tunnel during Run 1, confirmed this result. These measurements indicated that the innermost region of ECAL will have to be replaced during the LHC long shutdown 3 (LS3). Radiation tolerance of PMTs and CW bases was assessed with irradiation tests

conducted at a 50 GeV proton beam.[63] The CW bases remained operational up to doses of 15–20 kGy corresponding to 30–40 fb$^{-1}$ of integrated luminosity in the central ECAL cells. The replacement of CW bases can be easily performed during annual shutdown periods. Replacement of about 500 CW bases is estimated to be needed over the full upgrade programme. The transmittance loss of the PMT window in the wavelength peak range did not exceed 5%, ensuring that the radiation tolerance of the PMT entrance window will maintain sufficient transparency during the full upgrade lifespan.

As far as the HCAL is concerned, only the tile modules will suffer radiation damage effects, since the PMTs are installed behind the calorimeter and shielded by the iron. The radiation damage of the HCAL modules has been studied during Run 1, using calibration data obtained with a $^{137}$Cs source. The relative light yield of the scintillator tile rows (see figure 78) was measured with respect to the rearmost tile row, which receives the smallest dose and is not expected to suffer from significant radiation damage. A reduction of the light yield at the level of 15% was observed for the central HCAL cells after collecting 3.4 fb$^{-1}$ of luminosity. This result indicates that the HCAL innermost cells cannot survive the full lifetime of the upgrade. Since it is not possible to replace these cells in an easy way during the LHC runs, they have been removed and replaced by slabs of tungsten absorber, to mitigate the low energy particle background on the innermost regions of M2 and M3 muon stations (see section 9.5). The impact of this reduction in acceptance on the physics programme was estimated to be negligible as the information from the HCAL innermost cells will not affect significantly the software trigger analysis.

## 8.2 Electronics

The two LHCb calorimeters share the same electronics which consists of:

- a front-end board (FEB) (described in section 8.2.1), where the PMT signals are amplified, shaped and digitised and then, after proper formatting, shipped to the back-end electronics;

- a calorimeter control card unit (3CU) board (detailed in section 8.2.5) to distribute clocks and ECS commands to the FEBs;

- a set of calibration and monitoring boards (section 8.2.6).

The FEBs provide the digitised analog data in the form of $E_T$ measurements. They also perform simple data preprocessing, to send to the software trigger some precalculated quantities such as simplified energy clusters, calibrated energy measurements and global quantities such as total energy and hit multiplicity. This preprocessing is identified for simplicity as low level trigger (LLT), in analogy with the L0 hardware stage of the previous trigger scheme [7] and is discussed in section 8.2.1 and section 8.2.4. The information provided by the LLT is sent to the high level trigger (HLT) through the TELL40 for a possible use in event reconstruction or selection (see section 11.3.2).

In order to keep the same average anode current of the PMTs and minimise the ageing effects at the higher luminosity, the high voltage, and consequently the gain of the photomultipliers, must be reduced by a factor of five. Therefore, a partial gain compensation of a factor 2.5 was implemented into the FE electronics allowing also to double the $E_T$ dynamic range of the calorimeter system. The gain compensation can be tuned within a factor ranging from 0.5 to 2.0, by configuring the analog electronics which allows large flexibility in adapting to physics requirements. The equivalent input

---

[63]The beam is provided by the IHEP U-70 accelerator in Protvino, Russia.

**Table 10.** Summary of the requirements for the calorimeter analog FE.

| Parameter | Requirements |
|---|---|
| Energy range | $0 \leq E_T \leq 10\,\mathrm{GeV}$ (ECAL) |
| Calibration/Resolution | 4 fC/5 MeV $E_T$ per ADC count |
| Dynamic range | 4096-256 = 3840 counts: 12 bits |
| Noise | $\lesssim 1$ ADC counts (ENC < 4 fC) |
| Termination | $50 \pm 5\,\Omega$ |
| Baseline shift prevention | Dynamic pedestal subtraction |
| Max. peak current | 4–5 mA over 50 $\Omega$ |
| Spill-over residue level | $\pm 1\%$ |
| Integrator peak plateau | < 1% variation in $\pm 2$ ns |
| Linearity | < 1% |
| Cross-talk | < 0.5% |
| Timing | Individual (per channel) |

noise of the preamplifier has also been decreased in the FE analog design in order to maintain the same performances in spite of the larger gain of the electronics. Finally, the signal transmission protocol had to be adapted to the new 40 MHz readout scheme and to the new back-end electronics. As a consequence, the FEBs, 3CU and monitoring boards have been fully redesigned, although maintaining the same format and number of the previous system to minimise the cost of the infrastructure.

### 8.2.1 The front end board

The analog signals from the PMTs are clipped inside the CW base to keep them within the 25 ns LHC bunch crossing window, and then transmitted to the FEB through a 12 m long 50 $\Omega$ coaxial cable. Since the new FEB amplifier has a five times higher gain than the previous FEB, more stringent requirements in terms of noise ($\lesssim$1 ADC counts) were imposed, so that passive termination at the FEB level was not possible. Since each FEB receives the signals from 32 PMTs and an actively terminated input stage is required, an ASIC solution (the ICECAL ASIC [129]) has been implemented for the analog signal processing stage of the FEB (see section 8.2.2). Table 10 summarises the main requirements for the analog FE of the calorimeter system. Except for the PMT current and noise, the other requirements are similar to the ones for the previous calorimeter FE [1, 125, 130].

The calorimeter electronics [1, 125] is based on 246 FEBs, 192 for the ECAL and 54 for the HCAL [130]. This number includes FEBs needed for the measurement of signals of photodiodes of the LED monitoring system (4 FEBs for each subdetector, see section 8.2.6). Each board is connected to 32 PMT outputs. The region of the calorimeter which is covered by a FEB is a rectangle of $4 \times 8$ cells. Each FEB provides digitised data in the form of 32 transverse energy measurements to the BE electronics, and preprocessed information based on the ADC data of the board and on data received from neighbouring boards, to be used by the trigger farm. The digital section of the FEB is based on a 12 bit ADC. Each FEB handles 32 channels, thus the requested output bandwidth is $12 \times 32 \times 40$ Mbit/s. This load is distributed over four optical links driven by GBT chips in wide bus mode.

The location of the front-end electronics is unchanged for the upgrade. The 18 necessary 9U crates will be gathered in racks located on the calorimeter gantries, 14 on the ECAL platform and 4
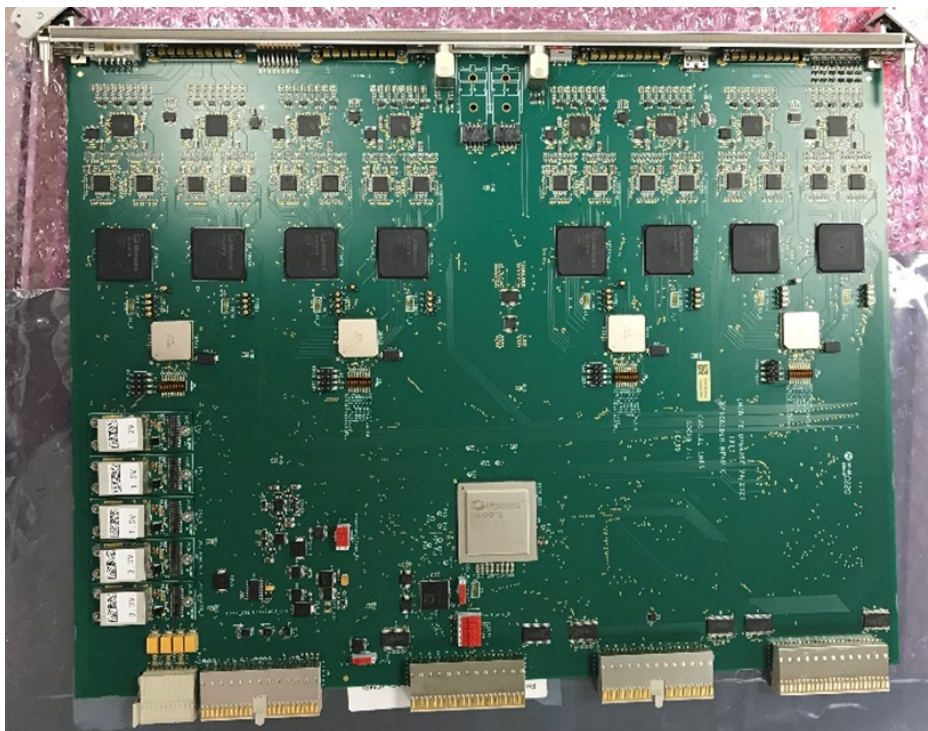
**Figure 79.** Picture of the FEB.

on the HCAL one. The clock, fast and slow control for the front-end electronics will be ensured by 18 3CUs plugged in the central slot of each crate. The 3CU boards are powered directly by the backplane of the crate as for the FEB and are connected to the BE electronics with bi-directional optical links.

There are 6 major sections in the FEB:

- 4 FEm blocks consisting of two ICECAL chips, four dual-ADCs, two FPGAs,[64] called FEm FPGAs in the following, and a GBTx component (driving 4 VTTx emitters), producing $E_T$ measurements for the data stream and calibrated $E_T$ values for the LLT processing;

- the trigger and sequencer FPGA (TrigSeq FPGA[65]), which is used to perform the LLT calculations but also sends additional global information (like e.g. BXID) over the optical links;

- the GBT-SCA that converts the e-links from the GBTx of the 3CU in the fast and slow control signals;

- the block of (de-)serialiser for the exchange of the LLT data between different boards in the same or neighbouring crates;

- the block containing the DC-DC converters and the protection delatchers;

- the VTTx transmitters that receive the data from the four GBTx and send them to the back-end electronics.

---

[64]Microsemi Igloo2 family type M2GL025-1FG484.
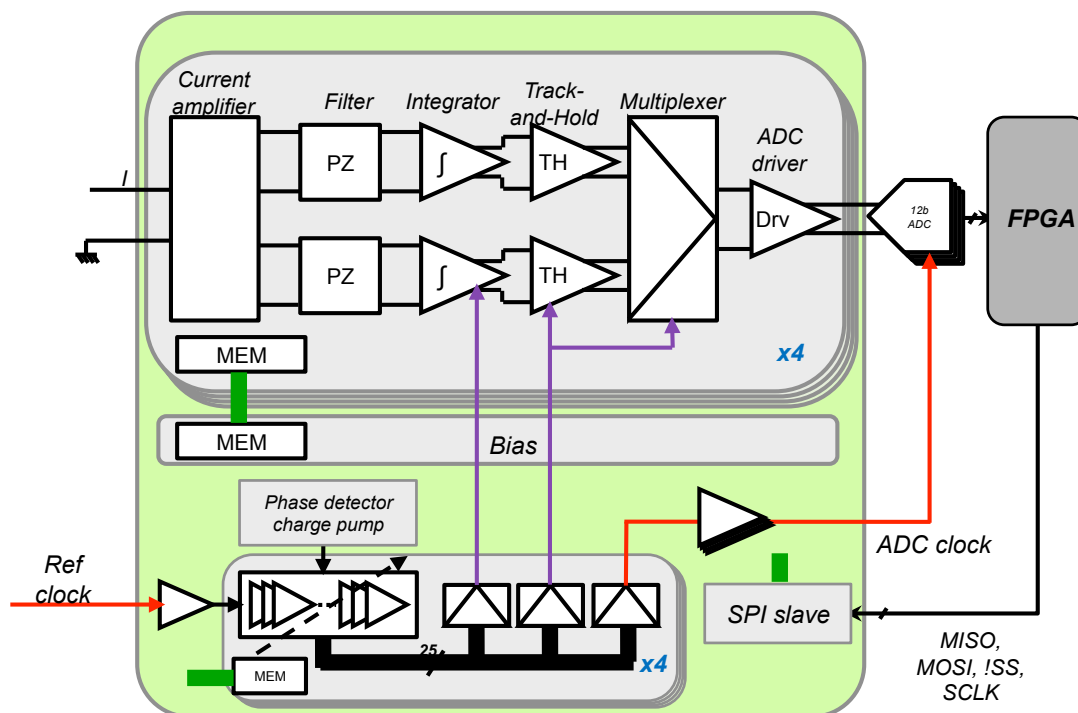[65]Microsemi Igloo2 family type M2GL150-1FC1152.

**Figure 80.** Block diagram of the ICECAL ASIC. Reproduced from [131]. © 2015 IOP Publishing Ltd and Sissa Medialab srl. All rights reserved.

The GBT-SCA ASIC is used to distribute control and monitoring signals to the FEB and perform monitoring operations of detector environmental parameters. It provides various user-configurable interfaces (I2C, SPI, GPIO) used by the 3CU and the FEBs. All registers storing configurations and permanent information are protected with the triple voting technique.

### 8.2.2 The ICECAL ASIC

The ICECAL is a four-channel fully differential amplifier implemented in SiGe BiCMOS 0.35 μm technology. The input stage consists of a current amplifier featuring an active line termination in order to avoid resistor noise. After the preamplifier stage, the signal is sent to two interleaved processing lines running at 20 MHz synchronous with the 40 MHz global clock. Each processing line shapes the signal with a pole zero compensation, integrates the signal, stores the integrated signal in a track-and-hold (T/H) module and finally sends it to the ADC driver through a multiplexer. A block diagram of the ICECAL ASIC is shown in figure 80. The charge integration is performed in the two processing lines by two switched fully differential amplifiers. In a first half-clock cycle, the first amplifier integrates the main part of the PMT signal, which is sampled by the T/H stage while the other amplifier is reset. In the second half cycle, the second amplifier integrates the tail of the signal which is sampled by the T/H stage, while the first amplifier is reset. The fully differential architecture largely eliminates the switching noise intrinsic in this scheme. The samples are then presented through the multiplexer to the ADC driver. The ICECAL receives its clock from the GBTx of the FE module it belongs to. It is configured through the TrigSeq FPGA with the SPI protocol. Among the most important configuration parameters there are the pole zero compensation parameters, the gain of the integrators and the clock phase to be used to integrate the PMT signals.

The four ICECAL channels are connected to two 12-bit dual ADCs.[66] The ADCs require a clock to properly sample the ICECAL output. The clock is produced by the ICECAL and directly injected into the ADC. The clock phase of each of the 32 channels of the FEB can be adjusted independently.

### 8.2.3 FEB digital processing section

As described in section 8.2.1, the 32 FEB input channels are grouped in four front-end blocks, each one consisting of 2 ICECAL, 4 ADCs, 2 FPGAs and one GBTx. The front-end FPGA processing is divided into three distinct functional stages.

The first stage processes the input ADC data, which needs to be resynchronised (each ADC channel has its own clock) and processed to remove the low frequency noise and to subtract the pedestal. A tunable latency is also introduced after the data synchronisation in order to correct for coarse bunch crossing misalignment between channels. The data are then sent to the GBTx for serialisation.

In the second stage, the LLT relevant quantities are calculated by first applying a conversion factor from 12-bit encoding to words of 10 bits. The data is then sent at 80 MHz to the TrigSeq FPGA for further processing or to the neighbouring FEBs for cluster calculations of the FEB border regions of the calorimeter (see section 8.2.4). The low frequency noise subtraction is implemented by subtracting signals previous to the current one, with the assumption that occurrence of consecutive hits in the same channel is extremely rare.

In the third stage, an SPI interface is implemented for the configuration and monitoring of the FPGA. Several FIFOs are used to inject digital patterns, store digital processing results, and store LLT calculations.

The data flow of the board can be tested by spy functionalities using either pattern or signal injection. Patterns of 12-bit words to simulate the ADC input can be injected very early in the processing. Each channel can be configured to receive the standard ADC input or in test pattern mode. The test sequence is controlled by the TrigSeq FPGA so that the patterns can be synchronised between all the FPGAs that are in the test pattern mode. Test sequences can run in single trains of 1024 clock cycles, in loops or in single steps started by a calibration sequence. Additionally, charge injection into the inputs of the amplifiers through small capacitors can also be performed. Charge injection can be made on all channels simultaneously or on any combination of individual channels. Charge injection of pulse pattern tests are fully configurable through the TrigSeq FPGA.

### 8.2.4 Data preprocessing

In order to facilitate the event reconstruction in the calorimeters by the full-software LHCb trigger, and to ease the electron, photon and hadron identification, a simple data preprocessing, which would be otherwise time-consuming in the software trigger, has been implemented in the FEBs. This is historically indicated as low level trigger (LLT), in analogy with the previous L0 hardware stage of the Run 1-2 LHCb trigger scheme, although this preprocessing does not provide any triggering mechanism by itself.

The output of the LLT calculations is fed to the trigger farm concurrently to the raw ADC data. About 7% of the bandwidth allowed by the four optical links of the FEBs is used by this stream without any loss for the ADC data stream.

---

[66]AD9238 from Analog Devices.

Electron and hadron candidates are defined as *clusters*, defined as sums of signals from $2 \times 2$ cells in ECAL and HCAL respectively. Since the Scintillator Pad Detector and the Pre-Shower systems have been dismantled for the LHCb upgrade, there is no way to separate electron and photon candidates at the early FEB stage, and the identification is deferred to the software trigger algorithms. Therefore, for this reason, in the context of the LLT, electron candidates indicate both electrons and photons.

Clusters are built either within a single FEB, if the $2 \times 2$ cells are contained in the 32 FEB channels, or using neighbour FEBs in a crate or FEBs belonging to neighbour crates. The communication between FEBs is ensured by 280 MHz serial links.

The first steps of the computations needed to obtain the electron and hadron candidates in the LLT are realised in the TrigSeq FPGA. The processing consists in a rough calibration of the energy deposited in the calorimeter cells (performed in the front-end FPGA) and in the computation of the $E_T$ of the $2 \times 2$ clusters in each FEB. In addition to the cluster transverse energy, a few more quantities are evaluated:

- the maximum transverse energy measured from the clusters built from $2 \times 2$ cells;

- the address of the cluster giving the largest transverse energy as measured in the previous calculation;

- the total transverse energy from the contributions of the 32 channels handled by the FEB,

- the number of cells on the region covered by the FEB for which the measured transverse energy is larger than a programmed threshold (hit multiplicity).

The results are added to the raw data on the optical links in order to be further processed in the event building farm or to be possibly used as electron, photon or hadron seeds in the software trigger. It is also planned to use the total calorimeter transverse energy as a fast luminosity counter for the experiment.

### 8.2.5 Monitoring and control: 3CU boards

In each calorimeter FE crate (figure 81), the central slot is reserved to the calorimeter control card unit (3CU) board. The main role of the 3CU boards is to distribute the signals from the LHCb control system to the FEBs contained in the crate. The crates are standard 9U VME-like crates with two custom backplanes. The lower one (3U backplane) provides the power supplies, the TFC commands and the clock distribution. The upper one (6U backplane) is reserved to the exchange of signals between the boards and with the other crates. The main role of the 3CU is to receive the GBT frame through the optical link and to extract the information which is needed by the FEBs inside a given crate: the 40 MHz clock, the TFC and ECS commands. The processing of the TFC commands and the control of the 3CU board is performed by a dedicated FPGA.[67]

The 3CU board contains a VTRx and a GBTx chip that can be used to implement multipurpose high speed bidirectional optical links. Logically, the link provides three distinct data paths for timing and trigger control, DAQ and slow control information. In practice, the three logical paths do not need to be physically separated and are merged on a single optical link. The GBTx component of the 3CU is configured to drive 17 GBT-SCAs: the on-board GBT-SCA and the GBT-SCAs of the 16 FEBs that can populate a crate. As mentioned in section 8.2.1, the FEBs are protected by delatchers which detect any current increase that could be due, for example, to a SEL. If such a
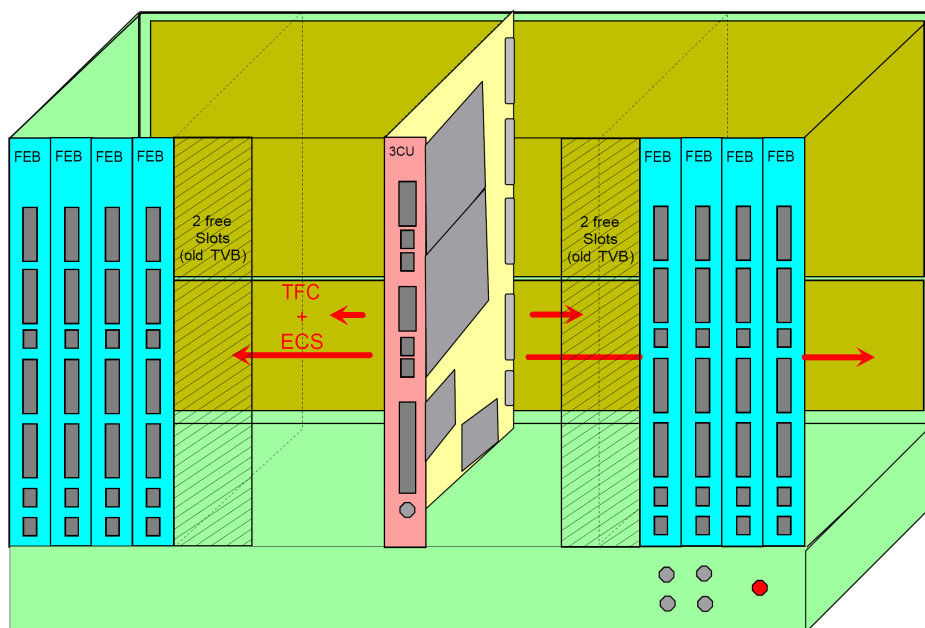
---

[67]Microsemi IGLOO2 family.

**Figure 81.** Schematic of the calorimeter crate.

situation occurs, the current is switched off for a few ms, and then switched on again. The 3CUs monitor the SELs and control the delatching. Additionally, the 3CU FPGA can enable FEB firmware reloading through the ECS if needed.

### 8.2.6 LED monitoring, high voltage and HCAL calibration systems

The LED monitoring and HV systems of LHCb ECAL and HCAL are described in details in [1]. The HV needed to bias the photocathodes and the ten dynodes of each phototube is generated by the CW generator boards, and is controlled by a single analog voltage in a range of 0–5 V applied to the control input of each CW base and generated by a dedicated DAC.

The monitoring of the PMT gain during data taking is performed by measuring the response to an LED flash of constant magnitude injected into the PMT entrance window. The LED flash magnitude is adjustable by applying a control voltage (0–5 V) to the inputs of each LED driver.

The control voltages of CW bases and LED drivers are produced in the common HV-LED DAC board [132], which is interfaced to the ECS. Every HV-LED DAC board provides 200 outputs for PMT HV control and 16 outputs for LED control. The 6016 PMTs of ECAL and 1488 PMTs of HCAL are served by 40 HV-LED DAC boards, 32 for ECAL and 8 for HCAL.

In addition to the LED flash magnitude, the LED monitoring system controls also LED flash timing, which is performed by a dedicated LED trigger signal board (LEDTSB) [133]. Each LEDTSB board provides 64 LED flash delays configurable via ECS. A total of 8 LEDTSB boards are used for ECAL and 2 boards for HCAL.

In order to account for possible instabilities, the LED flash magnitude is independently monitored using PIN photodiodes.[68] A fraction of the LED light is sent to these photodiodes and the corresponding signals are digitised by dedicated standard FEBs (see section 8.2.1). Since the HCAL is placed
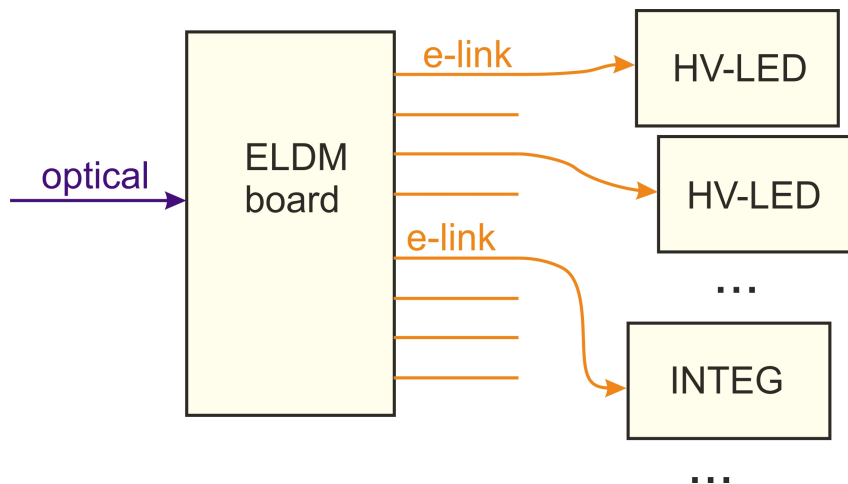
---

[68]Hamamatsu S1223-01.

**Figure 82.** Control board connection scheme with the new GBT protocol.

behind the ECAL, it is difficult to calibrate its energy scale on physics events. Therefore, the absolute HCAL calibration is performed using a ∼ 10 mCi $^{137}$Cs radioactive source that can be moved across every cell [134]. The gain is obtained by measuring the anode current of each PMT under source irradiation. The HCAL calibration is performed during LHC technical stops in dedicated data taking runs. The readout of the PMT anode currents during the HCAL calibration is performed via ECS by means of dedicated boards, named CsCalib, reading out a quarter of HCAL each.

The three types of control boards described above have similar features in their architecture. Namely, they consist of a motherboard bearing the elements executing the main function of the board (DACs for the HV-LED DAC boards, delay chips for LEDTSBs, ADC for the CsCalib boards), on which additional mezzanine boards are installed to implement data transmission and board configuration functions. The upgrade of these boards was dictated by the phasing out of the SPECS bus [135], used in the previous boards, and its replacement by the universal GBT protocol. The upgrade allowed also to replace the obsolete FPGAs of the configuration mezzanine with a new radiation tolerant FPGA.[69]

The transition from SPECS to GBT required also the introduction of an additional board, the e-link distribution module (ELDM) (see figure 82). This board, featuring a VTRx module and a GBTx chip, distributes the information from the optical duplex GBT line to several (up to 10) control boards connected via copper e-links.

### 8.3 Test beam results

To check the upgraded FEB functionality in more realistic conditions a test beam campaign was organised, mainly at the CERN T4-H8 beam line. An ECAL module equipped with photomultiplier tubes and their bases was exposed to a beam of electrons with energy ranging from 20 to 120 GeV. The signal was triggered using two scintillators and sent to a FEB prototype and, in parallel, to a charge integrator and a TDC, which allowed to verify ADC linearity and provided a precise time stamp. The tests were carried out with prototypes of the electronics at different development stages in 2012, 2015 and 2018. They allowed to study the analog electronics in realistic detector conditions (final design PMTs, bases, signal cables, etc.) and to verify that the performance, from the point

---

[69]Microsemi IGLOO2 series.

of view of energy resolution, linearity, noise and signal shaping was within specifications. A test with the final version of the FEB and a simplified version of LHCb DAQ system was made to test the whole acquisition and configuration chain.
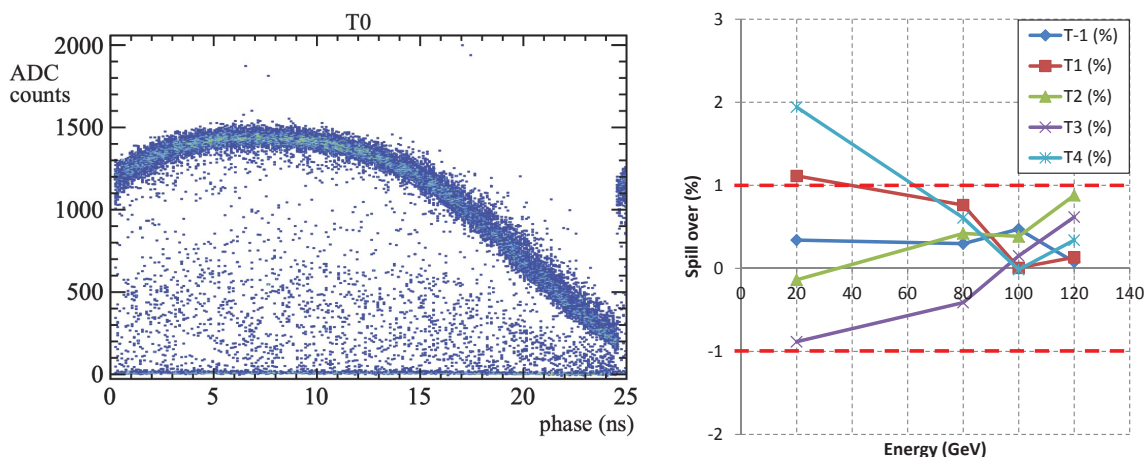


**Figure 83.** FEB test beam results. Left: integrated charge (ADC counts) as a function of hit time phase with respect to the internal 40 MHz clock. Right: spillover measured in different clock cycles as a function of the beam energy.

The electron energy range, from 20 to 120 GeV, allowed a detailed characterisation of the FE amplifier. The electron beam was obtained from the main H8 hadron beam by inserting a lead target into the beam line. Magnetic optics downstream the converter were used to select monochromatic electrons. With a 6 mm thick lead absorber, the electron beam had a pion contamination larger than 60%. Changing to a 12 mm thick absorber improved the electron purity. However, increasing the energy of the beam resulted in a decreased purity and made it difficult to discriminate between low energy electron events, muons and noise.

The CERN SPS beam is not synchronised with the 40 MHz LHC clock, hence hits are recorded by the FEB under test with a random phase with respect to the internal clock. Therefore, the relative phase between each event and the local clock was measured with a TDC. To measure the ICECAL integrator stability, signals were sampled in consecutive 25 ns windows and their time phase with respect to the 40 MHz clock was measured. A selection to separate electrons from pions was applied although it was complicated by the low beam purity. The measured ADC values as a function of the hit relative phase within the 25 ns window are shown in figure 83 (left). A plateau of ∼ 4 ns is clearly visible, with a stability better than 1% showing that if properly time-aligned, the ICECAL is able to integrate the full signal charge within the 25 ns LHC bunch-crossing.

By choosing events where the signal is fully integrated within 25 ns (*T0 cycle*), it is possible to look into previous and following clock cycles to measure the residual spillover charge not integrated within the main clock cycle. The result of this test is shown in figure 83 (right), where the spillover is measured in the previous (*T-1*) and following (*T1-T4*) clock cycles with respect to T0, at different beam energies. The measured spillover is within ±1% as required, with a maximal deviation observed at low energy, where however the separation between electrons and pions was very difficult. Note that the spillover can be positive as well as negative due to the signal shaping used to reduce the pulse width to equal or less than 25 ns inside the ICECAL ASIC.
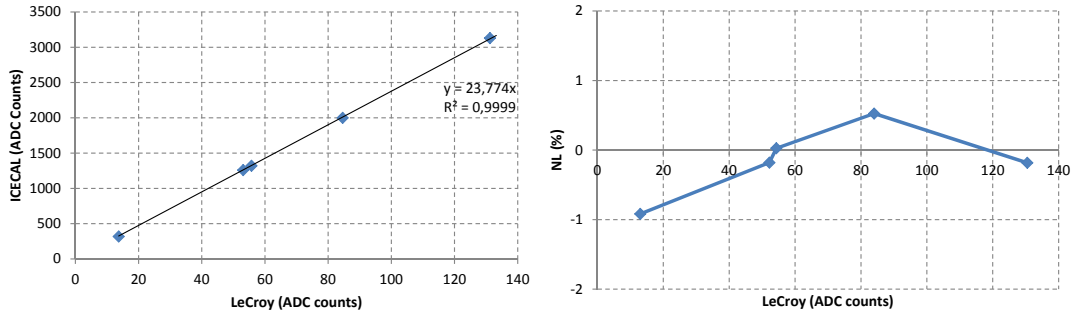
**Figure 84.** Left: energy values measured with the FEB prototype at a beam test with respect to the reference charge integrator values. Right: the nonlinearity deviation is shown to be less than 1%.
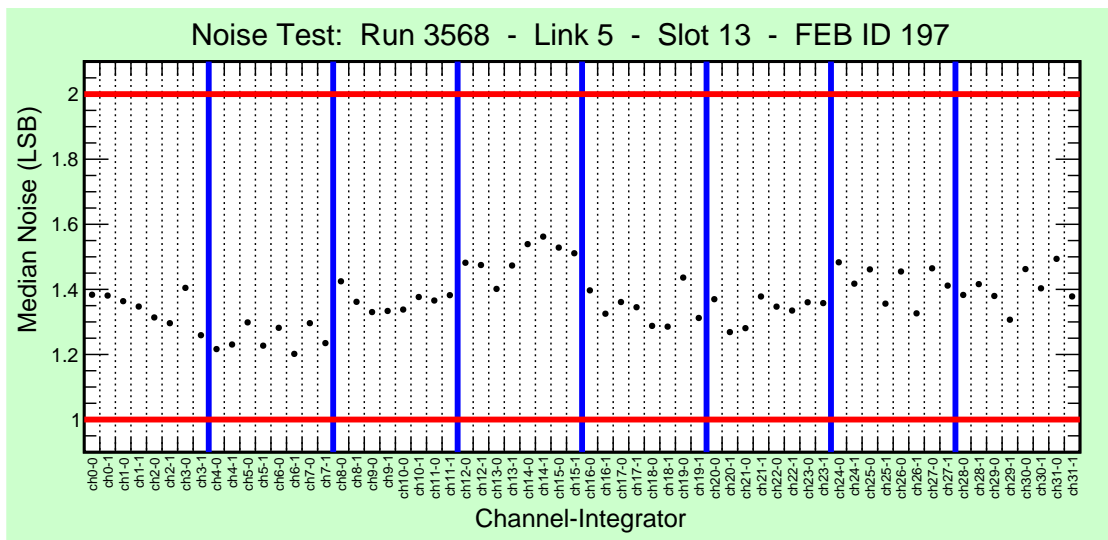


**Figure 85.** Noise of all channels of a FEB, measured in laboratory conditions.

The key parameters of the analog circuit were checked and showed a behaviour within specifications. Linearity was checked to be better than 1% by comparing the charge integrator and the FEB prototype readings for different electron energies. The results are shown in figure 84. As shown before, the integrator stability was better than 1% in 4 ns. The spillover was also tested for different PMTs, PMT biasing voltage, PMT bases and signal cables and was found to be stably below 1%.

The noise was checked in different conditions both in the laboratory and at the test beam. Results with and without pedestal and low frequency noise subtraction were also obtained. Test bench measurements without PMT biasing yielded a noise level of 1.3 and 1.6 ADC counts before and after pedestal subtraction, respectively, slightly above specifications. In the beam tests, the noise behaviour varied significantly, ranging from 1.6 to 3.4 ADC counts, probably due to setup changes and grounding imperfections. To obtain a more realistic estimation, a measurements was performed directly in LHCb, connecting to a channel at the crate level with proper grounding. After pedestal subtraction, the noise was measured to be about 1.4 LSB, in agreement with test bench results as shown in figure 85, a level that was considered acceptable.

# 9 Muon system

## 9.1 Overview

The LHCb muon detector [1, 136–138] has been successfully operated during LHC Run 1 and Run 2 with an excellent performance [139, 140].

The LHCb muon system is composed of four stations M2 to M5 comprising 1104 multi-wire proportional chambers (MWPC) for a total area of $385\,m^2$. The LHCb Run 1-2 muon system additionally included a station M1 located upstream of the calorimeters and comprising 12 gas electron multiplier (GEM)s in the innermost region and 264 MWPCs. Station M1 was utilised in the hardware L0 trigger and thus is no longer needed in the upgraded system. Its supporting structure has been maintained to host the SciFi Tracker neutron shielding as described in section 2.4.2. Each station is composed of two mechanically independent halves, Side A and Side C. The four stations M2 to M5, located downstream of the calorimeter system, are equipped with MWPCs and interleaved with 80 cm thick iron absorbers to filter low energy particles. Each station is divided into four regions, R1 to R4, of increasing area moving from the central beam axis outwards. The area and the segmentation of the four regions scale in such a way to uniformly distribute the particle flux and the channel occupancy across each station. The MWPCs are made up of four independent layers (or *gaps*), each consisting of anode wires between two cathode planes, to achieve a high efficiency and a high redundancy. The FE electronics host an amplifier-shaper-discriminator stage implemented in a dedicated ASIC as well as a digital section that allows time alignment of the signals and logical combinations of readout channels in so-called *logical channels*. The FE electronics was designed to be radiation tolerant up to 100 kGy which is expected to be adequate also at the new running conditions. Thus, the current FE electronics is kept unchanged. Digitised signals from the FE electronics were originally sent at 40 MHz rate to the L0 trigger but recorded for further processing only at a maximum rate of 1 MHz. To comply with the new LHCb readout scheme a complete overhaul of the readout electronics has been carried out representing the main upgrade of the muon system [108].

The monitoring and control electronics have also been completely redesigned to comply with the new 40 MHz readout rate and the new experiment's DAQ and control systems. Despite the significant changes required, the new electronics have been designed to be backward-compatible with the original architecture in order to minimise the cost and allowing to reuse of the original crates, cabling, and power supplies. The new electronics is presented in section 9.2, while section 9.3 describes the muon-system-specific processing implemented in the back-end electronics system. Section 9.4 presents the new ECS software developed to monitor and control the muon system.

The high particle rates expected in the upgraded muon system required the introduction of specific strategies to mitigate the otherwise unacceptable level of induced inefficiency. These are discussed in section 9.5.

The MWPCs have been operated for eight years and are expected to be left in place for the whole lifetime of the experiment. To secure the smooth operation of the muon detector, the number of needed spares has been estimated based on the operational experience during Run 1 and Run 2. In addition to the ones already available from past productions, 54 new chambers have been built (30 for M5R2 and M5R4 regions and 24 for M2R3, M2R4, M3R4, and M4R2 regions). Studies on expected long term operation of the chambers have been performed and are summarised in section 9.6.
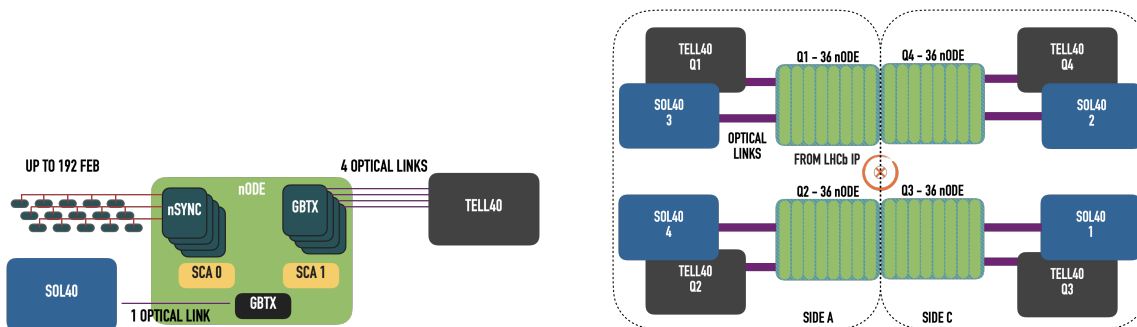
**Figure 86.** nODE communication scheme.

## 9.2 Electronics

The muon system electronics is designed to convert, format and transmit downstream the analog signals extracted from the detector. As introduced in section 9.1, the FE electronics is fully compliant with the upgraded experiment running conditions and will be therefore maintained. The MWPC signals are digitised by the front-end CARDIAC boards [141] which host two kind of chips, two CERN and Rio current amplifiers (CARIOCAs) [142] and a diagnostic, time adjustment and logics (DIALOG) [143], both implemented in a IBM 250 nm technology combined with a specific layout technique to be radiation tolerant up to 100 kGy. The CARIOCA is an eight-channel ASIC that implements a current mode amplifier, a three stage discriminator and a LVDS output driver. In each CARDIAC one DIALOG chip receives sixteen input channels from two CARIOCA boards. Each input channel can be configured with a programmable delay in steps of about 1.6 ns realised by a voltage-controlled delay line. Each channel can be propagated or disabled by means of a masking circuit. In order to optimise the number of readout channels, the CARIOCA outputs are combined within the CARDIAC into *logical channels* by means of selectable logic functions suitably chosen according to the position of the corresponding detector channels. The DIALOG provides also individual threshold settings for the CARIOCA discriminators and test pulses. The DIALOGs are fully accessible and configurable via an I2C interface.

During operation in Run 1-2, signals from CARDIAC boards, which are in some cases further combined by the Intermediate Boards, were received by the off-detector electronics (ODE) boards which provided a time stamp through a 4-bit TDC, formatted and shipped the data to the BE electronics. The TDC was implemented in a dedicated ASIC named SYNC which also allowed coarse grained time alignment of signals. Data was shipped by the ODE to the L0 trigger at a rate of 40 MHz and, upon L0 trigger positive decision, to the DAQ boards at a maximum rate of 1 MHz. This scheme can no longer be maintained in the upgraded readout system and the ODE has been redesigned and upgraded into the new off-detector electronics (nODE) boards. The TDC has also been fully redesigned and upgraded into the new SYNC ASIC (nSYNC). In the nODE boards, signals are synchronised with the master LHCb clock, compressed, formatted and sent to the downstream TELL40 readout boards via high-speed optical links. The nODE boards are also connected with the SOL40 boards to interface and distribute the TFC and ECS information. Figure 86 shows the communication scheme and the architecture of the new off-detector electronics system comprising 144 nODEs. Each quadrant will be served by 36 nODEs, installed in five crates.

The muon control and monitoring system [144], formerly provided by the service boards (SBs), has also been fully redesigned to comply with the new ECS and TFC systems and has been implemented
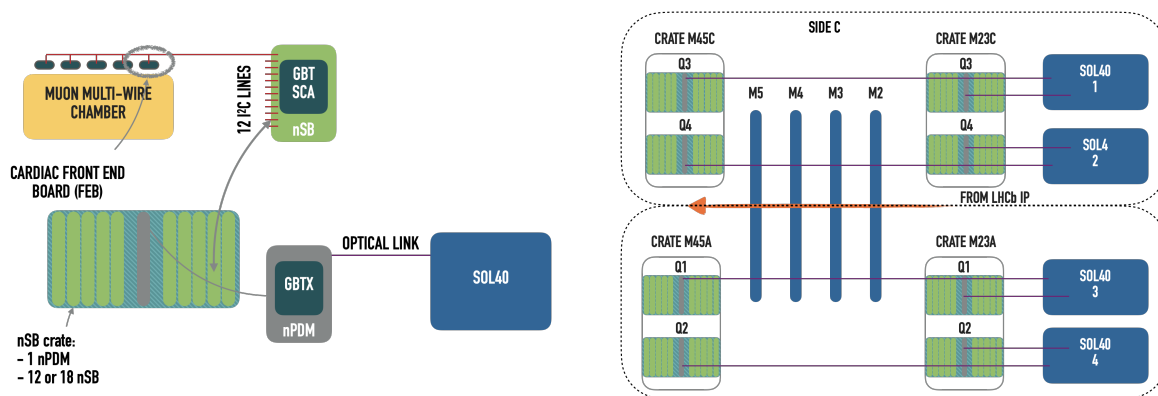
**Figure 87.** FE electronics communication scheme.

in the new service board system (nSBS). The nSBS is comprised of the new pulse distribution modules (nPDMs), the new service boards (nSBs), and the new custom backplanes (nCBs). Figure 87 shows the communication scheme and the architecture of the nSBS. Eight crates comprising 12 or 18 nSBs and one nPDM are installed for a total of 8 nPDM and 120 nSBs.

### 9.2.1 The nODE

The nODE conceptual design is similar to the previous ODE board and is intended to be backward-compatible with the previous architecture in terms of board format, crate occupancy and cabling.

The nODE receives up to 192 input logical channels which are processed by the four onboard radiation-tolerant nSYNC chips which provide clock synchronisation, bunch-crossing alignment, time measurements, histogram capability and buffering (see section 9.2.2). In order to optimise the output bandwidth and minimise the number of long-distance optical links the nODE transmits data frames with packed nonzero-suppressed hit maps of the corresponding input channels and zero-suppressed time stamps from the nSYNCs, thus transferring the decoding part to the BE TELL40 boards. Details of data formatting are discussed in section 9.2.2.

The logical block diagram of the nODE is shown in figure 88. Board functionalities can be grouped in three main logical blocks: TFC and ECS management block with related distribution sections, FE and data electronics block, and power converter and distribution block. The board layout, shown in figure 89, is divided in three functional sections: a master section implementing management and distribution of TFC and ECS information and two slave sections, UP and DW, implementing FE and data functionalities. Each stage includes dedicated power management circuits.

The master section is based on one GBTx chip, two GBT-SCA chips and one VTRx optical transceiver. The GBTx is configured in transceiver forward error correction mode and generates the reference clock synchronous with the master LHCb clock, which is then distributed to slave sections. The ECS/TFC path relies on two GBT-SCAs to control the slave GBTxs and nSYNCs. The TFC information of the nODE implements a subset of seven of the standard TFC commands of the LHCb experiment. The encoded TFC command is replicated four times in the data frame and distributed to each nSYNC via seven e-links [146]. The ECS interface is used to configure the electronics components of the board, to monitor their status and to download the time histograms.

The slave sections are comprised of a total of up to 192 input channels, four nSYNC chips, four GBTx chips and two VTTx optical transmitters. All such elements are arranged into two identical
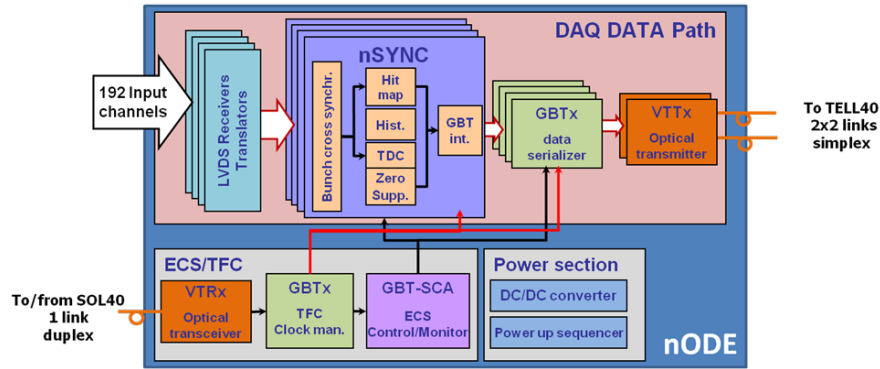
**Figure 88.** nODE block diagram. Reproduced from [145]. © CERN 2014 for the benefit of the LHCb collaboration. CC BY 3.0.
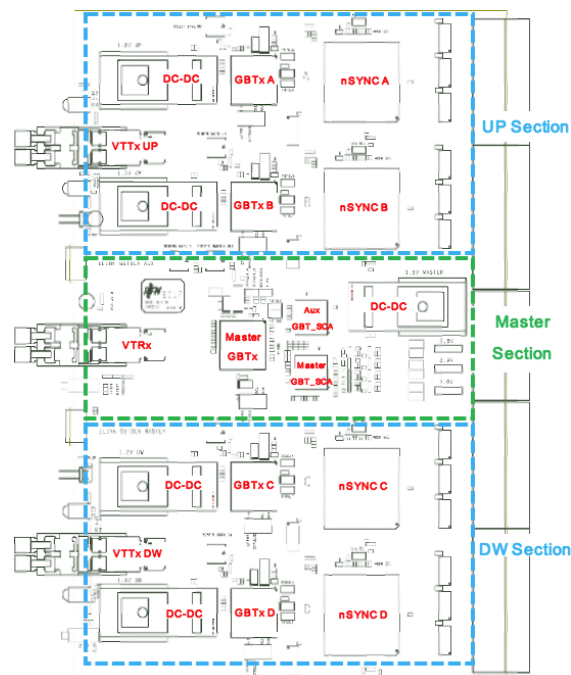


**Figure 89.** nODE board layout.

sections, UP and DW. Each section is further subdivided into two identical data paths comprising one nSYNC with 48 input channels, one GBTx and one channel of the VTTx module. The slave GBTxs are configured as simple transmitters in widebus mode with fixed header and fixed frame scheme. The nSYNC data frame is received by the GBTx via fourteen input e-links at 320 Mbit/s. The 112 bits data field of GBTx data frame is composed of a Header Field of 12 bits encoding the BXID, followed by four dummy bits and by the nSYNC data frame.

### 9.2.2 The nSYNC

The main building block of the new readout electronics of the LHCb muon system is the nSYNC chip [147], a radiation tolerant custom ASIC developed in UMC 130 nm technology as an evolution of the SYNC chip [148], used in the old ODE boards. The nSYNC architecture is composed of several
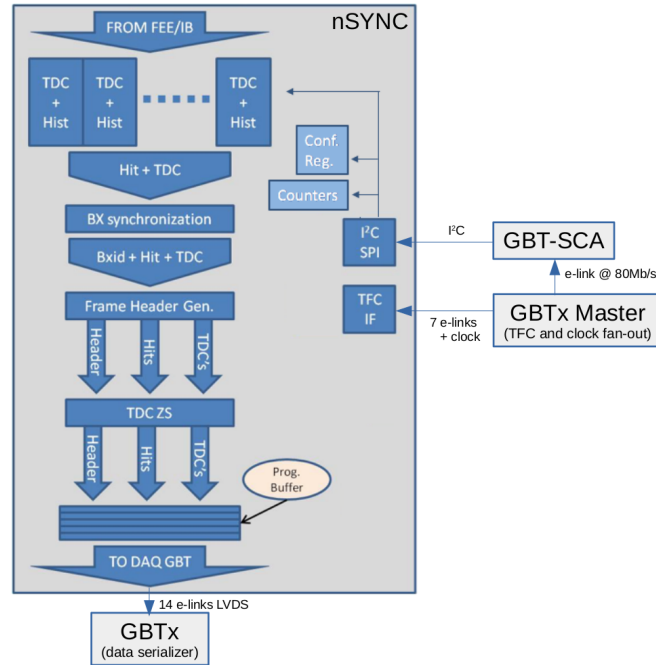
**Figure 90.** Schematic view of the nSYNC architecture and its interface with the GBT chipset. Reprinted from [147], Copyright (2019), with permission from Elsevier.

functional blocks, schematically shown in figure 90. The main purpose of the nSYNC is to integrate all the required functionalities for the upgrade of the readout system, such as clock synchronisation, bunch crossing alignment, hit map production, time measurements, histogram capability and buffers. The nSYNC handles also the zero-suppression algorithm for the time measurements and the interfaces to the DAQ and TFC/ECS systems. The nSYNC receives the digital signals coming from the muon chambers, through 48 LVDS input channels, and synchronises them with respect to the BXID. Concurrently, the phase of the arriving signal is measured by a TDC (one for each channel) with respect to the LHC 40 MHz master clock. This information is crucial for the time alignment of the whole muon detector. The BXID information, the hit maps and the TDC counts are then combined to build a frame that, after the TDC data zero-suppression, is transmitted to the DAQ through the GBT output interface.

The TDC is composed of a fully digital patented digitally controlled oscillator (DCO) [149], which produces a clock signal based on an input digital word. The TDC works with a nominal resolution of 1.56 ns. This is obtained by dividing the 25 ns LHC clock cycle into 16 slices, allowing to encode the TDC time stamp in 4-bit words. Lower (8 slices) or higher (32 slices) resolution can also be configured. Incoming signals trigger the DCO clock, whose periods are then counted until both the counter and the DCO are stopped by the arrival of the rising edge of the master clock. The phase measurement corresponds to the counter output and is stored in an output buffer. Since the DCO is composed of a digital delay chain, a systematic error is continuously accumulated during the measurements. Therefore a dithering system is implemented in order to add or remove a unit to the digital input word of the DCO, thus inverting the systematic error and preventing its accumulation. At nominal resolution, the output data of each TDC consist of a binary flag (corresponding to the hit/no-hit information) and a 4-bit-wide word with the measured phase, if a hit is present. Each nSYNC channel is also equipped with a histogram facility, comprised of 16 counters, in order to monitor the measured phases.

The TDC output data are sent in parallel through a pipeline in order to align all the hits belonging to the same bunch crossing. This allows to uniquely associate the 12 bit-wide BXID number to all the data. Finally the extended output frame is built by combining a header, the 48 bit-wide hit map and the TDC data. While the hit map is always sent nonzero-suppressed, the TDC data are instead subject to a zero-suppression algorithm: only the first nonempty block of TDC data are added to the frame, up to a maximum of 12 TDC measurements at nominal resolution. The TDC data address decoding is therefore deferred to the BE electronics, allowing an optimisation of the output bandwidth. The last eight bits of the frame are dedicated to the Hamming code, a feature that can be disabled to increase the TDC occupancy.

Output communication between nSYNC and GBTx chips is ensured by 14 LVDS links at a transmission rate of 320 Mbit/s. An intermediate asynchronous FIFO is implemented to interface the two clock domains. The high frequency clock is generated internally using a PLL with the 40 MHz master clock as reference. Moreover, an 8-steps programmable output pipeline allows the correct alignment from the receiver side. The 40 MHz clock and all other synchronous fast commands sent by the TFC system are received through the GBTx master chip. Asynchronous slow commands, like configuration commands, are received from the GBT-SCA chip through a I2C bus, as explained in section 9.2.1.

The radiation level at the muon readout crate location is not critical, being about 40 times smaller than in the detector region. For 10 years of LHCb upgrade operation the expected total ionising dose is 130 Gy, with an expected fluence of $2 \times 10^{12}$ 1 MeV $n_{\text{eq}}/\text{cm}^2$ [150]. Nevertheless, to ensure proper ASIC operation, two radiation hardness techniques have been implemented in the nSYNC design to mitigate the SEU rate: triple modular redundancy to protect the most critical registers, such as the first input stage of TFC commands and the configuration registers, and the Hamming code with the corresponding error detection and correction logic, to protect all the internal counters, buffers and output FIFO. The radiation hardness of the nSYNC has been verified with several tests, in particular using a 60 MeV proton beam and X-ray irradiation [151, 152], up to $2 \times 10^{12}$ 1 MeV $n_{\text{eq}}/\text{cm}^2$ and a total ionising dose of 1.3 kGy. The nSYNC showed an excellent performance under radiation with no failure or SEU behaviour after an accumulated total ionising dose ten times larger than the one expected for 10 years of LHCb upgrade operations. The chip current consumption and the internal clock jitter are expected to increase by less than 5% for the same time period. The SEU cross section per bit has been also measured to be $(0.53 \pm 0.04) \times 10^{-13} \text{ cm}^2$, corresponding to an expected SEU rate of less than 0.1 events per day for the whole muon detector for the most important registers, thus having a negligible impact on the overall muon detector performance and efficiency.

### 9.2.3 The new service board system

The main purpose of the nSBS is the management and distribution of TFC and ECS information to the CARDIAC FE boards. The nSBS consists of a crate hosting a master board, called the new pulse distribution module (nPDM), and up to 20 slave boards, called new service boards (nSB), which replaces the previous system. The nPDM communicates with the nSBs via e-links on a custom backplane, called new custom backplane (nCB). Each nSB is interfaced with up to 96 CARDIAC boards via 12 serial links. The block diagram of the nSBS is shown in figure 91. Control signals from the ECS and TFC are routed to two separate paths in order to reduce the complexity and the costs of the whole system. The functional representation of each path is shown in figure 92.
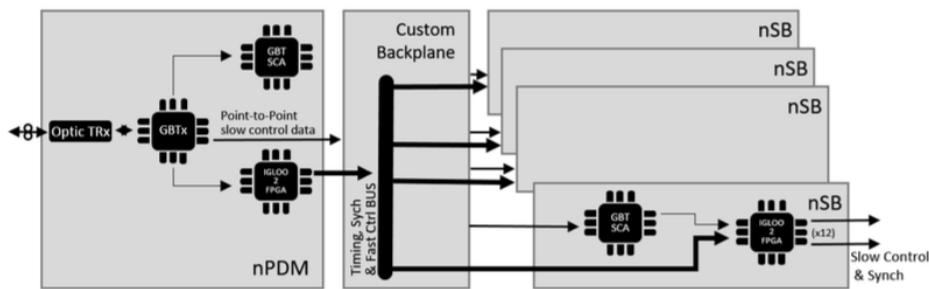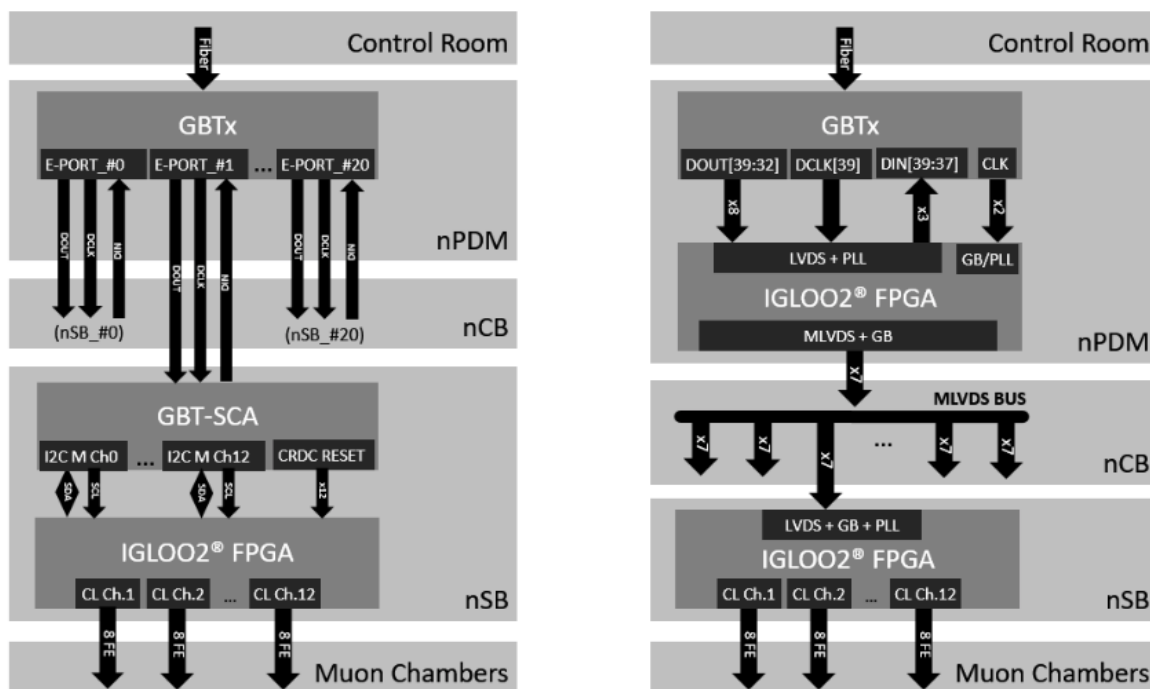
**Figure 91.** nSBS block diagram.



**Figure 92.** Left: functional representation of the nSBS ECS path. Right: functional representation of the nSBS TFC path. Reproduced from [144], with permission from Springer Nature.

The nPDM main components are a GBTx chip, a GBT-SCA chip, an IGLOO2 FPGA and a VTRx optical transceiver. A GBT-SCA and an IGLOO2 FPGA constitute also the main components of each nSB. The nPDM GBTx is controlled through the VTRx by the interface board SOL40. The GBTx is configured in transceiver mode with forward error correction and it generates the reference clock for the whole nSBS, synchronous with the master LHCb clock. The TFC interface relies on eight e-links at 40 Mbit/s from the same group and propagates the TFC commands to the nPDM FPGA, which in turn distributes them to the nSB FPGA via backplane lines. The local nPDM ECS interface uses the EC field of the frame to communicate with the nPDM SCA which controls and configures the nPDM FPGA. Twenty more e-links from the remaining groups of the nPDM GBTx working at 80 Mbit/s, routed through the backplane, are used to propagate directly the ECS information to the nSB SCAs of up to 20 nSBs. The nSB SCAs are used to configure and control the nSB FPGA which in turn distributes all the timing and control information to the FE CARDIAC boards via I2C links.

### 9.3 The muon readout board specific processing

As discussed in section 10.2.1, the TELL40 readout boards are provided with a generic firmware framework where subdetector specific firmware can be plugged in. The specific muon system firmware scheme is driven by bandwidth considerations and aims at optimising the number of long distance optical links and TELL40 boards.

The maximum theoretical bandwidth for one PCIe link is 64 Gbit/s. The PCIe protocol encapsulation and the DMA processing limit in fact the bandwidth to 54 Gbit/s but a conservative limit of 50 Gbit/s was imposed to each of the two output PCIe links of TELL40 boards [153]. In the GBT wide bus mode 96 bits of data can be transmitted per each link at the maximum trigger rate of 40 MHz, resulting in an input rate of 3.84 Gbit/s per link. Assuming that all the input data are transferred unmodified to the output of the TELL40 board using its generic output data format, this would limit the maximum number of links per board to 22. In addition, imposing a 70% limit on the FPGA resources of the TELL40, the maximum number of links that can be handled at maximum input rate is limited to 32. However, as a low occupancy is expected, in particular in the outer regions and in the most downstream stations of the muon detector, a specific data processing block with a zero-suppression algorithm has been implemented in the muon TELL40, in order to minimise their number and optimise the number of long distance optical links.

To estimate the expected output bandwidth several minimum bias events acquired in Run 2 at $\mathcal{L} = 3.7 \times 10^{32}\,\mathrm{cm^{-2}\,s^{-1}}$ have been superimposed. The output bandwidth per event obtained using the output data format described in section 9.2.2 is reported in table 11. The average rate is well below the maximum allowed of 50 Gbit/s.

**Table 11.** Maximum output bandwidth ( Gbit/s) per PCIe interface in the muon system stations at two different luminosity values when zero-suppression is applied; for comparison, also the output rate with no zero-suppression is reported.

| Station | # TELL40 | output rate ( Gbit/s) $\mathcal{L} = 2 \times 10^{33}\,\mathrm{cm^{-2}\,s^{-1}}$ | output rate ( Gbit/s) $\mathcal{L} = 4 \times 10^{33}\,\mathrm{cm^{-2}\,s^{-1}}$ | output rate ( Gbit/s) no zero-suppression |
|---------|----------|----------|----------|----------|
| M2 | 10 | 22 | 27 | 54 |
| M3 | 4 | 24 | 33 | 61 |
| M4 | 4 | 13 | 18 | 35 |
| M5 | 4 | 18 | 25 | 42 |

The number of readout boards per muon station and links per board are listed in table 12.

**Table 12.** Number of TELL40 per station and of input links per board.

| station | M2 | M3 | M4 | M5 |
|---------|----|----|----|----|
| # TELL40 | 10 | 4 | 4 | 4 |
| input links/TELL40 | 28 | 32 | 18 | 22 |

The muon specific data processing block consists mainly of a zero-suppression algorithm which decodes TDC data when available and formats the output data taking care of the different number of optical links connected to different TELL40s. The zero-suppression procedure is fully configurable via ECS, both at nSYNC level and in the TELL40. This implementation allows different options
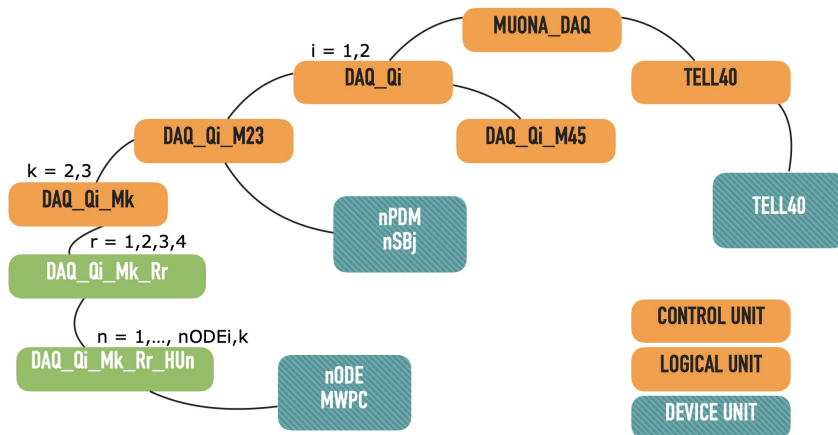
**Figure 93.** Finite state machine hierarchy scheme (example for the Side A).

to be used for boards covering regions of the detector with different occupancy levels if needed, although the baseline is to use the same settings everywhere.

## 9.4 Monitoring and Control

Control and configuration of the muon system is performed within the LHCb ECS system framework (section 10.4). Muon-specific software to monitor and control the muon system consists of a finite state machine organised in the hierarchical structure reported in figure 93, through which all the system can be configured and operated.

Moving down the tree from the top-level node, the system is partitioned into two nodes, Side A and Side C, serving the two mechanically and functionally independent halves of the muon detector. From each node it is possible to access specific subelements of the DAQ, high voltage, and slow-control ECS sector. The DAQ node is organised in quadrants, stations and regions coherently with the detector layout and is further subdivided in FE electronics, nSBS and off-detector electronics domains. The slow-control ECS sector is further subdivided to separately control and monitor low voltage, gas and cooling systems.

## 9.5 Mitigation of the high rate inefficiency

Consequences of the high rates expected at the upgrade running conditions on the performance of the muon system have been extensively studied [108]. From the analysis of data taken at $\mathcal{L} = 1 \times 10^{33} \, \mathrm{cm}^{-2} \, \mathrm{s}^{-1}$ during a test run in 2012 [154], it has been demonstrated that no space charge effects are expected at this luminosity. The only expected degradation of detector efficiency due to the increased rate is caused by the dead time induced by the FE electronics. The particle flux in the innermost region of station M2 is expected to be very high, resulting in a very unevenly distributed efficiency drop of about 7% on average in the region closest to the beam pipe. To reduce the inefficiencies an additional shielding has been installed around the beam-pipe before M2. In particular, as shown in figure 94, a tungsten *plug* shielding has been inserted in place of the removed innermost cells of the HCAL (see section 8.1.3). A rate reduction of about 30% in region R1 of M2 has been estimated by simulation with a consequent maximum rate estimated to be about 600 kHz/ cm $^2$. The high rate will induce inefficiencies in the FE electronics resulting from both the CARIOCA dead-time,
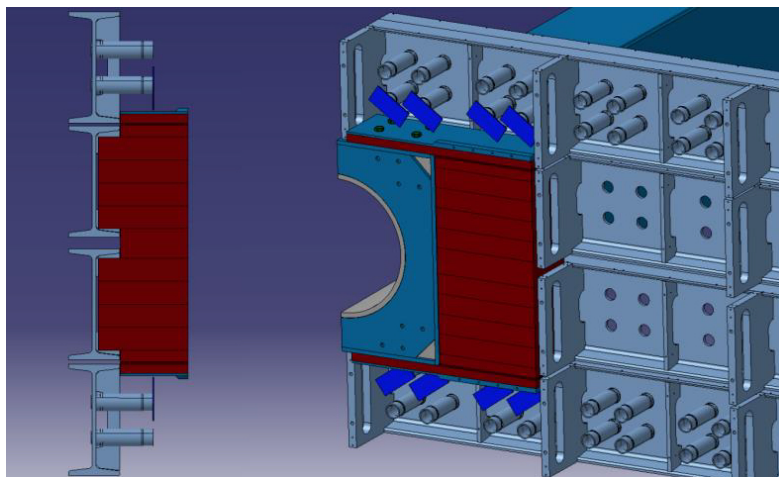
**Figure 94.** Mechanical drawing of the tungsten shielding around the beam pipe.

in the range of 70 ns to 100 ns, and the DIALOG signal formation time, which was successfully reduced from 28 ns at the beginning of Run 1 to 12 ns at the end of Run 2. The CARIOCA component is estimated to introduce the dominant contribution to the inefficiency, amounting to ~8% in M2R1 and ~4% in M3R1 and M2R2, respectively.

New pad detectors with increased granularity have been proposed [155] for these innermost regions to mitigate inefficiency from dead time. A first prototype for an M2R2 chamber, see figure 95 (left), has been designed and constructed in 2016. A prototype for region R1 of M2 and M3 stations, see figure 95 (right), was also constructed in 2020. Both prototypes successfully passed the necessary tests and were accepted for mass production. The new MWPCs with full pad readout could be installed during Run 3 or LS3, depending on the readiness of the chambers. To mitigate the effect of DIALOG dead time, the granularity of the logical channels has been increased by replacing some Intermediate Boards [136] with nODE boards in regions of particularly high rates, namely regions R2, R3 and R4 of station M2 (right behind the HCAL) and R4 of station M5 (where the rates are dominated by interactions with LHC materials placed behind the LHCb detector). The replacement of Intermediate Boards is expected to reduce the inefficiency by about 20%, which will be reduced by an additional ~ 40% when the three innermost regions will be equipped with new pad chambers.



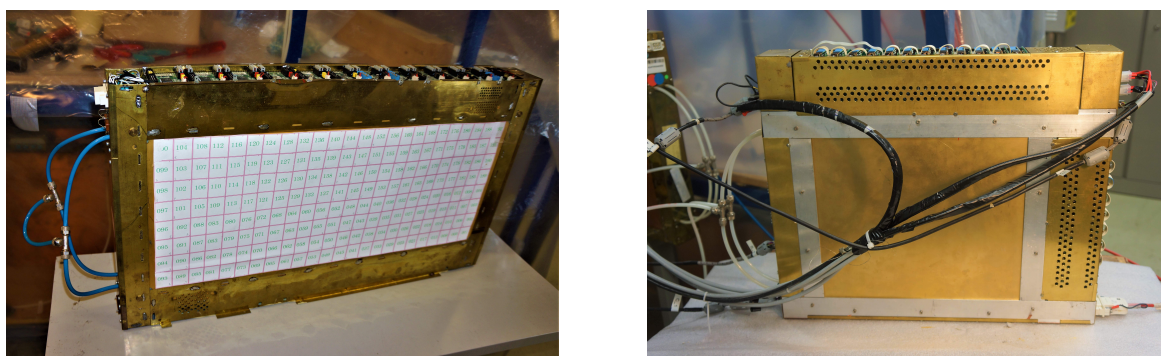**Figure 95.** Left: M2R2 new pad chamber prototype. Right: M2R1 new pad chamber prototype.

## 9.6 Long-term operation of the MWPCs

It is foreseen that the muon system MWPCs will be operated over the whole LHCb life time, for a total integrated luminosity of about $50\,\mathrm{fb}^{-1}$. Special care was therefore taken in limiting as much as possible severe damage to the wire chambers and strategies were developed to fix and recover MWPCs showing operational problems, in particular due to sparks or high currents.

In nine years of operation in a high radiation environment, the LHCb muon detectors did not show visible gain reduction or any other apparent performance deterioration. However, during this period, about 19% of the chambers were affected by sudden appearance of high currents in at least one of their wire layers, with a total of about 100 wire planes per year suffering HV trips during LHC operations. In order to prevent future severe damage and to ensure long-term operability of the MWPCs, a detailed review of the status of the chambers has been conducted during LS2. The findings of these systematic studies, described in details in ref. [156], are summarised in this section.

High currents observed in some of the chambers during operation were found to be originating from localised areas usually located on the cathodes. This has been verified by direct inspection of broken MWPCs, where carbonised or stained spots where observed on the cathodes, in some cases in correspondence of broken or damaged wires. These currents were found to be due to sustained discharges that, in addition to locally damaging the wires and the cathodes, generate noise and may lead to HV trips when exceeding the set threshold, causing temporary detection inefficiency. Inspection of the behaviour of these currents suggests that they are likely due to a Malter-like effect [157]. Even if the latter is often associated to ageing, the analysis of the HV trip distribution in time and in different regions of the detector, as well as direct inspection of damaged chambers, indicates that ageing due to prolonged irradiation is not the underlying cause and that most of the HV trips are connected with imperfections existing since the chamber construction.

A method for a noninvasive recovery on site has been developed and applied individually to all problematic gaps, consisting in HV training cycles carried out in presence of colliding beams and with MWPCs working with the standard gas mixture. The training procedure lasts typically many weeks (two months on average), before the affected wire plane is recovered and restored back to normal operation. A typical example of the appearance of a self-sustained current and of the recovery procedure during LHCb operation with beams is shown in figure 96. The four-fold redundancy of the muon system MWPCs and the HV training on site allowed to keep the whole muon detector continuously close to 100% efficiency for almost a decade. Moreover the recovery procedure developed and refined over the years has been shown to be effective, and less than 1% of the chambers had to be replaced because of HV trips in nine years of operation. The percentage of gaps treated with this method in the past and showing repeated high current problems decreased with time and was measured to be about 10% during the last two years of LHC operation.

In order to make the training procedure faster and even more efficient, a method for accelerated recovery has also been investigated during LS1. A set of four MWPCs removed from the apparatus because of persistent high currents, underwent the standard HV training procedure, but $\sim 2\%$ oxygen was added to the default gas mixture. Figure 97 shows the results of this procedure obtained for one of the wire planes. While no current decrease is seen after more than six hours of HV training with the standard gas mixture, consistently with observations described above, HV training in presence of oxygen is much faster and the discharge current is observed to drop down to zero in about four hours after a few HV cycles. After the training with oxygen, all of the four MWPCs were fully
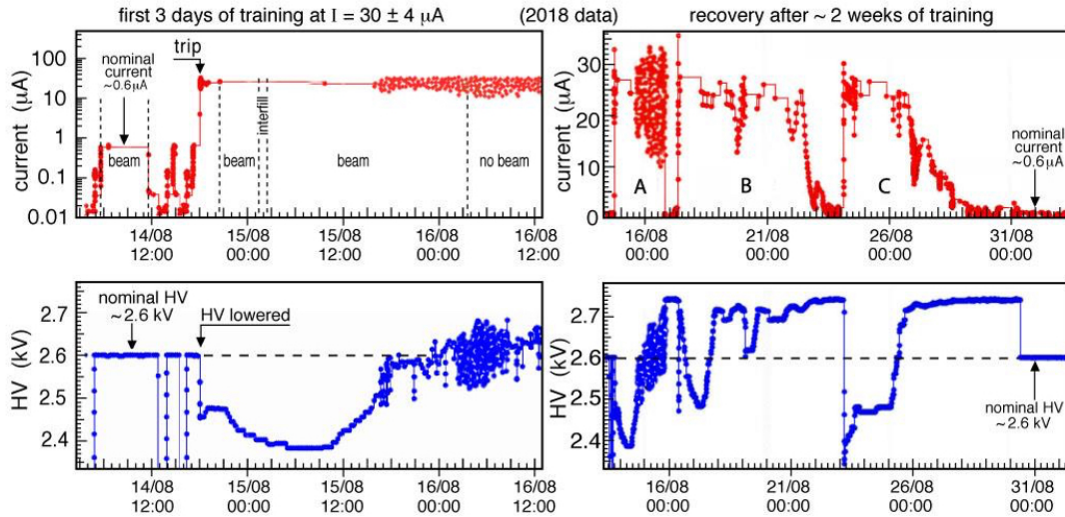
**Figure 96.** Typical recovery procedure for a gap (data are from a chamber in region M5R3). The plots on the left show (top) the current and (bottom) the HV setting during a period of about three days around the first appearance of the HV trip and the subsequent start of HV training. The plots on the right show (top) the current and (bottom) the HV setting during the full recovery procedure, which lasted about two weeks. The nominal HV setting for this gap is 2600 V. in normal conditions, the average current in presence of colliding beams is about 0.6 μA. Reproduced from [156]. © 2019 CERN. CC BY 3.0.
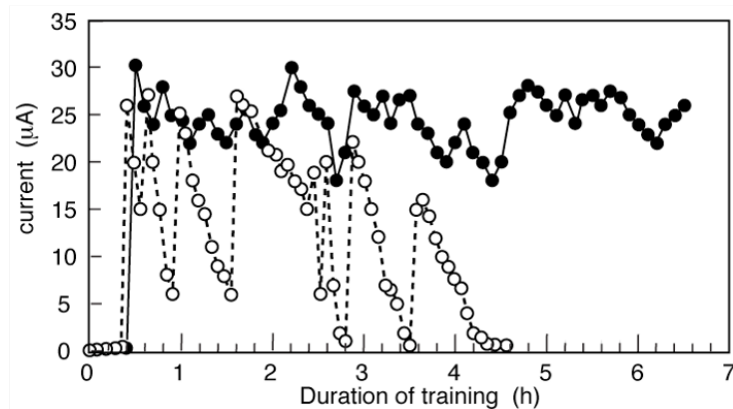


**Figure 97.** Current in the MWPC during the Malter-effect recovery training: default mixture (full circles) is compared with a mixture containing ∼ 2% of oxygen (open circles). Reproduced from [156]. © 2019 CERN. CC BY 3.0.

recovered and installed back on the apparatus being still operational today. The noninvasive character of the recovery techniques discussed in this section makes them an important ingredient for the long-term operation of the muon system. The results obtained through the fast recovery suggest that in the future a small amount of oxygen could be added to the working gas mixture during the LHC winter stops, either for targeted recovery interventions, or for conditioning while exposing the chambers to a high intensity radioactive source. A permanent use of a small amount of oxygen during detector operation is also being considered.

# 10 Online system

## 10.1 System architecture

The upgraded online system consists of a continuation and evolution of the successful experiment control system from Run 1 and Run 2 [158], a new timing and fast signal control for clock, synchronous and asynchronous commands distribution for the trigger-free readout, and a significantly increased data acquisition system. Further hardware and software subsystems are alignment and calibration frameworks, the online monitoring, storage and the infrastructure to operate all these systems and the event-filter farm. The system hardware consists of a single type of a powerful custom-made, flexible FPGA board and commercial off-the-shelf hardware.

Most of the system naming scheme is derived from the previous LHCb online system [158], often adapted to recall the 40 MHz master readout clock.

## 10.2 Data acquisition system

LHCb uses a synchronous readout where for *every bunch-crossing* all front-end elements participating in the data taking send data, possibly after zero-suppression. Therefore, from the point of view of the FE electronics, the LHCb online system can be justifiably described as "trigger-free", despite the fact that subsequent event selection mechanisms implemented in software are called "high level triggers". In the LHCb online scheme, the readout elements are grouped in *partitions*. A partition is a set of online resources, which can be read out, monitored and controlled independently, for example a part of a subdetector, a full subdetector or a group of subdetectors. Multiple partitions can run simultaneously, which is a very powerful tool for commissioning and testing. Partitioning is implemented by both the TFC and the ECS.

The LHCb data acquisition system is shown in figure 98. It consists of a farm of event builder (EB) servers hosting the back-end receiver boards (TELL40 boards) and the graphics processing units (GPUs) running the HLT1 application. Data processed by the EB and the HLT1 are then sent to the HLT2 for further processing and final storage.

Data are transported over half-duplex multimode optical fibres from the detector underground level through a service shaft up to a data-centre on the LHCb site surface. The radiation hard versatile link (VL) protocol is used, with most subdetectors using the gigabit transceiver (GBT) protocol at the OSI-layer 2 [160].

Specific to LHCb is that the links dedicated to data transmission are used in half-duplex mode; there is only a single fibre carrying data from the front-end to TELL40 boards, while control, configuration and monitoring are out of band on dedicated connections. The links dedicated to the control of the readout electronics, either to FE or BE, are in full-duplex mode; TFC and ECS are transmitted to the FE electronics sharing the payload on the links whereas only ECS are received back, still utilising the same optical duplex links.

The DAQ links are received by the TELL40 boards described below, which then push the data into the memory of the EB servers. After event-building the completed events are passed to the HLT1 running on GPUs installed in the EB servers. Accepted events are stored on the HLT1 buffer storage and then read by the HLT2 processes for final selection. Accepted events are consolidated into files and sent to permanent storage.

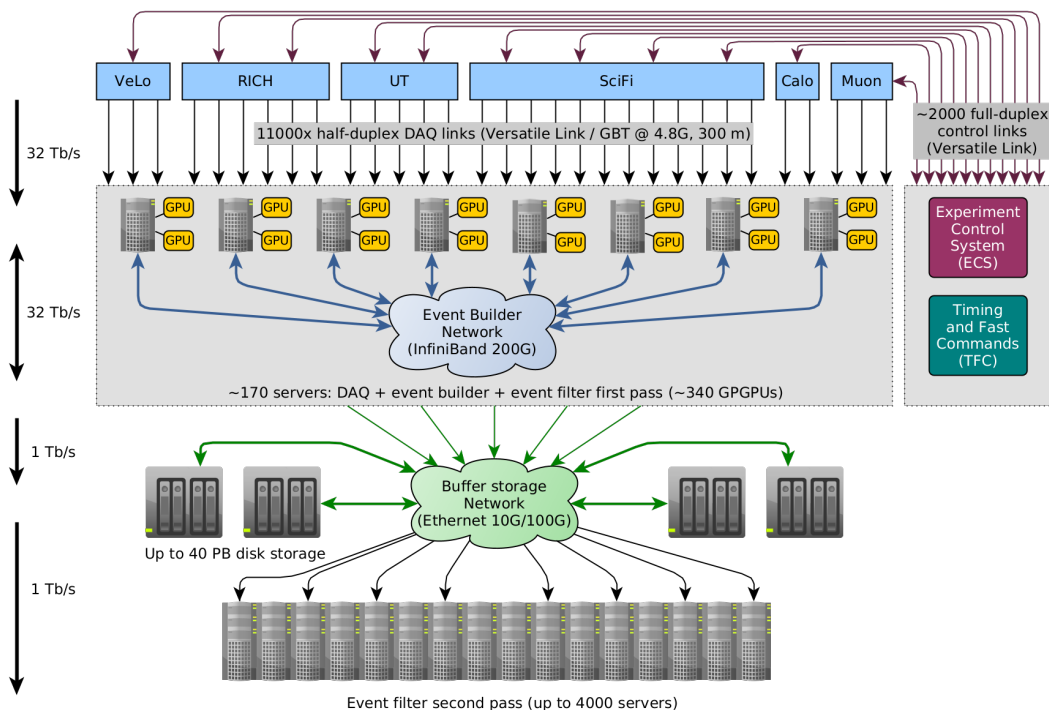The various components of the DAQ system are described in more detail in the following sections.

**Figure 98.** Upgraded LHCb online system. All system components are connected to the ECS shown on the right, although these connections are not shown in the figure for clarity. Reproduced from [159]. CC BY 4.0.

### 10.2.1  Common readout and control PCIe40 board

LHCb has chosen to use a single, custom-made board for data acquisition, slow control and fast, synchronous and asynchronous control. The board concatenates data, transforming input streams based on a custom protocol into an output stream based on a standard protocol used in the data centre. The input stream is the VL running the GBT protocol on top for all detectors except the VELO, which uses the GWT protocol. The output stream is PCIe Gen 3 with a bandwidth of 120 Gbit/s. This board has a PCIe form-factor (three-quarter length, standard height, dual slot) and it is designed according to PCIe specifications. Hinting at the board interface and at the 40 MHz experiment's driving clock, the board has been called PCIe40.

This board, shown in figure 99, is based on an FPGA[70] featuring 72 serialiser/deserialiser (SerDes) out of which 16 are used for the host-interface to the PC. The front-end electronics can be connected via up to 48 bidirectional optical links. The bandwidth of each link is up to 10 Gbit/s in both directions. The FE facing links are directly connected to optical transceiver modules (MiniPods) to achieve high density.[71] Physically, these links appear as eight MPO-12 connectors on the board. In addition, two serial links are reserved for the timing and fast control with a pluggable optical module. The cooling of the FPGA and MiniPods is obtained through a custom heat sink and relies on the air flow of the PC server. Components of the board require six different voltages which are provided by a daughter card connected to the main 12 V of the PC server. The card supplies up to 45 A on 0.9 V to the core of the FPGA, in a sustained regime. Finally, a series of phase-locked loops are

---

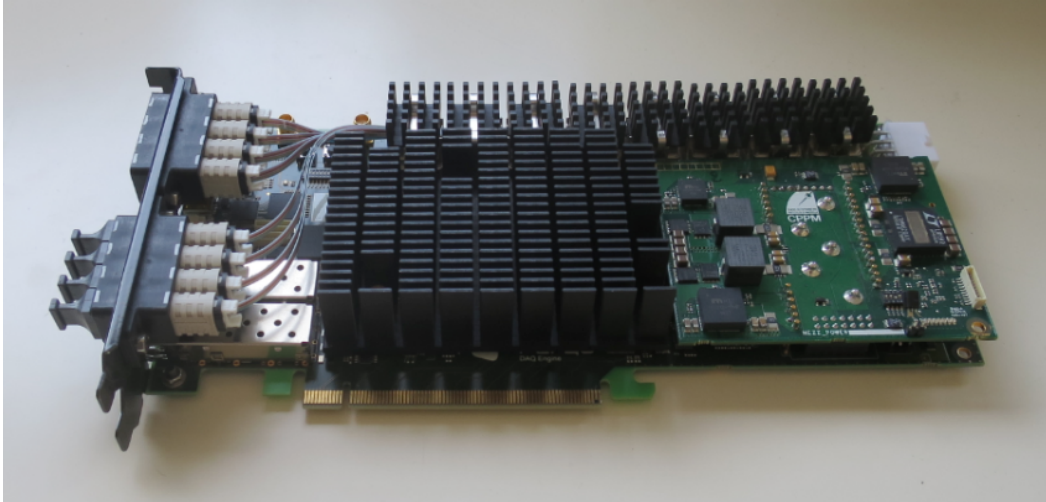[70]Intel™ Arria10.
[71]Broadcom™ MiniPod.

**Figure 99.** View of the PCIe40 board.

implemented to distribute clocks at 40 and 240 MHz with low jitter. Their phases are controlled at the level of 200 ps with respect to the LHC main clock.

When used for data acquisition, the PCIe40 interfaces the FE electronics with the EB. In that configuration, only the receiver direction is used and no transmitter MiniPods are fitted in order to save cost. In this configuration, the board is called TELL40. When used for slow and synchronous controls, the board is called SOL40 and all bidirectional links are active. Typically, 3 (8) TELL40 (SOL40) are installed per PC server respectively.

Different functionalities are obtained by charging different firmware versions in the FPGA, with at least one per subdetector (see section 10.2.2). A special firmware version is also implemented for the master-module of the synchronous readout, the SODIN board, which acts as readout supervisor.

### 10.2.2   Firmware

By reconfiguring the onboard FPGA with dedicated firmware, the PCIe40 can be used to serve very different roles within the upgraded LHCb experiment: fast control, clock and slow control distribution and data acquisition. The firmware of the board contains an interface layer code, common to all the boards, whose aim is to interface the hardware with the user firmware using common blocks, for example GBT decoder/encoder and PCIe IP cores. The actual user firmware defines the flavour of the board and its functionality within the upgraded readout architecture. This considerably reduces the number of components to be developed and optimises personpower as well as effort on firmware design and maintenance. The environment to develop the firmware for each configuration of the boards is common across the entire LHCb experiment, with only the user code being exclusive. The different board flavours are:

- SODIN (readout supervisor),

- SOL40 (main interface),

- TELL40 (data acquisition).

The SOL40 board serves three main purposes:

- interface all the readout boards to the SODIN by fanning-out the synchronous timing and trigger information and fanning-in throttle information [161];

- interface all the FE electronics to the SODIN by relaying the clock, timing and command information onto fibres towards the FE electronics [69];

- relay the ECS information by interfacing the slow control protocols at the FE electronics [162, 163].

The user code dedicated to the functionality of the readout of events from the trigger-free architecture faces considerable challenges. Events arrive from the FE to the TELL40 boards asynchronously across all input links due to the variable latency in compression/zero-suppression mechanisms. Thus the code must be able to handle a large spread in time between fragments of the same event. The readout code of the board must be able to decode the data frames from the FE, realign them according to their BXID, build an event packet and send it to the DAQ network. Figure 100 illustrates the architecture of the TELL40 firmware. Low level interfaces and features which can be common to all the subdetectors, like decoding or data time alignment, are developed centrally for the collaboration (coloured in blue in figure 100). Each subdetector's user will only develop a few specific blocks which will be fitted within the common architecture (coloured in red in figure 100). This optimises personpower and reduces the complexity of the development allowing for faster integration, commissioning and maintenance.
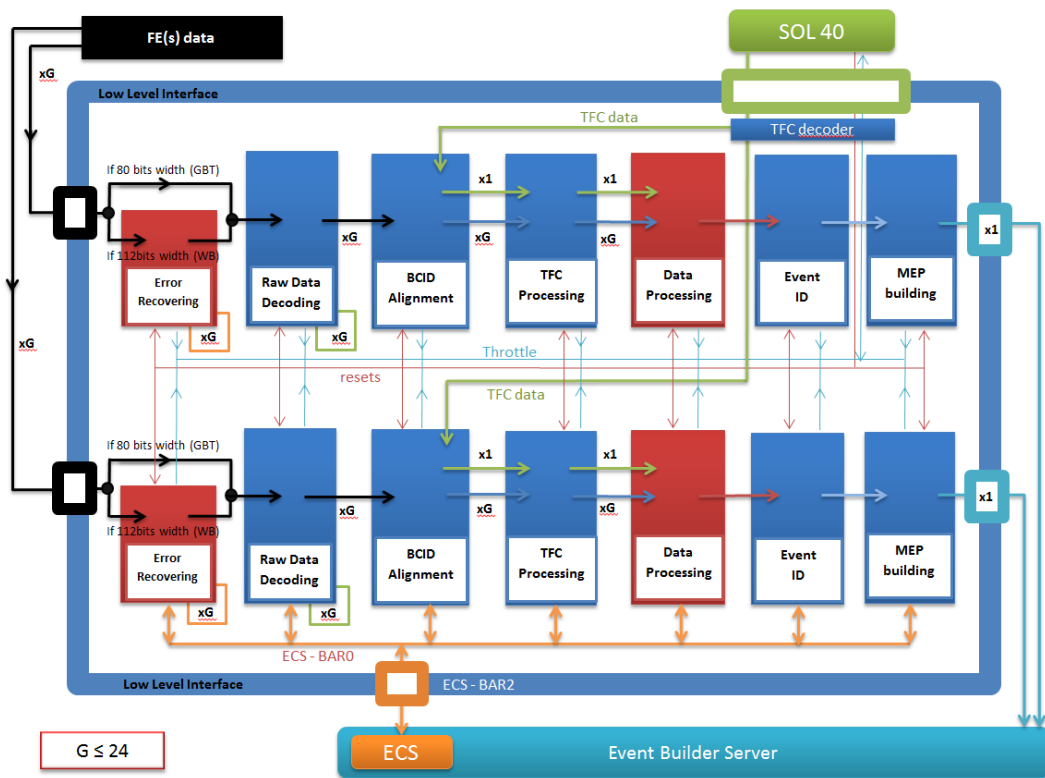


**Figure 100.** TELL40 firmware architecture. Common blocks in blue; specific subdetector blocks in red.

The LHCb online group provides engineering support and development to the different LHCb subdetector groups, mainly in relation to readout board firmware, low-level software for front-end

configuration and data acquisition, and control system components. A proper firmware/software framework based around continuous integration has been put in place to mitigate the very heterogeneous nature of the different subsystems. This approach reduces the time and effort required to detect, reproduce, and correct integration issues during the development cycle of the firmware.

### 10.2.3 Firmware and software frameworks and continuous integration

In addition to the online developer team, other members of the LHCb collaboration from several institutes participate actively in the firmware and software development. The benefits of automated integration testing are evident in such a geographically distributed structure, since issues like failing builds or failing tests can be flagged and reported automatically to the participants regardless of their location or timezone. Source control is organised using GIT and GIT repositories managed through the GITLAB infrastructure provided by the Information Technology department at CERN.

The readout FPGA firmware is organised into logically separate GIT submodules, allowing maintainers of the corresponding pieces of functionality to version their code independently. A dedicated top-level repository tracks the state of the firmware submodules, and provides a unified location from which different FPGA firmware codes can be built. When developers are ready to integrate their changes, they submit a merge request which automatically triggers a dedicated GITLAB continuous integration pipeline. This pipeline executes Questa RTL simulations according to several predefined firmware configurations and, if successful, executes the Quartus FPGA synthesis.

Given the complexity of modern high-capacity FPGAs, producing all required permutations requires of the order of 100 hours of computation. A series of optimisation steps are thus applied. Distributing synthesis jobs across a small computing cluster reduces this turnaround to a single day. As a further optimisation, the process keeps the result of each synthesis job in an internal cache and tracks changes in all firmware components across successive invocations of the same job configurations. This automated dependency tracking is used to selectively trigger only pipeline jobs whose dependencies have changed. For example, changes to a specific subdetector implementation, or to board-specific logic, will result in repeating only the jobs associated to that specific subdetector, or that specific board, respectively.

Both software and firmware are automatically packaged in RPM format and published to dedicated RPM repositories for distribution, as is customary on all Linux installations deployed at CERN.

### 10.2.4 Event-building

To perform the event selection, the full-software LHCb trigger requires the complete event information from all the subdetectors. Consequently, event-building, the assembly of all pieces of data belonging to the same bunch-crossing, is done for every collision of non-empty bunches at a 40 MHz rate.[72]

The EB system consists of a farm of servers hosting the TELL40 boards receiving data from the subdetectors and the GPUs running the HLT1 application. As each server receives only the data from the subdetectors connected to the corresponding TELL40s, all the EB nodes must be interconnected with a high performance network able to transmit the full information to the node which is in charge for the full event assembly. The architecture of this high-speed network is the

---

[72]The information for each bunch filling is transmitted synchronously to the TELL40s by the TFC system; data recorded during empty crossings are normally dropped.

main driver for the overall system design, which has been optimised such that the costs for handling the expected input event rate are minimised.

For any high-speed network the number and reach of the links are the main cost-drivers. Therefore, the network links are used bidirectionally so that their number is kept comparatively small. In addition, the system has been designed to be very compact and physically installed in only two modules of LHCb's data centre (described in section 2.4.4), keeping the physical length of the network links at a minimum.

The 200 Gbit/s high dynamic range (HDR) InfiniBand technology has been chosen for the EB network implementation, which not only offers comfortable data throughput but allows to fully exploit modern PCIe interfaces.

The PCIe40 cards needed to read out all the LHCb subdetectors are hosted in 162 servers with up to three cards per server and avoiding server sharing among different subdetectors. Each server is connected through two HDR ports to the EB network. Each server acts in turn as data-source and data-sink in the event-building process, where cyclically every node acts as full event builder (sink) and receives data from all other servers (sources). In order to achieve optimal network performance, data from several thousand bunch-crossings are packed together and these multievent fragment packets are treated as unit data blocks in the event building.

When a builder unit has received data from all sources, the sources are reordered to facilitate subsequent processing. Completed events are stored in a shared memory buffer and are handed over to the HLT1 selection process. The HLT1 application pushes the required event data to a GPU installed in each event-builder server. The events selected by HLT1 are then sent via a separate 10G/100G Ethernet network to temporary storage, from where they are accessed by the alignment and calibration processes and by the second-stage selection (HLT2).

## 10.3 Timing and fast control

The TFC is responsible for controlling and distributing clock, timing and trigger information, synchronous and asynchronous commands to the entire readout system as described in [68]. The system must maintain synchronisation across the readout architecture, provide the mechanisms for special monitoring triggers and manage the dispatching of the events to the trigger farm. It regulates the transmission of events through the entire readout chain taking into account throttles from the readout boards, the LHC filling scheme, calibration procedures and physics decisions if any, while ensuring a coherent data taking acquisition across all elements in the readout architecture. The specifications, functionalities and the full details of the system are described in ref. [161].

Generally, the signals generated and propagated by the TFC system to the entire readout system are:

- the LHC reference clock at 40 MHz, that is the master clock of all the electronics synchronised to the master clock of the LHC accelerator;

- commands to synchronously control the processing of events in the readout board [164] or front-end electronics;

- calibration and specific subdetector commands for the detector electronics.

In addition, FE electronics configuration is generated by the ECS and relayed by the TFC system to the FE boards. At the hardware level the TFC system is implemented using the common PCIe40 card. Details of the TFC aspects in this card are discussed in ref. [165].

### 10.3.1 Functionalities of the TFC system

The main functionalities of the TFC system are:

- *Readout control*: control of the entire readout system is made by one of the TFC Masters in the pool. The control of the readout implies controlling the trigger rate, balancing the load of events at the processing farm and balancing the occupancy of buffers in the electronics. The TFC system auto-generates internal triggers for calibration and monitoring purposes in a programmable way, as well as a full set of commands in order to keep the system synchronous. The details of the specifications for the FE and BE are described in detail in ref. [69].

- *Event description*: a data bank, containing information about the identity of an event as well as the trigger source, is transmitted by the central TFC Master to the farm for each event as part of the event data.

- *Partitioning*: this is achieved by instantiating a set of independent TFC Masters in SODIN, each of which may be invoked for local subdetector activities or used to run the whole of LHCb in a global data taking. An internal programmable switch fabric allows routing of the information to the desired destination.

- *Coarse and fine time alignment*: the TFC distribution network [166] transmits a clock to the readout electronics with a known phase, kept stable at the level of about 50 ps, and a very low jitter ($< 5$ ps). The latency of the distributed information is fully controlled and kept constant. Local alignment at the front-end of the individual TFC links is required to ensure synchronisation of the experiment. This alignment relies on the synchronous reset commands together with Bunch Identifier and Event Identifier checks.

- *Run statistics*: information about the trigger rates, run dead-time, number of events accepted, type of events accepted, bunch currents, luminosity and load of buffers is stored in a database to allow retrieving run statistics and information per run or per LHC fill.

### 10.3.2 TFC architecture, timing and control distribution

The upgraded TFC architecture and data flow are represented in figure 101. The readout supervisor SODIN is the TFC Master, responsible for generating the necessary information and commands to be interfaced to the LHC clock distribution system. The subdetector readout electronics comprised of FE and BE boards are connected to the SODIN via a network of bi-directional optical links via multiple SOL40 interface boards and passive optical splitters. These connections define the partition granularity and their topology defines a partition, controlled by the TFC to run any ensemble of subdetectors simultaneously.

These connections utilise different technology, protocol and bandwidth according to their destination and purpose:

- the connections between the SOL40 and the FE boards use the GBT protocol with Forward-Error correction enabled, at 4.8 Gbit/s (including the bits needed for error correction), for fast and slow control distribution to the FE and for slow control back from the FE;
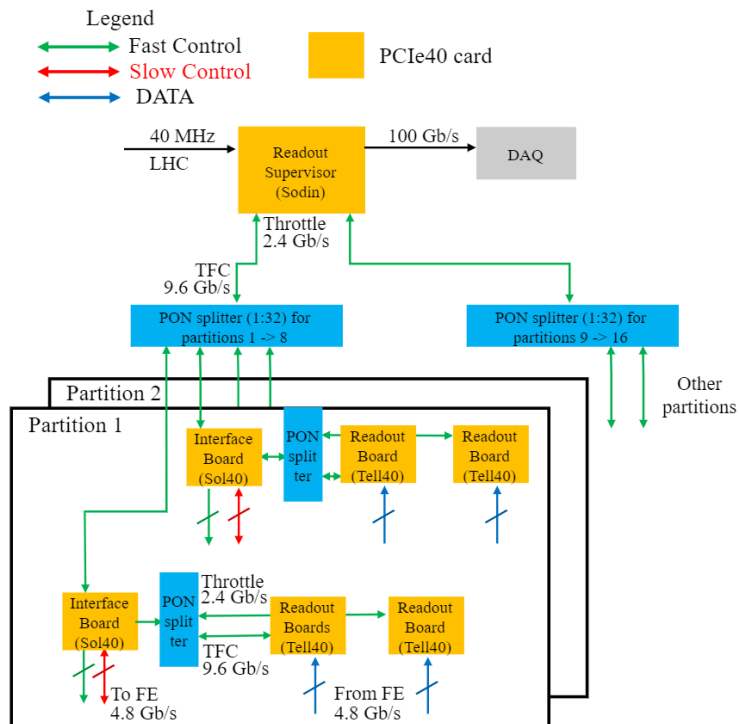
**Figure 101.** Logical architecture of the upgrade TFC system.

- the connections between the SODIN and the SOL40 as well as between the SOL40 and the TELL40s utilise Passive Optical Network technology (PON) using PON splitters to reach multiple destinations, running at 9.6 Gbit/s in the downstream direction for clock and command distribution and at 2.4 Gbit/s in the upstream direction for throttle information;

- the connection to the LHC clock is via an electrical interface (LVDS) from the LHCb clock reception and distribution system [164]. The SODIN receives the 40 MHz LHC clock as well as the 11.245 kHz revolution clock (commonly referred to as the orbit pulse) and uses this to synchronise the event information to the LHC collisions;

- lastly, SODIN is interfaced to the rest of the DAQ by its PCIe interface, transmitting roughly 100 Gbit/s of event bank information.

LHCb employs a unique mechanism to merge the TFC and ECS information in the same duplex link to the front-end, via their interface to the SOL40 boards. The TFC information is packed into the GBT verbatim at 40 MHz, while the ECS information can span multiple words and so make use of the available bandwidth as needed. At the front-end each subdetector developed its own architecture in order to decode this information and use it locally. The logical scheme of the merging is shown in figure 102, as a generic example. In fact, the SOL40 boards may be cascaded and configured differently to support different requirements in terms of number of links and bandwidth as well as supporting different architectures at the FE boards. This is entirely done via a mixture of configuration registers in the firmware as well as configuration parameters at compile time, while keeping the firmware development common to all subdetectors.
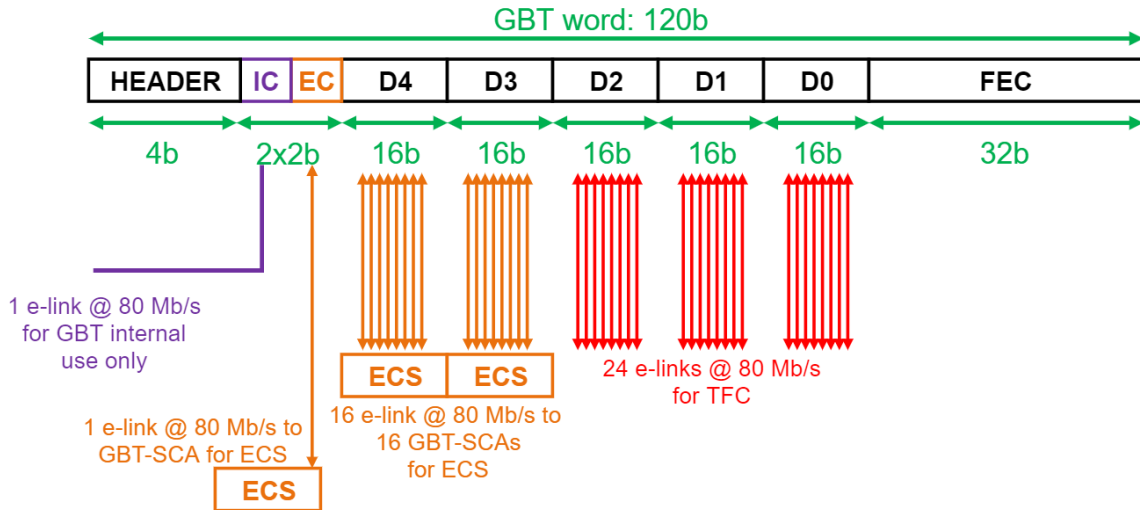
**Figure 102.** Schematic view of the packing mechanism to merge TFC and ECS information on the same GBT links towards the FE electronics. GBT words are subdivided into small e-links.

## 10.4 Experiment control system

The ECS is in charge of the configuration, monitoring and control of all areas of the experiment; this comprises classical slow controls of high and low voltages, fluid-systems, various sensors as well as monitoring and control of the DAQ and HLT systems. It provides an homogeneous and coherent interface between the operators and all experimental equipment, as shown in figure 103.



**Figure 103.** Scope of the ECS.

The ECS for the upgraded detector is an evolution of the current system, described in [158]. It is still developed in the context of the joint control project (JCOP) [167], a common development between the four LHC experiments and CERN. The project defined a common architecture and a framework to be used by the experiments in order to build their detector control systems.

### 10.4.1 Architecture

JCOP adopts a hierarchical, highly distributed, tree-like structure to represent the structure of subdetectors, subsystems and hardware components. This hierarchy allows a high degree of independence between components, for concurrent use during integration, test or calibration phases. It also allows for integrated control, both automated and user-driven, during physics data-taking. LHCb adopted this architecture and extended it to cover all areas of the experiment.
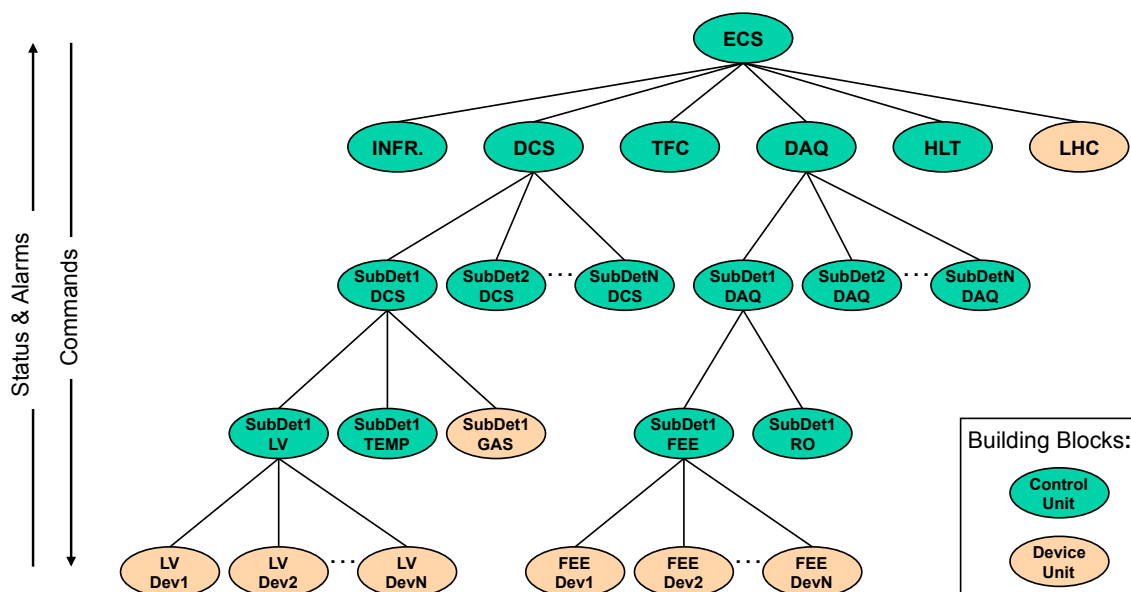


**Figure 104.** Simplified ECS architecture. Reproduced from [169]. CC BY 4.0.

Figure 104 shows a simplified version of LHCb's control system architecture. The building blocks of this tree can be of two types: Device Units (DU), the tree leaves, which are capable of driving the equipment to which they correspond, and Control Units (CU) which correspond to logical subsystems and can monitor and control the subtree below them.

### 10.4.2 Framework

The JCOP framework provides tools for the integration of the various components in a coherent and uniform manner. It is built upon a supervisory control and data acquisition (SCADA) system.[73]

While the SCADA system offers most of the needed features to implement a large control system, the CUs described above are abstract objects which are better implemented using a modelling tool. For this purpose SMI++ [168] was integrated into the framework. SMI++ is a toolkit for designing and implementing distributed control systems. Its methodology combines three concepts: object orientation, finite-state machines and rule-based reasoning. The JCOP framework is also

---

[73]Formerly called PVSS II, now WinCC-OA, http://www.etm.at.

complemented with LHCb specific components, providing for the control and monitoring of LHCb equipment or components such as DAQ electronics boards, power supplies or HLT algorithms.

### 10.4.3   DAQ and electronics control

The upgraded electronics are integrated into the control system following the philosophy described above. Standard LHCb components have been developed which allow users to configure, monitor and interact with their electronics. The upgrade electronics specifications document [68] contains requirements and guidelines for electronics developers, so that common software can be implemented.

As described in section 10.3, the ECS interface to the FE electronics is implemented via SOL40 interface boards, using the GBT system. This bi-directional link allows writing and reading of configuration and monitoring data. The GBT-SCA chip provides an interface between the GBT and standard protocols such as I2C, SPI or Joint Test Action Group industry standard (JTAG) and can be mounted on the FE modules, as shown in figure 105.
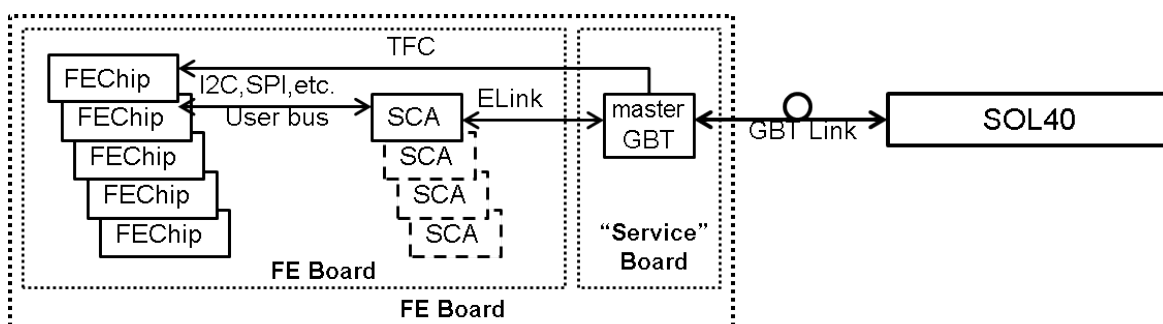


**Figure 105.** FE-ECS interface. Reproduced from [68]. CC BY 4.0.

A generic server process running inside the PC hosting the SOL40 cards provides the interface between the FE electronics and the SCADA system. Similarly to the FE electronics, the software for the configuration and monitoring of BE boards is maintained centrally in the form of JCOP components providing for the high-level description and access to all electronics components.

### 10.4.4   Guidelines and templates

Configurable framework components are distributed to the subdetector and subsystem teams in order to build their specific control systems. In order to ensure the coherence and homogeneity of the system, detailed guidelines specifying naming and colour conventions have been prepared and distributed. Whenever possible, the code necessary to implement the guidelines and conventions or the code to implement the finite-state machine behaviour specified for the different LHCb domains is also centrally provided in the form of templates.

### 10.4.5   Operations and automation

Like in the previous system, all standard procedures and, whenever possible, error recovery procedures are automated using the JCOP framework finite-state machine tools [169]. The experiment's operation, in terms of user interfaces, is again based on the JCOP framework and SCADA system, providing a global Run Control, control panels for detector systems, and alarm screens. As an example, the Run Control panel is shown in figure 106.

**Figure 106.** LHCb Run Control panel.

## 10.5 Monitoring, alignment and calibration

The software trigger applications move event data through the selection stages of HLT1 and HLT2. In between and at the end of these stages, events are stored in large global buffers. This mechanism makes data available also for monitoring and intermediate calibrations performed by dedicated applications and processing mechanisms, which are described in this subsection.

### 10.5.1 Monitoring

Automatic and interactive monitoring are essential tools to ensure the quality of the recorded data. Broadly speaking, monitoring works with two categories of input data. The first category is represented by histograms and counters produced by selection algorithms and other processes throughout the system. These are collected by so-called adder-processes, which accumulate quantities and histograms specific for partitions (see section 10.2), and are identified by run and fill number. Interactive analysis of these quantities is done using the MONET [170] web-based application. Data are acquired and published via LHCb's standard publish-subscribe framework, the distributed information management (DIM) system protocol [171]. Automated algorithms can access data with the same mechanism to generate alarms in the experiment control system. The second category consists of fully assembled events obtained after event building which are designated for monitoring by random selection or by the trigger processes. These events are tagged for monitoring when they are made available in the shared memory of the individual processing nodes. Monitoring-tagged events are duplicated

and picked up by the standard data processing, but also sent to the monitoring farm, via a separate network. In the monitoring farm they are first received by a distribution software layer implemented using the same shared memory architecture as the rest of the system. This layer implements various access patterns to monitoring data:

- every event is guaranteed to be delivered to a single monitoring application (*consumer*);

- every event is guaranteed to be delivered to one of a group of similar consumers for load balancing;

- events are offered to all interested consumers.

Guaranteed delivery is ensured by back-pressure. This mechanism can ultimately block the data flow, but an excess of input data is normally detected by the ECS before a complete stall. The distribution layer duplicates event data as needed. Consumers need to decode the raw-data, but otherwise have no limitations other than available computing resources. They are implemented using the same software frameworks as used by the trigger algorithms. They produce counters and histograms as output, which are monitored by the shift-crew and automated analysis systems.

### 10.5.2 Calibration and alignment

Calculating calibration and alignment constants is crucial for the optimal selection of events. As explained earlier, events selected by HLT1 are stored in a buffer storage. From this storage the events are accessed by sampling processes running on the event-filter farm nodes normally used for the HLT2 selection algorithms. Because the calculations are distributed, alignment and calibration algorithms can be run very quickly, even if large samples are required to obtain the required precision. The constants are typically available in a few minutes for the parameters which need to be updated as soon as possible in HLT1. The update is immediately propagated to HLT1 when the new parameters are available via a run change. Later, HLT2 consistently uses the same parameters. The other parameters used only in HLT2 processing are usually available less than one hour after data taking and can be immediately fed into the HLT2 algorithms, which can start event processing, organised by LHC run and fill, only after these parameters are available. This process is described in more detail in section 11.5.

### 10.6 Computing farm

While the first level selection (HLT1) is done on GPU cards installed in the event-builder nodes, the second level selection (HLT2), as well as calibration and alignment processes, are running on a large farm of general purpose CPU servers called the event-filter farm or computing farm. These servers are typically dense computing nodes, requiring half a rack-unit per server. They are connected via an Ethernet network, with a speed which tries to match approximately their relative processing power. The servers are quite heterogeneous because they come from many different procurement procedures. More than 3000 servers are available. Up to 80 can be packed into a single rack. When new batches of servers are acquired, the oldest servers are retired, which maximises reliability and optimises electrical power usage. Although the servers are currently not yet fully operated as a computing cloud, many characteristics of cloud computing apply. All nodes run the same operating system (currently, CentOS Linux 7), do not store anything relevant locally, and can be easily replaced, which makes it easy to use a wide variety of different machines of widely varying age and performance. From the point of

view of the ECS, their key characteristic is their processing power, in which they differ. Local storage is only used for scratch space and the OS installation. Management is done using industry standard tools.[74] Physically they are installed in the LHCb data centre described in this paper in section 2.4.4.

## 10.7 Infrastructure

The event-builder and event-filter farms are implemented directly on dedicated servers to simplify management and have more flexibility, although it can be expected that their deployment over a cloud system will grow over time. In contrast, the monitoring, the ECS, various infrastructure services such a logging, low-level system monitoring and others run on a typical enterprise infrastructure.

To ensure a uniform environment, a virtualisation cluster has been implemented, running on top of RedHat Enterprise Virtualisation hypervisors. The event-builder and event-filter applications, except as mentioned, are run on virtual machines. Many of these applications are determined by peak-load and many are very memory intensive for comparatively little CPU use. Logically, the virtualisation is divided into two clusters. The first cluster provides the core control system function, hosting all services needed for safe and controlled data taking. These services are provided with guaranteed resources which are fully redundant, ensuring that if a hypervisor or group of hypervisors fail, there are enough resources on the remaining hypervisors to completely take over after a restart. The second cluster hosts all services which are not critical, such as development or general purpose machines and certain control functions not required during data taking. Enterprise virtualisation requires a shared storage system, which is provided by a commercial filer.[75] An independent filer provides all the shared file-systems for the entire cluster.

To accommodate more appropriately the growing number of services originating in the cloud computing world, a Kubernetes [175] cluster is also available. Services running on the cluster must provide redundancy at the application level, because in this setup no attempt is made, beyond basic best practices, to provide hardware redundancy.

The computing infrastructure is comprised of almost 5000 physical servers and more than 10000 network ports. In addition to physics data-taking, the infrastructure is used to serve other computing needs of LHCb.

## 10.8 Performance

The performance of the event-building network has been studied using simulated data generators to emulate the TELL40 boards. Figure 107 shows the scaling of the total event-builder data throughput as a function of the number of builder units per event-builder server (each server implements two independent builder-units). The deviations from ideal (i.e. linear) scaling behaviour are believed to be related to remaining hardware and low-level tuning issues, in particular the underlying fragment (message) size, which will be addressed during the commissioning period. However, already at this level the performance requirements are met with a sufficient margin.

## 11 Trigger and real-time analysis

The objective of the trigger system is to reduce the data volume, which reaches 4 TB/s at the nominal instantaneous luminosity in $pp$ collisions, to around 10 GB/s which can be recorded to permanent

---

[74]Among them Foreman [172] and Puppet® [173].
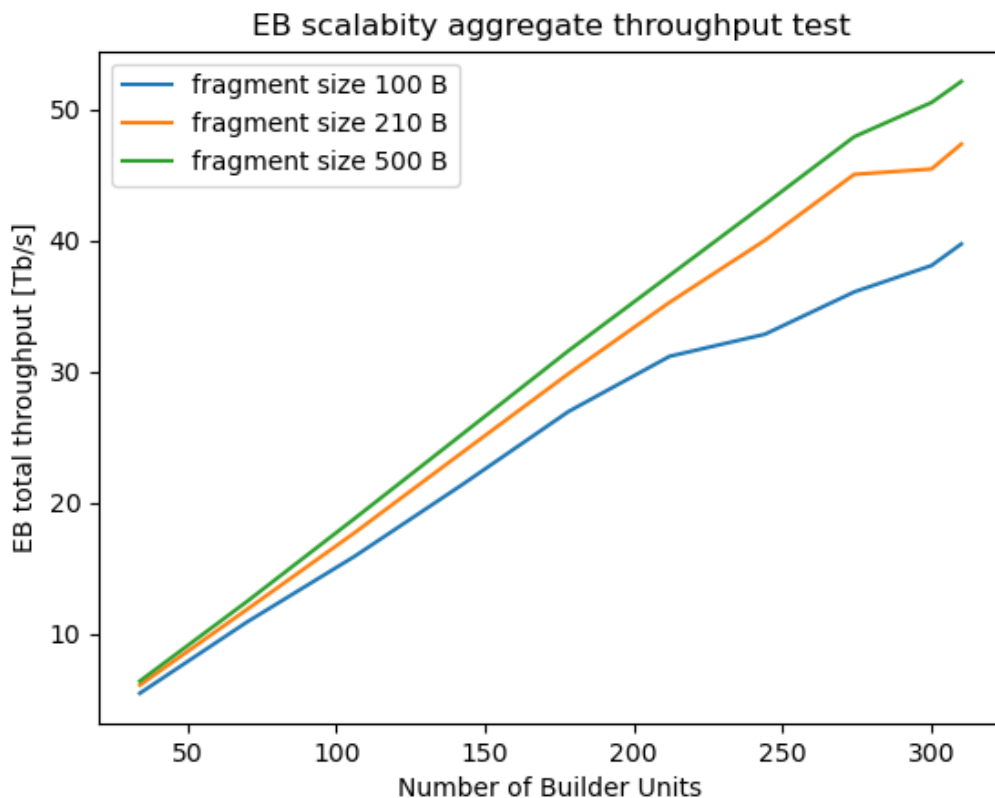
[75]NetApp® [174] filer.

**Figure 107.** Aggregated data-rate in the LHCb event-builder network as a function of the number of builder units for various nominal event-fragment sizes.

offline storage. As discussed in section 1, this reduction of a factor 400 is complicated by the high rate of potentially interesting signals which can be at least partially reconstructed in the detector acceptance. Over 300 kHz of bunch crossings contain a partially reconstructed beauty hadron and almost 1 MHz of bunch crossings contain a partially reconstructed charm hadron [9].[76] It is therefore not possible to fully pursue a traditional inclusive trigger strategy in which events containing the signals of interest are identified using a small set of generic signatures.

Instead, the majority of LHCb physics channels must be selected by the trigger by fully reconstructing and identifying the specific signals of interest, while saving only a limited subset of information about the rest of the event [176]. This *real-time analysis* approach, pioneered by LHCb in Run 2 data taking [177], necessarily requires that the trigger performs a full offline-quality reconstruction, enabled by an alignment and calibration of the detector performed in quasi-real-time. The physics motivations behind real-time analysis have been documented in detail refs. [8, 178]. They result in a two-stage trigger system: a first inclusive stage or high level trigger,[77] the HLT1, based primarily on charged particle reconstruction which reduces the data volume by roughly a factor of 20, and a second stage, the HLT2, which performs the full offline-quality reconstruction and selection

---

[76]A particle is considered partially reconstructed when only part of the decay products are reconstructed and identified.
[77]This terminology is a legacy of the previous experiment where the trigger had a hardware level (L0) and a software level (HLT).

of physics signatures. A large disk buffer is placed between these stages to hold the data while the real-time alignment and calibration is being performed, as described in section 10.5.

## 11.1 Physics requirements

Because of the necessity to perform real-time analysis, the physics requirements of the HLT2 trigger stage are easily expressed though demanding: it must perform the full offline-quality detector reconstruction on all events, using the offline-quality alignment and calibration provided in real-time, and support on the order of 1000 independent selections for signals of interest without any restriction on the topology of those physics channels (see section 11.4 for more details).

The bulk of the non-trivial physics requirements concern HLT1, where a partial reconstruction is performed with the attendant trade-offs between speed, efficiency, and output rate. The objective of HLT1 is to reduce the event rate to a level at which the data can be buffered to disk for real-time alignment, calibration, and further processing in HLT2, while maintaining high efficiency across the LHCb physics programme. This rate depends on the throughput of the HLT2 full detector reconstruction and on the maximum write-speed possible to the LHCb disk buffer.[78] It is therefore helpful to approach the HLT1 physics requirements by looking at the signal rates expected in the LHCb upgrade, which have been documented in ref. [8] and served as the initial motivation for the development of LHCb real-time analysis processing model.

As discussed before, the signal rates are dominated by charm and beauty hadron decays which are at least partially reconstructible within the LHCb acceptance. Signal rates for other areas of the LHCb physics programme such as electroweak physics or studies of quarkonia are significantly lower. While some short-lived hadronic resonances are produced at higher rates, measurements of these states are limited to the percent-level by luminosity knowledge or data-driven determination of reconstruction efficiencies. Therefore the full data set is only needed in sparsely populated regions of the kinematic parameter space where the event rate is negligible. Strange hadrons have signal rates which are more than an order of magnitude higher than charm and beauty. However, accepting strange hadrons with high efficiency would conflict with the main physics aims of the experiment, which concern charm and beauty hadrons. Therefore preference is given to rare strange hadron decays, whose signal rates are negligible compared to the singly-Cabibbo suppressed charm hadron decays at the centre of LHCb charm studies.

Since the signal rates are so high, it is of paramount importance that the HLT1 selects events due to the presence of signal as opposed to fake ("ghost") tracks or random combinations of tracks.[79] If this requirement is met, at a luminosity of $2 \times 10^{33}\,\mathrm{cm}^{-2}\,\mathrm{s}^{-1}$ a maximal HLT1 output rate of $2\,\mathrm{MHz}$ is from first principles sufficient to cover the full LHCb physics programme.

Having defined the maximal HLT1 output rate which the trigger system should be able to support, it is important also to define the smallest HLT1 output rate at which it remains efficient for LHCb signals of interest. This is important because at times of very high collider efficiency the HLT1 output rate might need to be reduced in order to avoid overflowing the disk buffer. Studies [179] show that it is possible to reduce the HLT1 output rate to around $500\,\mathrm{kHz}$ while remaining relatively efficient for LHCb charm physics channels. This rate is thus considered as a lower limit for the HLT1 output rate.

---

[78]Throughput is the number of events processed per unit time and is expressed in Hz.

[79]Fake tracks are the result of the pattern recognition algorithms that may reconstruct false trajectories from random combinations of hits.

All these requirements imply that the HLT1 must be able to reconstruct and select at least the following physics signatures:

- Tracks or two-track vertices displaced from the primary $pp$ interaction, or $v$. This signature can be used to select any event containing a long-lived hadron or $\tau$ lepton, which covers the vast majority of LHCb analyses;

- Leptons, particularly muons, regardless of their displacement from the $v$. Displaced leptons can be selected as any other tracks, although the efficiency can be kept higher for the same output rate by using lepton identification criteria to allow displacement- or $p_T$-based criteria to be loosened. Non-displaced (di)leptons are particularly important for spectroscopy studies, exotic searches, and electroweak physics.

Based on these constraints, requirements on the track reconstruction itself can be defined:

- HLT1 must reconstruct all tracks in the VELO detector acceptance in order to accurately determine the position of primary vertices and calculate the displacement of tracks and secondary vertices produced in the decay of long-lived particles;

- HLT1 must be able to reconstruct tracks regardless of whether they are displaced from the $v$ or not;

- HLT1 must be able to reconstruct track momenta at the percent level, in order to ensure a fast rise of the turn-on curve of the HLT1 efficiency versus the HLT2 reconstructed particle momentum (where the resolution is around 0.5%);

- the HLT1 reconstruction must provide an accurate and precise covariance matrix of the track measurement closest to the beam line, in order that the turn-on curve of HLT1 efficiency versus the HLT2 reconstructed particle displacement remains sharp;

- HLT1 must be able to identify tracks as muons or non-muons.

It is crucial that HLT1 is able to reconstruct tracks above a certain kinematic threshold which depends on the physics in question. Experience from Run 1 and Run 2 shows that a $p_T$ threshold of 500 MeV is adequate for most of the beauty and charm physics programme, although being able to reconstruct tracks down to $p_T = 200$ MeV would be highly beneficial for many areas of charm physics studies. For rare strange hadron decays it is important to reconstruct tracks with a $p_T$ as low as possible, and a more natural cutoff is to require 3 GeV of momentum, which corresponds to the lowest momentum at which it is possible to identify muons in LHCb.

While there are studies which could benefit from the reconstruction of calorimeter quantities (like e.g. energy deposits by photons) in HLT1, this is not considered as a mandatory requirement because a good portion of the signal efficiency can be achieved using tracks alone. A similar argument applies to electron identification or the reconstruction of downstream tracks, specifically $K^0_S \to \pi^+\pi^-$ and $\Lambda \to p\pi^-$ decaying outside the VELO. Both would significantly benefit certain areas of the physics programme but they do not represent a strict requirement for a functioning HLT1.
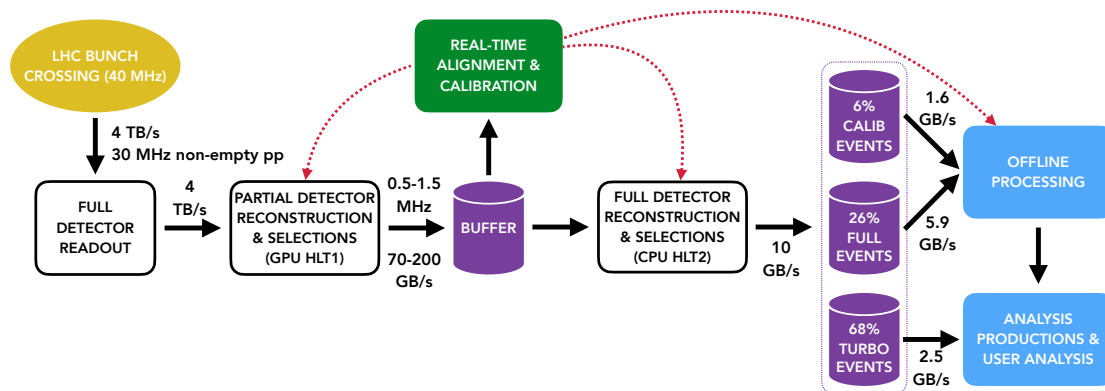
**Figure 108.** Online data flow. Reproduced with permission from [180].

## 11.2 Design overview

### 11.2.1 Hardware design

As described in section 10 the centrepiece of LHCb trigger is the full detector readout and event building at 30 MHz.[80] In addition to its inherent flexibility, this approach also matches the described physics requirements for HLT1, most notably the requirement that all the tracking algorithms must reconstruct the events at 30 MHz. Since the tracking detectors constitute approximately two-thirds of LHCb overall data bandwidth, the cost of a hardware trigger allowing the other one-third (principally RICH and calorimeters) to be read out at a lower rate does not justify the added complexity and reduced flexibility of such a system. The trigger data flow from detector to offline storage is illustrated in figure 108.

The online disk buffer serves two purposes: it stores events selected by HLT1 while real-time alignment and calibrations are being performed, and it allows events selected by HLT1 to be buffered for processing between LHC fills. This in turn effectively increases the processing power of the HLT2 computing resources. Simulations of the LHC fill structure [181] result in an optimal buffer size of around 30 PB, which allows around 80 hours of LHC collisions to be buffered at an HLT1 output rate of 1 MHz. The disk buffer architecture is constrained not only by the total size but also by I/O limitations of individual disks, which impose a minimum number of disks in the system.

Because the detector is fully read out upfront, there is in principle total freedom to choose the computing architecture with which to process the data; the only constraints are the available budget and the capacity (racks and cooling) of the data centre. LHCb has sought to take advantage of this freedom since about the end of Run 1, pursuing the development of high-throughput reconstruction algorithms on both CPU and GPU architectures. Eventually, a full cost-benefit analysis [182] led to the choice to implement HLT1 on GPUs while remaining with the same CPU architecture used during Run 1 and Run 2 for HLT2. Because the GPUs are hosted in the event-building servers as described in section 10, there is a limit of around 500 GPUs which can be installed for running HLT1. As the baseline HLT1 described in the remainder of this paper requires only around 200 latest generation GPUs, this limit does not introduce a significant constraint into the system.

---

[80]Although the LHC bunch crossing frequency is 40 MHz, the rate of events with some signal visible inside the LHCb acceptance (*reconstructible* events) is ∼ 30 MHz.

### 11.2.2 Software design

The software design of the LHCb upgrade trigger is guided by the principle that LHCb's real-time and offline processing should be as similar as possible. More concretely, the upgrade code base is set up so that any offline reconstruction and selection is a specific reconfiguration of the same underlying algorithms as the real-time version. Much of the quality assurance and validation machinery is also shared, easing the maintenance burden. This convergence is partially a result of the physics requirements, in particular the need to run the full offline-quality reconstruction within the trigger. In fact, the LHCb code base has been evolving in this direction throughout Run 1 and Run 2, so the almost total convergence achieved for the upgrade can be seen as a natural completion of this process.

Most of the code base is a mixture of C++ used for data processing and algorithms, PYTHON for job configuration and CUDA in the GPU-based HLT1. The CPU code base is built on top of the GAUDI [183, 184] framework.

Both are optimised for multithreaded execution. The GPU code base is implemented within the cross-architecture ALLEN framework [185]. ALLEN is based on a single source code, that can be compiled for both CPUs and GPUs, supporting different GPU architectures through CUDA and HIP languages. The framework also provides a custom memory manager to optimise GPU memory usage, and a multievent scheduler to handle the data and control flow of events processed in batches during data-taking In the context of event simulation HLT1 ALLEN algorithms are converted to GAUDI algorithms automatically and event processing is controlled by GAUDI to fit into the offline processing chain (section 12.6).

In order to ensure thread-safety, the design of algorithms ensures that there are no race-conditions and that results are not altered by the execution order of the threads. This in turn allows the configuration to largely define the data flow, with the control flow automatically deduced by the scheduler; manual overriding remains an option to resolve ambiguities or enforce a preferred execution order. The detector geometry and conditions are read in from dedicated databases, as described further in sections 12.2 and 12.3. Although the scheduler is fully flexible, in what follows, for simplicity of exposition, the description of the reconstruction and selections will be separated as if the trigger executed a monolithic reconstruction followed by multiple selection algorithms.

The code base is maintained by a common effort between the developers of LHCb's various software projects, with continuous integration, testing, and code review before any changes are made. A rotating team of (less experienced) shift takers supported by (more experienced) maintainers supervises this process and carries out the task of merging and releasing specific versions of the code for use in production.

### 11.3 First trigger stage

The HLT1 first level trigger design follows the requirements described in section 11.1. The baseline implementation in terms of reconstruction and selections is described here. Key performance figures are documented in detail in ref. [186] and are not repeated here.

### 11.3.1 Reconstruction

The baseline HLT1 reconstruction focuses on finding the trajectories of charged particles which originate within the LHCb vertex detector and traverse the rest of the LHCb tracking system (specifically the UT, SciFi Tracker, and muon chambers). The objective is to measure their momenta with percent-
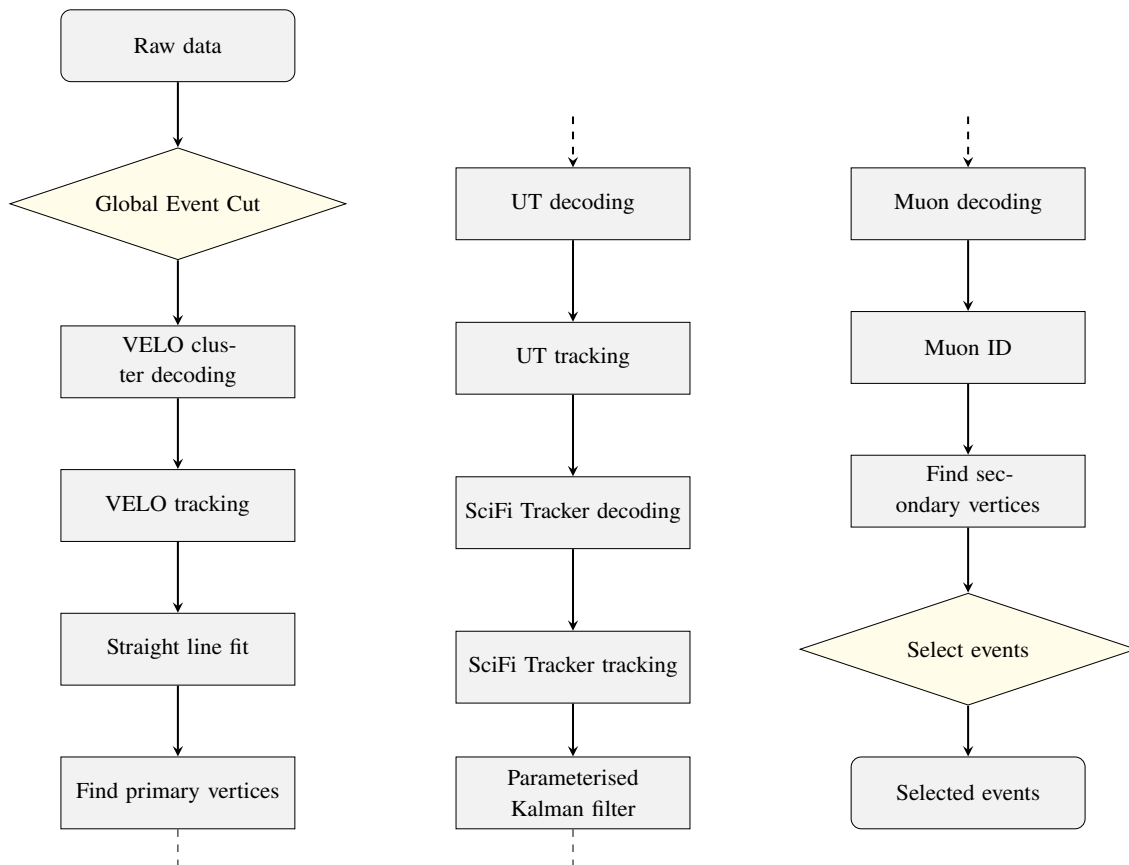
**Figure 109.** Baseline HLT1 sequence. Reproduced with permission from [186]. Rhombi represent algorithms reducing the event rate, while rectangles represent algorithms processing data.

level precision, associate each particle to the $pp$ collision where it was produced and measure its displacement from that $pp$ collision, and identify the particle as a muon or non-muon.

The reconstruction sequence is shown schematically in figure 109. It is preceded by a Global Event Cut (GEC) which removes a fraction of the events with a very large number of tracks, which cost a disproportionate amount of computing time to reconstruct while having a worse detector performance. The cut can in principle be applied based on the occupancy of any combination of subdetectors. The baseline criterion is to reject the 7% of busiest minimum bias events based on UT and SciFi Tracker occupancies. This will be revisited during detector commissioning to make sure that the cut uses information from those subdetectors whose data occupancies best match simulation. As necessary, special reconstruction and selection sequences can be deployed which bypass this GEC, e.g. for electroweak physics or very forward high transverse momentum exotic signatures. These are not discussed further here for brevity.

From the algorithmic point of view, the reconstruction sequence is straightforward:

1. tracks are reconstructed in the VELO and are used to locate the positions of the primary vertices where the beam collisions occurred;

2. tracks are extrapolated to the UT, and subsequently to the SciFi Tracker detector, based on a minimum allowed momentum and/or transverse momentum. A magnetic field parametrisation is

used to predict the track position and speed up the reconstruction. In the baseline configuration, the minimum $p_T$ threshold is set to 500 MeV, but this value will be updated once the performance is evaluated with real data. Reducing this threshold requires improved computational performance but increases the number of tracks available to physics selections and consequently the signal efficiency. This is particularly important for signatures typical of strange and charm quark physics;

3. above the $p_T$ threshold, the track momentum is known with a precision better than 1% across the full momentum range of interest. The momentum is thereafter used as input to a parameterised Kalman filter [187, 188] which estimates the position and covariance matrix of the particle at the beam line. This information is in turn used to calculate the particle's displacement from all primary vertices with a resolution similar to the one achievable in HLT2 with a fully aligned and calibrated detector. Then, selections typically apply requirements based on a minimal displacement from the primary vertex which is equivalent to at least two or three standard deviations;

4. tracks are identified as muons or non-muons and fitted to a common origin to form two-body displaced vertex candidates, which are then input to the selections.

This reconstruction sequence largely matches that of LHCb Run 2 HLT in conceptual terms. The main difference is that only a simplified Kalman filter is used in order to reduce processing time. Further detailed documentation can be found in refs. [186, 189].

The sequence processing time is further reduced by exploiting prereconstructed hit clusters from the VELO, built at an early stage into the VELO TELL40 FPGA firmware by a 2D cluster finding algorithm. The FPGA algorithm identifies VELO clusters in a fully parallel way by pattern matching, and computes their topology and position centroid in real time providing then the results to HLT1 [57].

### 11.3.2 Selections

The HLT1 selections are broadly divided into four categories: the primary inclusive selections for the bulk of LHCb physics programme; selections for calibration samples essential to a data-driven evaluation of the reconstruction performance; selections for specific physics signatures not covered by the inclusive triggers; and technical triggers for luminosity determination, monitoring, calibration and alignment. Here, a brief description of the physics logic of these selections is given, while their performance is described in section 13.3.

The primary inclusive selections are:

- a *two-track vertex trigger*, requiring large transverse momentum and significant displacement from all primary vertices in the event; this trigger provides the bulk of the efficiency for nonmuonic signatures;

- a *displaced single-track trigger*, requiring large transverse momentum and significant displacement from all primary vertices in the event [190]. This is the most inclusive of the HLT1 triggers as well as the least biasing (because it fires on a single particle) of the hadronic triggers. However, because of the abundance of genuine displaced hadrons in upgrade conditions, particularly from charm decays, its main role is to provide redundancy for the two-track vertex trigger;

- a *displaced single muon trigger*, which operates similarly to the single-track trigger, but requires that the track be identified as a muon; this in turn allows the transverse momentum and displacement requirements to be significantly relaxed;

- a *displaced dimuon trigger*, which operates similarly to the two-track vertex trigger, but requires that both tracks be identified as muons; the transverse momentum and displacement requirements can be relaxed even further compared to the single muon trigger;

- a *high-mass dimuon trigger*, which does not make any track displacement requirements, but requires that the dimuon invariant mass be above 2900 MeV; this trigger is particularly important for the study of charmonia, whether in prompt production or from the decays of beauty hadrons;

- a *very high transverse momentum muon trigger*, which does not make any GEC requirement and is used for electroweak physics and searches for exotic signatures;

When selecting calibration samples, it is particularly important to ensure a good coverage across the kinematic spectrum. Therefore, while charmonia are already selected by the high-mass dimuon trigger, the calibration selections focus on the decays of charm hadrons such as $D^0 \to K^- \pi^+$, which are critical for a data driven evaluation of particle identification performance. Analogous selections are foreseen for other charm hadron species. Additional selections are used to create samples enriched in tracks which traverse regions of the detector particularly critical for the tracker alignment or RICH mirror alignment. HLT1 throughput scales to $O(100)$ selections. It is therefore likely that numerous triggers for specific hard-to-select signals will be implemented as data taking progresses. Examples from previous LHCb data-taking periods included diproton triggers, diphoton triggers, decay-time unbiased charm triggers, and others.[81]

## 11.4 Second trigger stage

The second level trigger (HLT2) uses the information provided by the real-time alignment and calibration of the detector to perform an offline-quality reconstruction, followed by $O(1000)$ selection algorithms which decide whether or not to retain any given event. If an event is retained, the selection algorithms specify which portions of the full event are recorded to permanent storage following the real-time analysis paradigm developed during Run 2 [176, 177].

### 11.4.1 Reconstruction

The LHCb reconstruction is divided into four main components: charged particle pattern recognition, calorimeter reconstruction, particle identification, and the Kalman fit of reconstructed tracks which allows their parameters to be measured with the best possible precision and accuracy (see section 11.4).[82] As the upgrade detector design is not fundamentally different from that of the previous LHCb detector, its reconstruction is also conceptually similar to what was used during Run 1 and Run 2. Nevertheless some changes and improvements have been possible in specific areas, as described here.

---

[81]Although there is no full calorimeter reconstruction in the HLT1, a rudimentary calorimeter clustering has been implemented in ALLEN, allowing to reconstruct electrons and photons. This was possible due to computing resources freed by the calorimeter preprocessing (low level trigger) discussed in section 8.2.4.

[82]Within LHCb, *Kalman fit* identifies a fit of track parameters based on a Kalman filter which combines the coordinates of an ensemble of selected hits and additional information like magnetic field and material distribution maps.
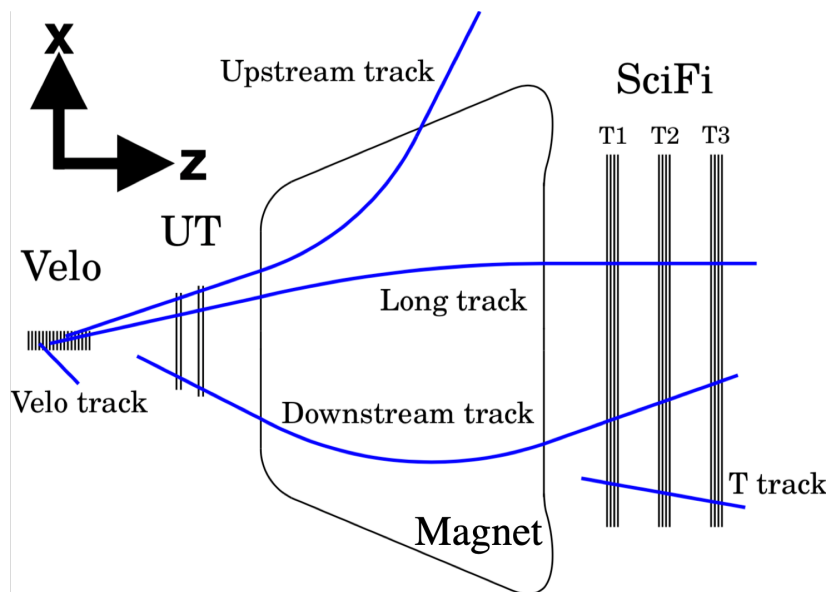
**Figure 110.** Track types in the LHCb detector bending plane. Reproduced from [1]. © 2008 IOP Publishing Ltd and Sissa Medialab. All rights reserved.

**Charged particle pattern recognition.** Different tracking algorithms exist to reconstruct different track types, illustrated in figure 110. Tracks which originate in the vertex detector (*VELO tracks*) are used to determine the positions of the primary $pp$ collisions, a process known as primary vertex finding. The combination of $v$ positions and track trajectories, in turn, allows tracks which originate from the decays of long-lived particles and are therefore displaced from the $v$ to be precisely identified. As there is effectively zero magnetic field inside the VELO, these tracks must be extrapolated into the region covered by the UT (*upstream tracks*) and SciFi Tracker (*long tracks*) in order to measure their momentum. Long tracks have the most precise and most accurate momentum determination and are used in nearly all LHCb analyses. In addition to the *forward* algorithm which extrapolates VELO tracks to the SciFi Tracker, a second redundant reconstruction path (seeding) performs a standalone reconstruction of track segments in the SciFi Tracker (*T tracks*) before matching them to VELO tracks and optionally UT hits. In addition, SciFi Tracker seeds are extrapolated to the UT and used to form *downstream tracks* in order to reconstruct particles which originate outside the VELO but before the UT. Downstream tracks provide the bulk of LHCb statistical power for the study of decays involving strange hadrons. The track extrapolations used in all of these pattern recognition algorithms are, for reasons of speed, based on parametric models of trajectories in the LHCb magnetic field. Duplicated tracks (*clones*) can be formed when different algorithms reconstruct the same track segment in one of the subdetectors, for example when a long track and a downstream track share a T-station seed. These are filtered by removing duplicates within individual pattern recognition algorithms. Following the Kalman fit, a global clone-killing algorithm uses the fit quality to perform a final arbitration between overlapping VELO, long, and downstream tracks and removes the remaining duplicates.

The charged pattern recognition algorithms have undergone significant evolution from the Run 1 and Run 2 code in order to make them better able to efficiently use modern multicore CPU architectures. An example of such optimised algorithms is described in detail in ref. [191], while their performance is documented in section 13.
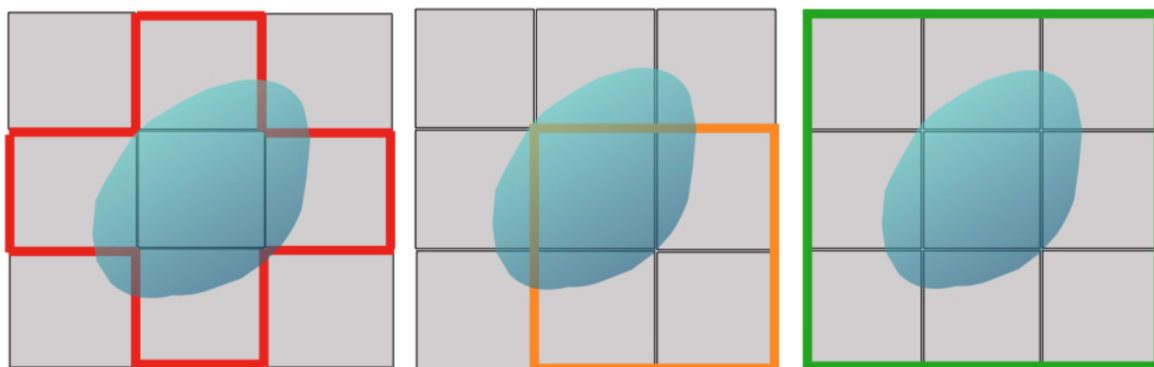
**Figure 111.** Shapes used for the upgrade calorimeter reconstruction, referred to as *cross* (left), *2 × 2* (centre) and *3 × 3* (right).

**Kalman fit.** While the pattern recognition is used to accurately group detector hits into collections corresponding to individual charged particles, a separate step based on a Kalman filter (referred to as the *Kalman fit*) is required in order to determine the properties of charged particle trajectories with maximum accuracy and precision. Several approaches exist: a detailed Kalman fit which uses lookup tables to describe the magnetic field and material distribution and Runge-Kutta methods to propagate the particle trajectories, an intermediate solution where the propagation is performed as in the detailed Kalman fit but the interactions with the detector material are parametrised, and a fully parametric Kalman fit also parametrising the particle propagation through the LHCb magnetic field. The last option is described in more detail in ref. [192]. Despite a significant investment of time and effort, a full parallelisation of the detailed Kalman fit has remained difficult. However, significant gains in speed have been achieved by simplifying the data structures used inside the algorithm and reducing the amount of output data to what is strictly necessary for physics analysis. In addition, the software framework is sufficiently flexible to allow the use of a mixture of parametric and detailed Kalman processing, depending on the specific application. While the parametric Kalman fit achieves a performance suitable for nearly all physics applications, the detailed Kalman fit remains necessary for use in the detector alignment procedure, as well as for certain analyses which require the ultimate precision on track states. The detailed Kalman fit as well as the solution with the parametrised material description will be therefore both available in HLT2 and in offline data processing and will be used depending on the specific needs.

**Calorimeter reconstruction.** Although the LHCb calorimeter system comprises a hadronic and an electromagnetic calorimeter, during Run 1-2 the HCAL was used almost exclusively for the first level hardware trigger which has now been removed for the upgrade. At present there are no upgrade analyses which foresee using HCAL information in the real-time processing and therefore only the ECAL is reconstructed.

ECAL clusters are formed from $3 \times 3$ cells, as shown in figure 111. Because of the higher pileup expected in upgrade conditions and the limited spatial granularity of the calorimeter, a combination of cross and $2 \times 2$ shapes is needed to obtain the most accurate position and energy resolution for each individual cluster. Multivariate algorithms based on the shower shapes and individual ECAL cell energies are used to distinguish single-photon clusters from those formed by the merger of multiple photons, most notably in highly boosted $\pi^0 \to \gamma\gamma$ decays. Electron clusters are identified by

extrapolating tracks to the ECAL region and subsequently matching them to ECAL clusters. Dedicated algorithms for calorimeter reconstruction have been prepared for the upgrade, with completely new structure and inherent algorithmic logic to improve the computational and physics performance; the calorimeter reconstruction performance is documented in section 13.

**Particle identification.** LHCb uses a combination of its two RICH detectors, the ECAL, and the muon system in order to identify the five basic long-lived charged particle species — electron, muon, pion, kaon, and proton. Unlike in central detectors, tau leptons are considered as any other composite particle in LHCb and there are no centralised identification algorithms for them; the same holds for neutral particles decaying in the detector. The efficiency to correctly identify the charged particle corresponding to any given track is heavily dependent on the charged particle species and on the dominant backgrounds. In general, as shown quantitatively in section 13, performance depends most strongly on the particle momentum, then on its pseudorapidity, and then on the detector occupancy. The different subdetectors dominate particle identification performance in different momentum regimes, except for muons where the muon system plays a dominant role in all cases. The particle identification performance depends critically on an accurate knowledge of the track trajectory within each particle identification subdetector, and therefore requires the use of Kalman fitted tracks in order to achieve the best results.

The standalone RICH reconstruction is described in detail elsewhere [105]. It has undergone few conceptual changes but the software has been rewritten and reoptimised for computational efficiency and speed.

The standalone muon reconstruction has been gradually improved throughout Run 1-2, and for the upgrade, new variables have been introduced which further improve the performance in the high-occupancy regime of the upgrade by considering correlations between the hits in the different layers of the muon detector [193].

While optimal absolute performance is achieved by combining the information provided by each subdetector into global multivariate classifiers, it is equally important to have a particle identification performance which can be accurately calibrated using data tag-and-probe samples, as described in section 11.5.3.[83] This is particularly true of the precision measurements which make up the bulk of LHCb physics programme and which require permille-level control of particle identification efficiencies and misidentification rates in order not to be limited by systematic errors. For this reason, multiple multivariate classifiers, trained on simulation and tuned to have a better and more stable performance in different kinematic regions, exist. The choice of which classifier is optimal for any given analysis and where particle identification information is used at the trigger level is left to the analysts which have to ensure the relevant samples exist which can calibrate these classifiers in the kinematic regions of interest.

### 11.4.2 Selections

Unlike HLT1, where a limited number of largely inclusive selections are sufficient to select the most interesting events, HLT2 relies on about one thousand different selection algorithms, each tuned for a particular signal topology and/or physics analysis. Although the software framework provides full

---

[83]In the tag-and-probe method typically a two-body decay of a resonance (e.g. a $J/\psi$) is used to identify a track independently of PID algorithms. One of the two tracks is used to *tag* the identity of the other (*probe*) which can be used to calibrate the PID algorithm.

flexibility to schedule interleaved sequences of reconstruction and selection steps, in practice almost all selection algorithms are executed once the complete offline-quality reconstruction has been performed.

In order for these selections to achieve the necessary computing throughput, tracks and neutral objects are zipped together with particle identification information into Structure-of-Array data structures which can be efficiently processed in parallel. Both rectangular-cut-based and multivariate or artificial intelligence-based selections can be deployed.

In addition to identifying which events should be recorded to permanent storage, each selection algorithm identifies which subset of event data to record. If multiple selection algorithms decide to record an event, the superset of information requested by them is recorded. This *real-time analysis* paradigm (also named Turbo analysis) allows the rate of recorded events to be increased by decreasing the volume of information recorded for each event. The Turbo concept was already deployed during Run 2 and a detailed description can be found in refs. [176, 177]. The Turbo mechanism allows full flexibility on the amount of event information that is stored (*selective persistence*), from the bare minimum of two tracks and vertex coordinates for a two-body decay, up to the full event information, depending on the specific physics channel under study, as described in ref. [194]. As also shown in figure 108, while a majority of triggered events will be saved in the reduced Turbo format, the majority of the data volume will consist of the calibration and traditionally triggered (*full*) events. In order to minimise the overall data volume, HLT2 selections will be grouped into *streams*, with all selections belonging to a stream sharing an underlying physics logic and recording similar sets of event information. Streams can be configured according to broad physics channels like e.g. charm physics, hadronic beauty decays, leptonic decays, electroweak physics, etc. and will evolve during the experiment's lifespan, as needed.

### 11.5 Alignment and calibration

The fact that HLT2 performs a full offline-quality event reconstruction and selects the majority of events based on the real-time analysis paradigm also necessitates an offline-quality alignment and calibration of the detector in real-time. This serves two separate purposes. The first is to provide the most accurate alignment and calibration parameters to the real-time reconstruction and selections. This ensures that the physics parameters of interest, such as particle masses or decay-times, are computed with the best possible resolution, maximising the selection efficiency. In addition, calibration parameters can be stored and made available for offline physics analysis without (in most cases) the need for further calibrations, simplifying the analysis workflow. The second is to provide large tag-and-probe samples which can be used offline to calibrate the difference between the detector performance in real data taking and its performance in simulation. This is particularly critical for physics studies which require a permille-level accuracy of single-particle reconstruction and identification efficiencies for each particle type.

Alignment and calibration procedures have been designed to maximise the physics reach and analysis flexibility rather than imposing centralised calibrations. For example, analyses which can benefit from offline recalibration can be stored in the FULL stream. A notable example is electroweak physics in which the very high transverse momentum of the signal decay products are very sensitive to any residual misalignment. The optimal alignment needed for electroweak physics can be obtained using large samples of very high momentum tracks (typical of $W$ and $Z$ boson decays) which can be collected only after a few months of data taking, thus requiring an offline calibration which
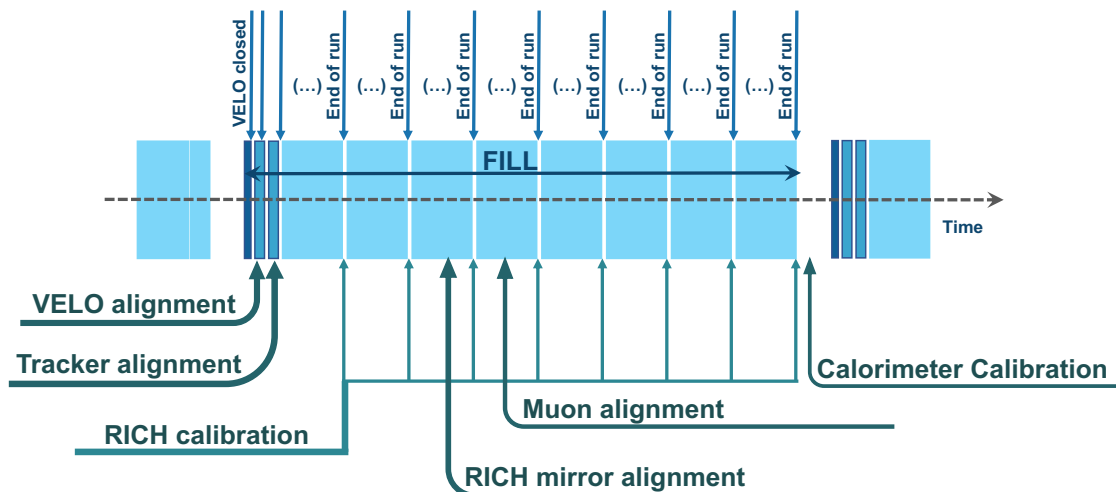
**Figure 112.** Schematic view of the real-time alignment and calibration procedure starting at the beginning of each fill.

can be performed for example once per year. Even these analyses, however, benefit from having the best possible alignment and calibration available in real time, as it minimises the difference in measured quantities used in online and offline selections and therefore simplifies the modelling of corrections for such effects.

Each step of the real-time alignment and calibration procedure uses different input samples and is performed at a different frequency. This is illustrated in figure 112 based on Run 2 operations. While the strategy will remain the same for the upgraded detector, the details will naturally evolve with commissioning experience.

### 11.5.1 Global alignment

A detailed description of LHCb global alignment procedure and strategy can be found in refs. [195, 196] and [139, 197]. The alignment of the LHCb detectors proceeds in a sequence, with the VELO aligned first, followed by the UT and SciFi Tracker detectors, the RICH mirror alignment, and finally the muon detector alignment. Extensive studies of alignment stability [198] have been performed during Run 2 in order to make sure that all possible lessons from the Run 1-2 LHCb detector are applied to the upgrade.

For the upgrade, the optimisation of the alignment configuration has been studied extensively using simulated samples, taking into account the survey measurements of each subdetector and its mechanical and thermal behaviour. It will be further tuned on the first data and its performance followed over time to identify and correct any trends.

The VELO alignment requires a track sample traversing all the modules in any azimuthal and radial position. Since the residual magnetic field in the VELO region is negligible, a large sample of minimum bias events can be used for this purpose. This sample is enriched with a sample of beam collision on the residual gas upstream of the VELO region, to select tracks originating far from the collision region and crossing several VELO modules. This data sample is collected in a few seconds at the beginning of the fill. The VELO alignment is performed using an additional constraint that tracks should come from their associated $v$ [196]. This alignment is performed at the start of each fill to account for the opening and closing of the VELO detector. In Run 2, only the alignment

parameters related to the relative displacement of the VELO halves were varying after the closing, and similar behaviour is expected also in the upgrade.

The alignment parameters used by the reconstruction are only updated if the observed changes exceed a certain tolerance; in this case the new alignment parameters become the reference values for the next fill, and so on. During Run 2, it was observed that updates were required on average every few fills. When parameters are updated, they are picked up in HLT1 via a run change, and the alignment used in HLT1 and HLT2 is always kept consistent to minimise systematic uncertainties associated with the reconstruction and selection efficiencies. Since the VELO alignment takes a few minutes to run, while a fill lasts around 8-10 hours, the fraction of data which is treated with an imperfect alignment in each fill is rather small.

For the alignment of the UT and SciFi Tracker detectors, a selection of signal tracks from the decays of well known resonances (notably $D^0$ and $J/\psi$) are used. The constraint of the daughter kinematics to the mother particle mass [196] is a powerful tool to significantly improve the alignment in the tracker system. Around $2 \times 10^4$ reconstructed resonance decays are required, which are collected by dedicated HLT1 selection algorithms and streamed to the alignment tasks. As the rate of fully reconstructed $D^0$ and $J/\psi$ signals will be larger than $10\,\mathrm{kHz}$ during the upgrade, such samples can be collected in a matter of seconds and the tracker alignment can therefore be executed at the start of every fill immediately following the VELO alignment. Once again, actual updates to the alignment parameters only occur when tolerances are exceeded. In Run 2 this typically occurred only for some elements of the detectors after magnet polarity changes or interventions on the detector during technical stops.

The RICH mirror alignment requires a sample of tracks selected such that their Cherenkov photons are distributed equally among the different RICH1 and RICH2 mirrors; this effectively means down-weighting tracks in higher occupancy areas of the detector and preferentially selecting high-purity tracks in the more peripheral areas of the RICH system. It consequently takes significantly longer to select the alignment samples, which may take up to a few hours. These corrections are not used in HLT1 and thanks to the disk buffer there is enough time to evaluate RICH alignment before running the HLT2 reconstruction.

Finally the muon alignment is performed using $J/\psi$ decays. Although very important for the performance of the L0 hardware trigger during Run 1 and Run 2, the muon detector alignment has a small impact on the overall tracking system performance. As the upgraded detector no longer uses a hardware trigger, the muon detector alignment should only need to be updated at the beginning of the data taking and in case of an opening of the muon stations due to a hardware intervention.

### 11.5.2 Calorimeter calibration

The performance and calibration of the calorimeter system is described in detail in ref. [128]. As the ECAL and HCAL are unchanged for the LHCb upgrade apart from their readout, no significant changes to the calibration procedure are foreseen. Due to the removal of the hardware trigger stage, only the ECAL requires accurate calibrations to ensure efficient electron and photon identification. The gain of PMTs is first calibrated in-situ using the LED monitoring system. Subsequently this LED system is used to monitor changes in the gain due to ageing. After each fill, the PMT high voltages are adjusted so that the LED signal of the ongoing fill matches the LED reference value. This procedure, introduced during Run 2 data taking, is able to control the ageing-induced PMT gain variations at the level of 1–2%. Secondly, a fine-grained calibration is performed for each ECAL cell based upon the reconstructed $\pi^0$ mass in that cell, achieving the ultimate possible accuracy. This calibration is performed once per month using $\sim 3 \times 10^8$ minimum bias events.

### 11.5.3 RICH gas refractive index calibration

In addition to the RICH and muon alignment and ECAL calibration, optimal particle identification performance also requires the calibration of the RICH gas radiator refractive index. As the refractive index is in particular highly sensitive to temperature and pressure variations, these calibrations are performed on a per-run basis using dedicated calibration samples selected by HLT1 [105]. Since the RICH is not used in HLT1, calibrations can be applied for each run while the data is stored on the disk buffer and made available to HLT2.

### 11.5.4 Tag-and-probe samples for offline data-driven corrections

The calibrations described so far ensure that LHCb takes data with an optimally performing detector at all times. However they do not address differences between the detector performance in data and simulation, which are of central interest to physics analyses. These differences are studied by collecting large tag-and-probe samples which allow for single particle reconstruction and identification efficiencies to be measured in the same way using data and simulation. The tag-and-probe methods are applied to pion, kaon, proton, electron, and muon samples within the LHCb acceptance; they also enable the study of charge asymmetries in reconstruction and particle identification.

The samples collected to measure particle identification efficiencies are described in ref. [199]. Extensive studies have been undertaken during LS2 in order to further optimise the kinematic coverage of these samples, particularly taking advantage of the removal of the hardware trigger to enhance coverage at the edges of the kinematic and geometric acceptance of the LHCb detector. During Run 1 and Run 2, these methods were able to measure particle identification efficiencies with an accuracy of a few permille, and it is expected that the same performance can be maintained.

For the study of track reconstruction efficiencies, large samples of displaced $J/\psi \to \mu^+\mu^-$ decays, in which one muon is fully reconstructed while the other is reconstructed in only part of the tracking system are used [200]. Tracking efficiencies measured with muons have historically been used also as a proxy for hadronic reconstruction efficiencies. However, alternative dedicated methods to directly measure hadron track efficiencies have been proposed and the relevant trigger selections implemented. All these methods require a dedicated reconstruction for the probe track, which must be carefully optimised to fit within the real-time timing budget. Because electrons lose a great deal of their energy to bremsstrahlung radiation, it is not possible to measure their reconstruction efficiencies with displaced $J/\psi \to e^+e^-$ decays alone. Thus, samples of $B \to XJ/\psi(\to e^+e^-)$ decays, where $X \in K^\pm, K^{*0}, \phi$, are reconstructed using only the VELO segment for the probe electron. The additional kinematic constraints of the $B$ mass and the distance of flight between the production $pp$ collision vertex and $B$ decay vertex give sufficient information to suppress backgrounds and measure the per-electron reconstruction efficiency with subpercent level systematic uncertainties [201].

## 12 Software and computing

As discussed in the previous section, the software trigger inherits the two-stage model already implemented in Run 1 and Run 2 and relies heavily on the mechanism of selective persistence to optimise the amount of data stored. Moving the event reconstruction and selection to the online domain implies that centralised offline data processing involves only the final data preparation for physics analysis. This process consists of data *skimming* and *slimming*, if required, and data streaming

according to physics content in order to optimise the amount of data that users need to access for their analysis.[84] This offline data processing is globally known as `sprucing` (see section 12.6). The last step of the data processing flow is represented by the physics analysis which may or may not be centralised. The `sprucing` utilises the same selection framework employed in HLT2, while physics analyses proceed through *analysis productions*, a centralised way to produce artefacts (e.g. ntuples) that are subsequently utilised for physics measurements.

The demanding processing rate of the trigger applications implies a redesign of the core software framework GAUDI, with the goal of optimising it for speed on current computing architectures, without compromising physics performance. The main lines of development in this major redesign consist of the usage of multithreading and vector registers, the optimisation of the data model, the extensive modernisation of code and algorithmic improvements, the adoption of modern technologies for the detector description and conditions data.

The offline computing model follows the distributed computing paradigm already exploited in Run 1 and Run 2. The amount of storage needed for the recorded data is driven by the trigger output bandwidth of 10 GB per second of LHC collisions. The offline computing work is dominated by the production of simulated events. Several avenues are exploited in order to mitigate the resource requirements. The following sections describe in more detail the concepts outlined above.

## 12.1 Core software

The GAUDI core software framework [183, 202] is the common infrastructure and environment for the software applications of LHCb. It was designed and implemented before the start of the LHC and the LHCb experiment, and it has been in production without major modifications ever since. Its main design principles remain still valid, however a review and modernisation were needed to meet the challenges posed by the LHCb upgrade and make it flexible to adapt to forthcoming challenges.

Previously, the major limitations of GAUDI were its weak scalability in RAM usage and its inefficient handling of CPU-blocking operations. These limitations have been addressed by introducing multithreading, a task-based programming model, and a concurrent data processing model where both interevent and intraevent concurrencies were considered. The introduction of a multithreaded approach changes the programming paradigm and introduces new guiding design principles, such as thread safety and reentrance, the declaration of data dependencies (needed for scheduling and concurrency control) and the immutability of data (that simplifies concurrency control and allows to focus on control- and data-flow rules).

The implementation of a concurrent, task-based framework according to the above design principles required a deep revision of the LHCb code base and a refactoring of many existing components. In particular, the declaration of data dependencies between GAUDI algorithms implied the explicit declaration of the input and output data requirements of algorithms. Categorising algorithms accordingly, and providing common interfaces and implementation rules for each category allowed code-developers to integrate their algorithms in the new task-based framework with minimal effort.

The introduction of multithreading in GAUDI had significant effects on the memory footprint of applications. As an example, figure 113 shows the memory utilisation of a prototype HLT1 application

---

[84]*Skimming* is the process of a further event rejection before the analysis, while *slimming* is the reduction of the amount of information stored for a certain signal event. Skimming reduces the number of events but not their size on storage, while slimming reduces the event size but does not reject events.
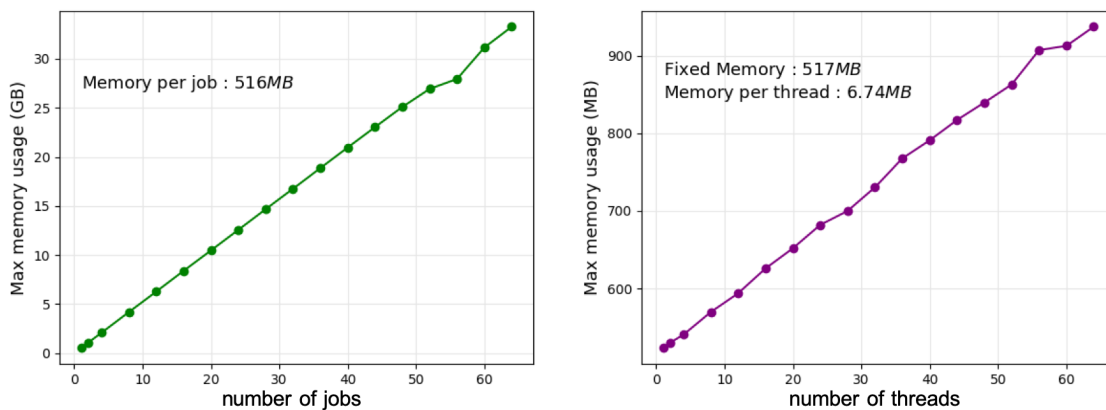
**Figure 113.** Memory consumption of a prototype HLT1 application when run in (left) multijob and (right) multithread modes. The application was run on 3000 events per thread on a reference server node with 20 physical cores and a factor 2 hyper-threading. Note that the *y*-axis scale of the right plot does not start at zero.

when executing many single-threaded jobs or one multithreaded job. Tests were run on a reference server node.[85] The memory increase is about 0.5 GB per job in the former case, while it is about two orders of magnitude smaller (6.7 MB/thread) in the latter case. At large number of processes or threads, the memory usage in the multithread approach was found to be a factor of about 40 smaller than the usage in the multiprocess approach.

In addition to the introduction of multithreading and to the improvement of the core software framework, the performance of the LHCb software has been optimised by following complementary approaches and techniques. The software stack has been modernised by using the latest C++ versions, introducing a full code review and suppressing code that was no longer used. The computationally intensive parts of the code have been adapted, and data structures have been reviewed, to exploit architectural features such as vector registers and nonuniform memory access (NUMA) domains, and to effectively use memory caches. In addition to these purely computing-related aspects, algorithmic improvements have been made, in some cases by completely changing the strategy of the most time-consuming algorithms. In these cases, it has been carefully verified that the performance of the relevant data analyses was not affected. An example, which summarises the improvements on software performance due to the points mentioned above, is shown in figure 114. In this figure, the evolution of the throughput of the HLT1 application between the autumn of 2018 and summer of 2019 is shown, as measured on a reference server by using simulated minimum bias events in nominal upgrade data-taking conditions. The key changes in the reconstruction algorithms during this period are colour-coded and described in the legend. The throughput improvement due to the introduction of single-instruction multiple-data (SIMD) instructions and data structures suitable to be used on vector registers are clearly visible.[86]

---

[85]Equipped with two Intel™ Xeon E5-2630 CPUs.

[86]Although eventually HLT1 has been implemented on GPUs, the example is illustrative of the improvements introduced by the software upgrade concepts that are implemented in CPU-based HLT2 and offline analysis applications as well as on GPU applications.
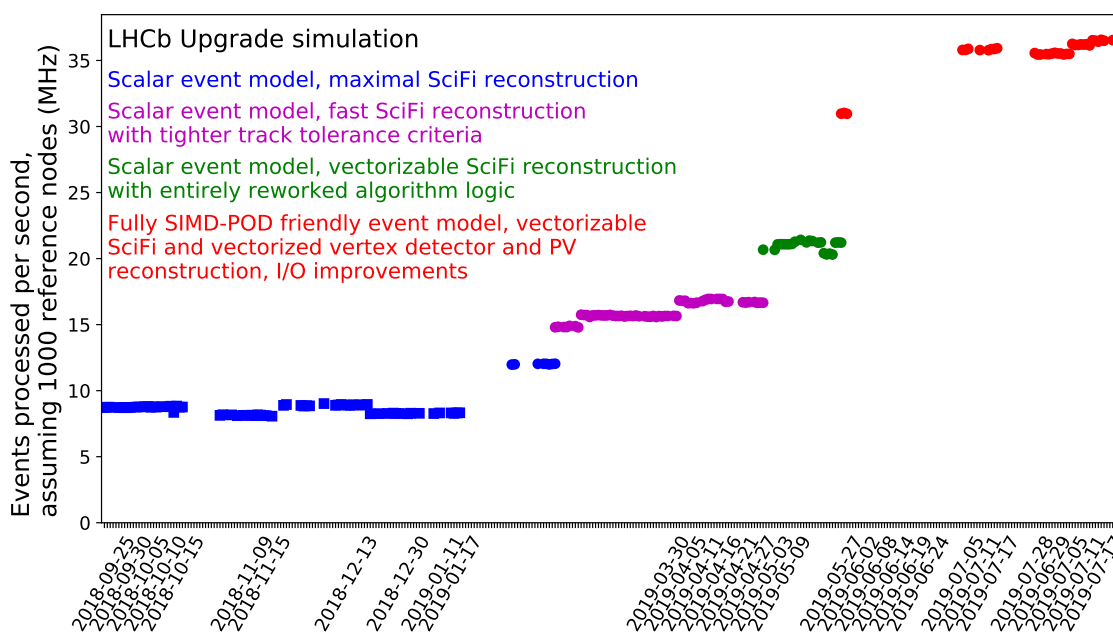
**Figure 114.** Evolution of the throughput of the CPU-based HLT1 prototype application between autumn 2018 and summer 2019, as measured on a reference server.

## 12.2 Conditions database

Conditions data describe the information about the detector that is required for data processing (e.g. detector configuration, alignment or calibration constants, environmental parameters). Conditions data may have a fine granularity, down to the level of individual detector elements. The space of conditions data has a three-dimensional structure defined by a geographical location, the time evolution, and a versioning applied to the entire condition data set. LHCb uses a Conditions Database to keep track of the time-dependent, nonevent data required to process the data collected by the experiment. For many years, conditions have been stored in a set of SQLite files that were accessed by means of COOL [203], until a Git-based solution was developed and commissioned in 2017. The Conditions Database is stored in a GIT repository hosted on the CERN GITLAB instance. Conditions data are accessed by a library, which also adds the time dimension (i.e. the interval of validity for a given condition) to the filename and versioning dimensions provided by a GIT repository. The YAML standard is used to persist condition data. Conditions values in the Conditions Database are used in different contexts and have to support different use cases, in particular simulation and real data processing. In simulation, the response of the detector for fixed sets of conditions is exploited. In real data processing, the time evolution of conditions is followed, with the GIT versioning of data used to track the evolution of the response of subdetectors, or improved alignment algorithms. Conditions data can be generated automatically in the online environment, for example calibration and alignment parameters or temperature and pressure values retrieved from probes, or produced manually for example in offline workflows that require expert analysis before inclusion in the database. In the former case, conditions are automatically pushed to the GITLAB project and published to the CERNVM file system (CVMFS) [204–206], while the offline conditions have to be proposed in the form of GITLAB merge requests, to undergo a review before being integrated and published.

## 12.3 Detector geometry description

The description of the geometry of the LHCb detector serves a wide range of applications, from simulation to detector alignment, visualisation, and computation of material budget. It must fulfil the LHCb needs in terms of flexibility, management, integration with the conditions database, and speed of navigation between detector elements. The Detector Description for HEP (or DD4HEP) toolkit [207] is used for the LHCb detector description. It replaces an XML-based framework [208], designed and implemented within the collaboration, but not optimised for modern computing architectures and no longer maintained. The DD4HEP toolkit builds on the experience gathered by LHC experiments and aims to bind the existing tools for detector description, simulation and visualisation to produce a consistent toolkit. Its structure is modular, with a generic detector description model as its core, and a set of nonmandatory extensions or plugins that can be used when needed. The object model and visualisation are provided by the ROOT framework. For simulations, the Geant4 toolkit is used.

## 12.4 Software infrastructure

The LHCb software stack depends on many software packages developed at CERN, in other scientific institutions, or in industry. Many of these packages are open source. External packages must be tracked, versioned, compiled and distributed. LHCb uses the *LCG releases* [209], which are prepared by the EP-SFT groups at CERN, for packages that are common among the LHC experiments. These releases are prepared using the LCGCmake [210] tool.

As many developers are involved in the development of the millions of lines needed for the experiment framework, a strong version control system and good practices are crucial to ensure the quality of the software. The LHCb code base is split into several projects, versioned and managed independently, each having a distinct goal. Each of these projects is managed using the Git [211] version control system. This allows keeping the full history of all code changes. Code reviews are prevalent in the software industry, as a way to improve code quality and harmonise development practices across organisations. LHCb Git projects are hosted on the CERN GitLab server, which also provides features allowing better collaboration between the members of the teams. New features to the code are peer-reviewed before being merged into the main code base. The projects are organised around the Jira [212] task management and the GitLab issue management systems, as deployed by the CERN IT department. This allows the developers within the LHCb collaboration to follow the evolution of the projects and collaborate in a constructive manner.

A Software Configuration Database has been built to track all projects and their dependencies [213] using the Neo4j graph database [214]. This information is crucial to the management of the software in the short term, but is also necessary to identify the software needed to continue analysing the LHCb experiment data in the long term.

In order to ensure the quality of the software produced by the developers, automatic builds of the software are performed, as described in ref. [215]. This infrastructure relies on the industry-tested Jenkins [216] automation server as deployed in the CERN IT OpenShift service. The build nodes are provided by the CERN IT department in the form of OpenStack [217] virtual machines. Results are gathered and displayed on a custom made web application [218].

Unit tests are run straight after the builds and the results are published to the LHCb Dashboard. Integration tests requiring more resources are run using LHCbPR, the LHCb performance and regression testing service [219]. LHCbPR is responsible for systematically running regression tests,
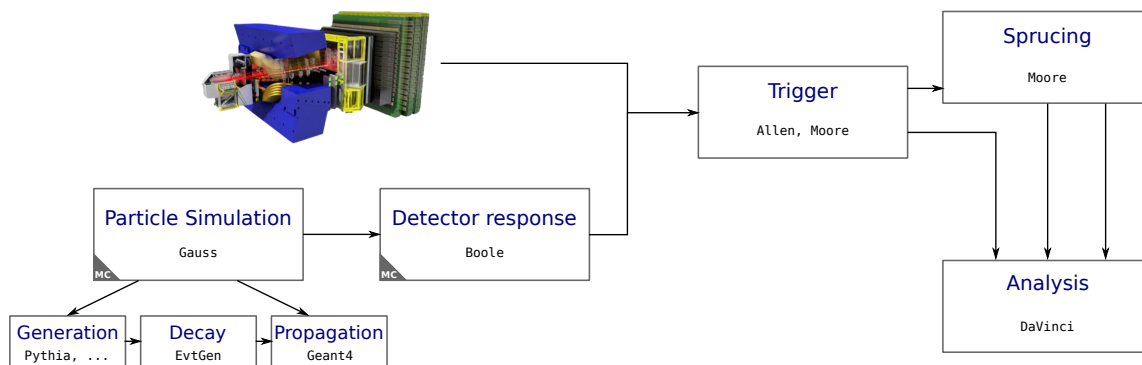
**Figure 115.** Schematic representation of the LHCb upgrade data flowand the related LHCb application, with an emphasis on simulation.

collecting and comparing results of these tests so that any changes between different setups can be easily observed. The framework is based on a microservice architecture where a project is broken into loosely coupled modules communicating with each other through APIs. The test service requests from LHCbPR information on how to run tests, then runs them and finally saves the results back to LHCbPR. Users have the ability to retrieve and analyse these test results.

The LHCb software is distributed using the CERN software deployment service CVMFS. Private installations of the software stack, as well as full management of the installed packages, are also possible. In this case, the applications are packaged in the RPM package manager [220] format, which allows specifying dependencies for applications relying on external packages (e.g. an installation of the analysis package relies on more than a hundred different packages).

The LHCb software stack is supported on 64-bit Intel architecture. Other architectures (ARM, IBM Power, GPUs) are handled by the LHCb build and release infrastructure, providing that the operating system used by LHCb can be installed, that packages using LCGCMAKE are ported and released, and that nodes for that architecture can be managed by the experiment's instance of JENKINS.

The LHCb software is preserved in the long term by the same tools used in the LHCb development process, more specifically the version control system as well as the CVMFS repositories, virtual machines and containers. Databases containing information about the software artefacts, their dependencies and their use in physics analysis ensure that applications can be rebuilt and rerun if necessary.

Several collaborative tools are used to coordinate the efforts and increase the proficiency of software developers. Rather than developing ad hoc solutions, the general strategy is to monitor the available tools and trends and adapt them to the specific use cases of LHCb.

## 12.5 Simulation

Monte Carlo samples have been essential for the detector design and preparation of the data processing and will be instrumental to the physics analysis of the LHCb upgrade. In order to define, tune and validate the reconstruction and selection algorithms, Monte Carlo samples are processed through a data flow identical to that of real data, as shown in figure 115. Simulation is followed by the online HLT and offline data processing described in sections 11 and 12.6 exactly as for real data.

The software to generate simulated events in LHCb is conveniently encapsulated in two separate applications. The first is the GAUSS package, responsible for the event generation and simulation of

particle interactions in the detector volumes which result in energy deposits or *hits* in the subdetector sensitive elements. The second is the BOOLE package, in charge of modelling the detector and readout electronics response by converting the hits into specific subdetector signals. This broad modularity is essential to profit from the experience continuously gained during the commissioning and operation of the subdetectors.

The generator-level information is transparently propagated through all the data processing steps to allow detector or physics performance studies comparing reconstructed and *true* (i.e. as produced by the physics generator) quantities.[87]

Like all other LHCb applications, GAUSS and BOOLE are built on the GAUDI core software framework and were heavily modernised with the adoption of parallelisation and multithreading and by implementing a new geometry description.

In addition, a new GAUDI-based experiment-independent simulation framework, named GAUSSINO [221, 222], was introduced to decouple widely used simulation software packages like GEANT4 or physics generators from LHCb-specific developments.

Simulation is the main consumer of LHCb computing resources as demonstrated by the fact that about 80% of the CPU resources made available to the collaboration during Run 2 were employed to generate, simulate and process Monte Carlo events. Therefore, to be ready to cope with the rapidly growing demand for simulated events due to the increase in LHCb integrated luminosity, code optimisation and modernisation has been complemented with the adoption of approximated simulation methods, referred to as *fast* and *ultra-fast simulations* (see also section 12.8.3). On a similar line, developments aiming at reducing the storage resources for simulated events were also pursued, as discussed in section 12.5.3

The software and physics performance of the simulation suite is monitored exploiting the software infrastructure described in section 12.4.

### 12.5.1 Gauss

The GAUSS application facilitates modelling of the physics processes occurring in the $pp$ collisions and takes care of the transport of the resulting particles through the experimental apparatus, including their interactions with the magnetic field and the detector and infrastructure material.

GAUSS has been extensively used by LHCb since the early development phases of the experiment and has undergone various changes to support evolving needs [101].

Figure 116 presents an overview of the structure of the GAUSS application. The simulation process consists of two subsequent phases. In the *generation* phase the physics process is obtained by dedicated generators, such as PYTHIA8 [26, 223] and EVTGEN [224]. In the *simulation* phase the generated particles are transported through the experimental apparatus, relying on the GEANT4 toolkit, on resampling techniques or on custom parametrisations, finally providing particles, vertices and energy deposits (*hits*) in the LHCb event data format. The HepMC [226] format is used to transfer information between different tasks within the generation phase and to pass the generated particles to the simulation step.

The design of GAUSS is based on GAUDI. Selection of components and steering of processing phases is achieved through a unique and coherent configuration system which represents a fundamental

---

[87]Physics generators are dedicated software packages that statistically generate particle four-momenta and decays based on first-principle physics quantities such as cross sections and branching ratios, using Monte Carlo techniques.
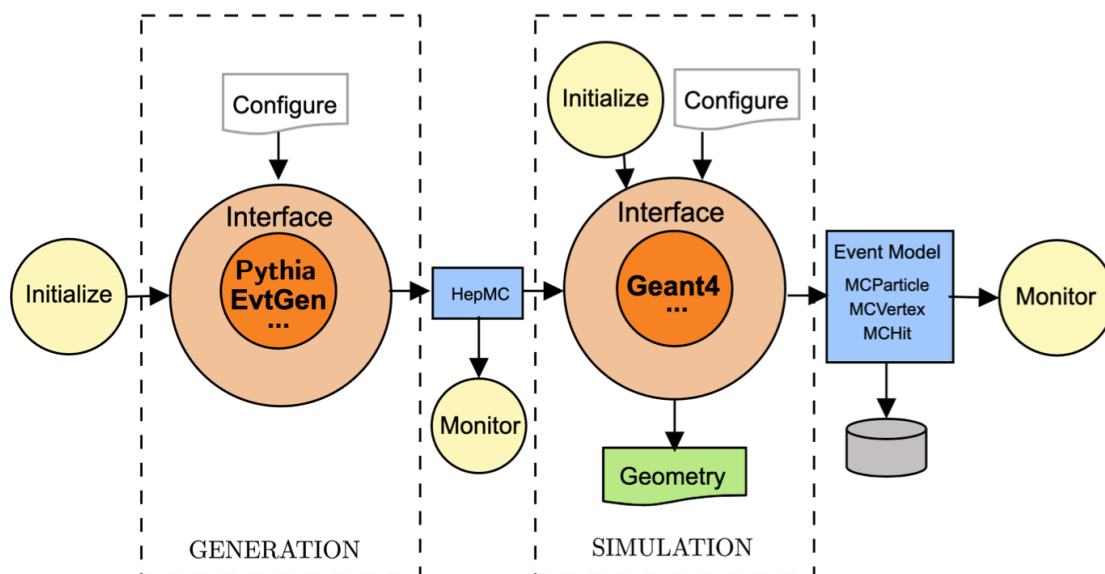
**Figure 116.** Schematic structure of the Gauss application. Reproduced from [223]. © 2022 IOP Publishing Ltd. CC BY 3.0.
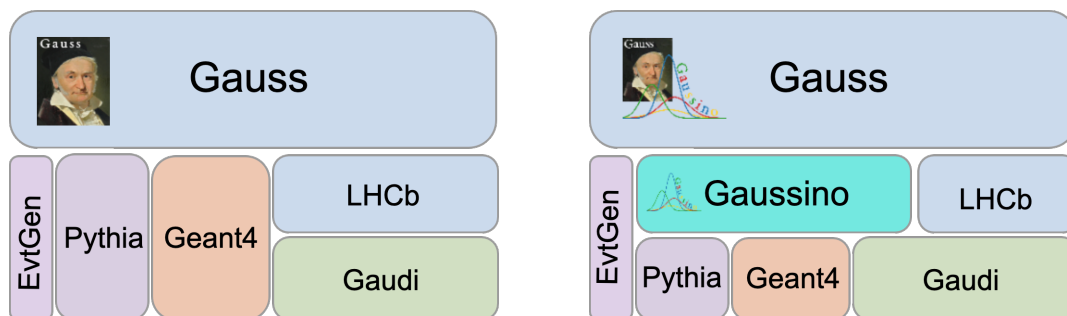


**Figure 117.** Graphical representation of the dependencies in the simulation software stack in (left) Run 1-2 and (right) in the upgrade, where the experiment agnostic package Gaussino decouples Gauss from Pythia and Geant4. Reproduced with permission from [225].

building block of the Gauss application. The modular architecture of Gauss makes it an excellent candidate from which to model a Gaudi-based core simulation framework for future experiments.

The experiment-agnostic simulation core features were packaged in Gaussino, providing interfaces to widely used packages such as Pythia or Geant4. Gauss is built on Gaussino and adds the LHCb-specific simulation functionalities as depicted in figure 117.

To be compatible with the multithreading model adopted for the LHCb upgrade, Gauss must be compatible with a multithreaded scheduling. This required the migration from the in-house geometry description software package used for Run 1 and Run 2 to the DD4HEP package [207]. Multithreading is handled by Gaussino which is also equipped with a general interface to steer the transfer of geometry information to Geant4.

The excellent modularity introduced through Gaussino and the upgraded Gauss allows also to easily simulate Run 1-2 events thanks to its support for the legacy geometry and data persistence.

**Figure 118.** Throughput and memory scaling for the generation of $pp$ collisions producing at least a $D^0$ meson with beam conditions as found in the 2016 data-taking period in LHCb. Shown are the curves for a shared (P8) and a thread-local (P8MT) interface to Pythia8, followed by the Geant4-based simulation. The contribution of the simulation phase, as obtained by reading the generated events from file ("Sim only"), is also shown. Reproduced with permission from [227].

### 12.5.2 The new Gaussino experiment-agnostic core simulation framework

The migration of the LHCb simulation to a multithreaded computing model implied the development of a general thread-safe interface with external simulation packages, such as Geant4 or physics generators, a task which was one of the main drivers of the development of Gaussino, explicitly designed for a much wider use than for LHCb only.

Gauss and Gaussino have an identical structure, shown in figure 116. The same architectural design was also retained, but using more modern Gaudi features in the implementation.

The Gaussino generator phase is essentially derived from that of Gauss [223], which was repackaged extracting the parts not specific to LHCb. Gaussino provides a thread-locking infrastructure to encapsulate and protect the execution of external tools not designed for multithreading, such as those including FORTRAN dependencies. Thread-safe libraries, such as Pythia8, can also be executed in multiple independent thread-local instances achieving faster execution at the expenses of a larger memory footprint.

The simulation phase, closely tied to Geant4, required a full redesign of the interfaces to various components to make them more experiment-independent and compatible with Geant4 multithreading. The same design choice taken by ATLAS was adopted to make the different Gaudi and Geant4 concurrent models work together.

As an example of the performance of this approach, the evolution of the memory occupation and event throughput as a function of the number of threads involved in the simulation of $pp$ collision producing at least a $D^0$ meson [227] is reported in figure 118. While the throughput increases almost proportionally with the number of threads, a large fraction of memory is shared, enabling the adoption of high-performance computing resources with several tens of cores to generate Monte Carlo samples.

Gaussino provides an additional interface to steer the interaction with Geant4 giving the possibility of replacing its detailed description of physics processes with fast simulation models for a specific detector. This interface was developed with the aim to minimise the work needed to implement fast simulation models for LHCb and facilitate their integration in Gauss. A facility to ease the production of training data sets for fast simulation models is also available in Gaussino.

### 12.5.3 Filtered, fast and ultra-fast simulations

While crucial to optimally adapt the LHCb simulation to modern computing architectures, the software modernisation effort described above is not sufficient alone to match the amount of necessary simulated events with the computing resources pledged to LHCb. Therefore, considerable effort has been devoted to the development of simulation techniques and technologies with a lesser impact on the computing and storage resources. Some of the developed strategies have already been deployed and widely adopted during Run 2, enabling a 70% reduction of the average computing power and a 60% reduction of the average storage occupation per simulated event [228].

In particular, *filtered* Monte Carlo simulation and a new output format were introduced during Run 2. In a filtered simulation, produced events are rejected by analysis-specific criteria before storing them. By the end of Run 2, only 13.9% of the simulated events was preserved, drastically reducing the impact of simulation on the storage resources [229]. Filtered Monte Carlo productions are used also for the upgrade, where the analysis-specific criteria are defined by the HLT or by the sprucing selections (section 12.6). The new output format, with selectively persisted event information, complemented by the corresponding Monte Carlo generator information, well matches the Turbo mechanism (see section 11.4.2) and is expected to be used for the majority of simulated samples in most cases in combination with Monte Carlo filtering.

Once a careful software optimisation is achieved, to further reduce the computing resources needed for event simulation it is necessary to compromise on the accuracy of simulated event features by replacing the detailed simulation based on GEANT4 physics models with fast simulation techniques based on either resampling methods or parametrisations of the detector response. The choice of the features on which a degradation of the accuracy is acceptable, the level of reliability expected from a simulated sample, as well as the statistical precision needed or, correspondingly, the number of events to simulate, depend on the specific needs of the data analysis and no unique solution can make optimal usage of the computing resources. A palette of simulation options spanning from GEANT4-based *detailed simulation* to fully parametric options has thus been developed to efficiently cover as many use cases as possible [230].

The most widely adopted fast-simulation option (see for example refs. [231–234]) is implemented in the REDECAY package [235]. Once a signal process is identified, for example the production of a heavy hadron, multiple random instances of its decay products are propagated through the detector while the detector response to the rest of the event is computed only once and superposed to the decay products of the signal particle. An overall decrease of the computing power per simulated event by a factor of 10 to 20 is achieved using the REDECAY technique, introducing an effect of event-to-event correlation which is only relevant for a minor fraction of analyses.

To study detector-induced effects that can be considered independent of the overall event particle multiplicity, it is often useful to simulate large samples of a specific process, e.g. the decay of a heavy-flavoured hadron. In these cases, the full simulation of the $pp$ collision can be avoided by generating the heavy hadrons according to parametrised spectra of their kinematic variables. These simulations, referred to as *particle guns*, are widely adopted for example to study background contributions and trigger effects in charm physics, or to analyse specific detector effects such as charge detection asymmetry [236]. Using particle guns, an overall increase in speed by a factor of 50 with respect to a detailed simulation is achieved [237].

A third fast-simulation option already deployed in production relies on a simplified geometry where entire subdetectors are removed, possibly complementing the simulation by sampling the
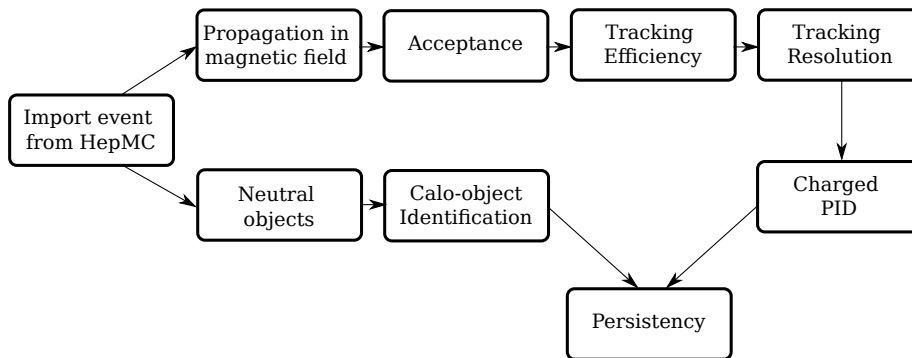
**Figure 119.** A flow-chart representing the LAMARR project as a pipeline of parametrisations.

nonsimulated features from parametrised distributions at analysis level. For example, *tracker-only* simulations can save up to 90% of the computing power per event by skipping the simulation of the optical photons in the RICH detectors and of the electromagnetic and hadronic showers in the calorimeters. The PID information is then added *a posteriori* using data-driven techniques originally developed to calibrate the simulation on unbiased data sets [199].

However, in some cases, partial detector simulations may require ad hoc solutions, for example in hardware trigger emulation. A more general solution, with similar CPU speedup, is obtained by replacing the simulation of time consuming processes with specific parametrisations obtained from dedicated simulated samples and deployed as part of the simulation software stack. In particular, techniques based on resampling from hit libraries [238] and querying generative models [239] have been developed to speed up the simulation of the energy deposits in the calorimeter, taking advantage of novel features of the GEANT4 package designed to replace part of the propagation with custom statistical models and available via GAUSSINO as described earlier in this section.

Machine learning techniques are also widely adopted to build statistical models of the detector response [240] and make it feasible to define parametrisations of extreme complexity to reproduce the whole detector simulation and reconstruction procedure by combining simple random generators and deterministic formulas. Figure 119 depicts the extension of GAUSS devoted to such an ultra-fast simulation option, named LAMARR, designed as a pipeline of parametrisations, taking in input generator-level quantities and producing reconstructed, analysis-level variables on output. The LAMARR framework aims at easing the deployment of machine learning models in GAUSS, providing a common infrastructure to data preparation and persistence configuration. Models are developed in ROOT TMVA [241, 242], SCIKIT-LEARN [243] or KERAS [244] and deployed as compiled shared objects, possibly relying on the SCIKINC package [245].

Original machine learning models were developed to enable training on real data, introducing statistical background subtractions in the training procedure [246], and adopted for the parametric simulation of the PID features obtained from the RICH, calorimeters and muon detectors, and their combinations [247].

### 12.5.4  Digitisation

Except for LAMARR and its ultra-fast simulation approach, the output of simulation applications must be processed by BOOLE to simulate the detector response (*digitisation*).

**Figure 120.** Offline data flow. Reproduced with permission from [180].

Boole converts the simulated hits into the same format as obtained from the DAQ of each subdetector, so that they can be processed by the common reconstruction and selection applications. Firstly, Boole simulates the subdetector technology and electronics response, including imperfections such as noise, cross-talk and dead channels. Then, the output of this stage is packed in *banks* which are passed to the event builder emulator to produce a raw data buffer identical to the output of the LHCb DAQ chain. Information of each hit digitisation history is preserved in the Boole output to allow detector and physics performance studies. The Boole application has been completely overhauled to match the upgraded DAQ and detector technologies.

### 12.5.5 Monte Carlo production

A correct and efficient production of Monte Carlo samples requires a well developed strategy and accurate prioritisation to match the physics analysis needs [229]. To maximise the CPU usage efficiency, the Monte Carlo production system has been upgraded to enable continuous integration tests in order to verify the consistency of simulation requests before submitting them to the distributed computing system. Simulation quality assurance tools were also improved and include the monitoring of the event simulation process by checking their physics output.

### 12.6 Offline data processing and analysis

The offline data processing flow can be seen in figure 120 and is detailed in this subsection.

### 12.6.1 Sprucing

The offline data processing is globally known as `sprucing`. The `sprucing` code base uses the same application as HLT2, namely the Moore application. Furthermore, the same algorithms and tools are shared between HLT2, `sprucing` and the offline analysis software project DaVinci [248] (see section 12.6.3), namely the functor based selection and combinatorial algorithms.

The `sprucing` performs three main functions: it applies further selections to data saved into the FULL stream (data skimming); it enables the tuning of the amount of event information to be persisted in the final output files through the `Turbo` mechanism (data slimming); it streams the data into a number of different physics stream files and create file summary records (FSR) that store metadata about the file content in the output Root files (see section 12.8.1. Data skimming and slimming can also be combined as needed. ' A typical example of data skimming is represented by events selected by inclusive topological HLT2 trigger selections [181], which were instrumental in Run 1 and Run 2. These HLT2 selections persist the full event information into the FULL stream. These events are further processed by exclusive `sprucing` algorithms that perform again the particle reconstruction and apply further selections to reduce the data volume to be saved to disk (data skimming). The same procedure can be applied to specific samples, such as data sets used to derive data driven calibrations in particle identification algorithms

Data in the FULL stream can also be used to tune the amount of event information that the selective persistence of HLT selections saves on output, in view of a future implementation in HLT2, reducing the size of the events saved on disk for further physics analysis (data slimming).

Normally, data which are saved in the `Turbo` stream by HLT2 are not further processed except for possible conversion to the final offline data format, creation of FSRs and streaming (*pass-through* `sprucing`). These steps correspond to the tasks performed in Run 2 by the Tesla application [177].

Irrespective of the intermediate processing (skimming, slimming or pass-through), data from both FULL and `Turbo` streams are eventually distributed into physics *streams* optimised to allow analysts to access reduced data sets categorised by physics topic. This optimises data processing time for a specific analysis as well as the use of computing resources such as disk access.

Compared to HLT2, offline `sprucing` benefits from less strict limits on CPU time for selection algorithms. While no recalibration and no rerunning of the pattern recognition is planned, the `sprucing` selections will have more time for example for the analysis of complex cascade or many-particle final state decays where the number of track combinations to be tested increases very rapidly. In addition, offline `sprucing` allows detailed analysis of physics topics where the full event information has to be taken into account such as electroweak physics or jet reconstruction.

Both exclusive and pass-through `sprucing` will run concurrently with data taking while global re`sprucing` campaigns will take place in end-of-year shutdowns whereby data will be staged from tape.

### 12.6.2 Distributed analysis productions

The output of the `Turbo` stream and of the `sprucing` is split into multiple streams which are directly accessible to analysts. In the Run 1-2 (known as *legacy*) data model these data sets would be processed by submitting user jobs to LHCbDirac that filter one of these streams to select physics-quantities of candidates for a specific analysis, typically resulting in a reduction in data volume of $O\left(10^3\right)$ (see section 12.7 for all details on distributed analysis). While this model works well for smaller data sets, scaling has been problematic with legacy analyses requiring many thousands of jobs. This causes the majority of analysts to be affected by site downtime, infrastructure instabilities and other distributed computing issues. These problems are compounded by the imperative nature of user jobs where each one has exactly specified input data and cannot be adjusted to adapt to current grid conditions. To deal with these issues and with the foreseen much larger data volume to be processed

by analysts in Run 3 and Run 4, a new strategy for analysis data processing has been developed, based on centralised *analysis productions*.

Analysis productions are an extension of the LHCbDIRAC transformation system, which has been primarily used for the centralised processing of LHCb data and simulation. Productions are submitted declaratively by providing the GAUDI configuration and bookkeeping query for the input data, which enables LHCbDIRAC to automatically handle failures and adjust the way in which files are grouped. Information about productions and the provenance of files is permanently stored in the LHCbDIRAC bookkeeping, enabling high quality analysis preservation and additional safety checks to be performed. Interactive analysis work can directly interact with LHCbDIRAC to obtain the location of the analysis production output data, thereby reducing the need to copy data manually and further supporting analysis preservation efforts.

Good testing of productions is essential as invalid productions have the potential to waste computing resources and cause instability in LHCbDIRAC itself. To provide assurance that GAUDI configurations provided by analysts are correct, extensive tests are run in GITLAB continuous integration prior to submitting the productions and the results of these tests are summarised on a dedicated website.

### 12.6.3 Offline analysis

While in the previous LHCb data processing flow, the trigger software was largely based on the previously designed offline analysis framework, the opposite is true for the upgraded LHCb data taking [249]. As much as possible, the new offline analysis framework is based on the software developments made for the trigger (see section 11.2.2). This approach avoids duplication of effort, profits from the major developments outlined in section 11 and guarantees a similar look-and-feel of the software to be used by analysts. Most importantly, all basic analysis building blocks are shared between the trigger, the `sprucing` and the offline analysis, thus guaranteeing that the same software is used to compute the same quantities online and offline.

Following this approach, the DAVINCI analysis software is descoped (compared to the versions designed to process the legacy data) to be only used for producing output tuples from input data (including simulation). Its core is an algorithm to produce *ntuples* that stores measured quantities using *functors* provided by the THOR framework [250] within MOORE.[88] These are the same functors that are used in the trigger or `sprucing` algorithms, thus ensuring a one-to-one correspondence between applied selection requirements and observables used offline. It should be emphasised here that since no full offline reconstruction is foreseen in the computing model, the input data objects (tracks, clusters, etc.) are identical online and offline [9, 194].

### 12.6.4 Data and analysis preservation

The centrally produced samples are preserved on the distributed computing infrastructure and catalogued in a dedicated LHCb bookkeeping system (see section 12.7). In 2020 LHCb ratified the CERN Open Data policy [251]. In accordance with the access policies outlined in this document, LHCb will make the output of the `Turbo` selections as well as the output of the `sprucing` available to the public through the CERN Open Data portal. The software necessary to read these files will be preserved as CVMFS releases. In order to enable secure access for third parties to the replicas of

---

[88]Ntuples are collections of variables related to an event, typically stored in a format suitable for ROOT or PYTHON packages. Functors are C++ objects used in the LHCb code to return calculated kinematic variables.

the data stored on the grid, a web-based interface is under development, which will allow users to configure analysis production jobs with minimal knowledge of the LHCb software.

To facilitate flexibility and creativity in the end-user analysis, implementing a statistical interpretation of the filtered data, only minimal constraints are placed on the necessary analysis code. While commonly used libraries, such as RooFit, ScikitHEP [252], etc. are available through CVMFS and can be managed using the Conda package manager, analysts can write custom scripts and custom routines to best answer their specific analysis tasks. For each publication the respective code is preserved in the respective physics working group GitLab repositories. Intermediate files, in particular filtered ntuples, are usually stored on the EOS storage system at CERN.

## 12.7 Distributed computing

The distributed computing of LHCb is based on the Dirac interware [253], and its LHCbDirac extension. The Dirac project is developing interware to build and operate distributed computing systems. It provides a development framework and a rich set of services for both workload and data management tasks of large scientific communities. The LHCbDirac infrastructure relies on database backends and services. The databases are provided by the CERN/IT database infrastructure (*database on demand* [254]).[89] The services run on a dedicated computing infrastructure also provided by CERN/IT.

Dirac was started by LHCb as a project for accessing computing grid resources. In 2009, following interest of other communities, the project has been open sourced. Now, it is a truly open source project hosted on GitHub and released under the GPLv3 license. Within the following years, Dirac has been adopted by several experimental communities both inside and outside high energy physics, with different goals, intents, resources and workflows.

In order to accommodate different requirements, Dirac has been designed with extensibility in mind. The core project (Dirac) can be *horizontally* extended by adding projects that are independently versioned but nevertheless strongly interdependent, as they concur to form a Dirac release. All projects are hosted together on GitHub, share the same license, and are maintained by the same set of users. The Dirac core project is the glue that keeps the satellite projects together. It also depends on software not maintained by the Dirac consortium, which is collected in *externals*. Horizontal extensions include: WebAppDirac, the Dirac web portal [255]; RESTDirac [256], which extends some Dirac services by providing a representational state transfer (REST) interface; VMDirac [257], which allows to create virtual machines on clouds; COMDirac, which extends the Dirac user interface. The Pilot [258] project is instead independent from all the other ones. This scheme is illustrated in figure 121 The *vertical* extensibility of Dirac enables users and virtual organizations (VOs) to extend the functionalities of the basic projects, in order to provide specific functionalities.

LHCb uses fully the functionalities provided by Dirac, but has customised some of the Dirac systems by implementing the LHCbDirac [259] extensions. New `Bookkeeping` [260] and `Production Management` [261] systems have been created, both also providing GUI extensions within the LHCbWebAppDirac, (the LHCb extension of WebAppDirac). The Dirac Pilot project is extended within the LHCbPilot project.

It is essential that Dirac provides a transparent and uniform interface for VOs to access resources that are more and more heterogeneous which implies that IaaS (Infrastructure as a Service) and

---

[89]The CERN/IT database infrastructure is provided by Oracle[TM].
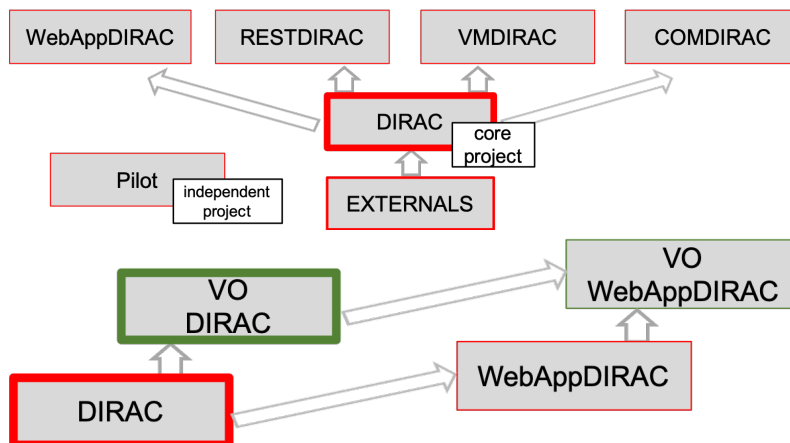
**Figure 121.** Diagram of a DIRAC release components. The concepts of (top) *horizontal* and (bottom) *vertical* extensibility are illustrated. Reproduced from [253]. © 2022 IOP Publishing Ltd. CC BY 3.0.

IaaC (Infrastructure as a Client) models must also be supported. This is realised in DIRAC by a generic pilot model [262], where a plugin mechanism enables easy adaptation on a wide range of computing resources [263–265], including cloud resources, high performance computing (HPC) centres and the servers of the LHCb online farm.

The DIRAC system scales in terms of traffic and data set growth, and maintainability. In terms of traffic growth, DIRAC closely follows the most modern architectural directives. Nevertheless technological updates such as the usage of message queues, PYTHON 3, multiprocessing and centralised logging systems add the required robustness to the system. In addition, the relational databases used in DIRAC are adequate to ensure full scalability for the data set growth.

System and software maintainability has also been taken in due consideration in the constant evolution of DIRAC, by implementing a proper monitoring, easily maintainable code with increasing functionality tests, the use of continuous integration tools and performance tests, better documentation and user support.

The DIRAC jobs are handled by the workload management system (WMS) and by the DIRAC data management system (DMS). Data sets are retrieved through a bookkeeping tool. The LHCb bookkeeping is a metadata and provenance catalogue used to record information about the data sets. In order to cope with the rapidly increasing data size, the main bookkeeping tables and indexes are partitioned. This allows to run more efficient queries using, for example, partition-wise joins.

## 12.8 Computing model

The new paradigm for the trigger selection process, described in the previous section, implies necessary changes in the offline computing model. Owing to the five-times higher instantaneous luminosity and higher foreseen trigger efficiency, the LHCb upgrade has a signal yield per time unit approximately ten times higher than that of the Run 1-2 LHCb experiment. The pileup also increases at upgrade instantaneous luminosity, resulting in an average event size increase by a factor of three. As a consequence a data volume larger by more than a factor of 30 is expected. A corresponding necessity to generate significantly larger samples of simulated events arises, as the number of events to simulate

is proportional to the integrated luminosity. The computing resource requirements are substantially mitigated by the novel real-time data processing model and by the massive use of fast simulation techniques, as discussed in the previous sections.

### 12.8.1 Data processing flow

Building on the experience developed during LHC Run 2, in the LHCb upgrade most of the activities related to data processing, such as event reconstruction and calibration and alignment of subdetectors, are performed online. The output produced by the HLT is stored on tape through three streams: FULL, Turbo and calibration (TurCal). The Turbo stream undergoes only minimal offline processing before being stored to disk. The FULL and TurCal streams instead requires further offline filtering.

The online reconstruction and trigger selection process execute the order of thousand trigger lines, each of which is associated to one of the three (Turbo, TurCal or FULL) streams that are subsequently stored offline. The offline processing of the Turbo stream, which comprises the bulk of the events, is performed by the TESLA application that converts the information from the raw format into ROOT I/O objects such as tracks, calorimeter clusters and particles, ready to be used for physics analysis, adds the luminosity information, and persists them on disk in the appropriate format. Turbo events are also classified into streams for an easier access at analysis stage.

Events in the FULL and TurCal streams are further processed offline by the sprucing application, which reduces the event size and performs a further event selection before storing the events on disk (data slimming and skimming, see section 12.6). Each of the sprucing selections is associated to a specific analysis, or group of closely related analyses, and the output information can be persisted at the appropriate level. An average event retention of 80% is obtained.

The offline data and processing flow is described in figure 122. The Turbo, FULL and TurCal streams are exported from the LHCb data centre. One copy of all data is stored at the CERN tape system. One additional copy of the FULL and TurCal raw data is stored at another Tier 1 tape system. All data are also copied to intermediate buffer disk storage. Data are then immediately processed by the appropriate stream dependent applications, as previously explained, and saved on disk. The Turbo data are simply streamed and put onto disk storage and as a second copy on a Tier 1 tape system. The first sprucing pass happens synchronously with data taking and can be prescaled if needed. Two replicas of the data will be kept on disk after this first processing pass. A second processing pass (resprucing, implementing updated selections and calibrations) is typically performed after the data taking period, usually during the LHC winter shutdown. When the second processing pass has been performed, the number of copies of the previous processing saved on disk is reduced. One copy of each sprucing pass is kept also on tape archive.

In the streaming scheme described above each user typically analyses a small fraction of the whole data set. In order to avoid bottlenecks due to each user chaotically running jobs on individual streams as desired, the data processing for user analysis is organised in centrally-managed productions (analysis productions), further described in section 12.6. In addition, users are allowed to submit jobs to offline resources, using the GANGA framework [266, 267], for analysis prototyping and testing purposes and other cases, such as running parametrised pseudo-experiment simulations, and performing fits or other further stages of the analysis.

The production of Monte Carlo events is described in detail in section 12.5. The simulated events are produced by two applications, GAUSS [101] and BOOLE [268], taking care of the event
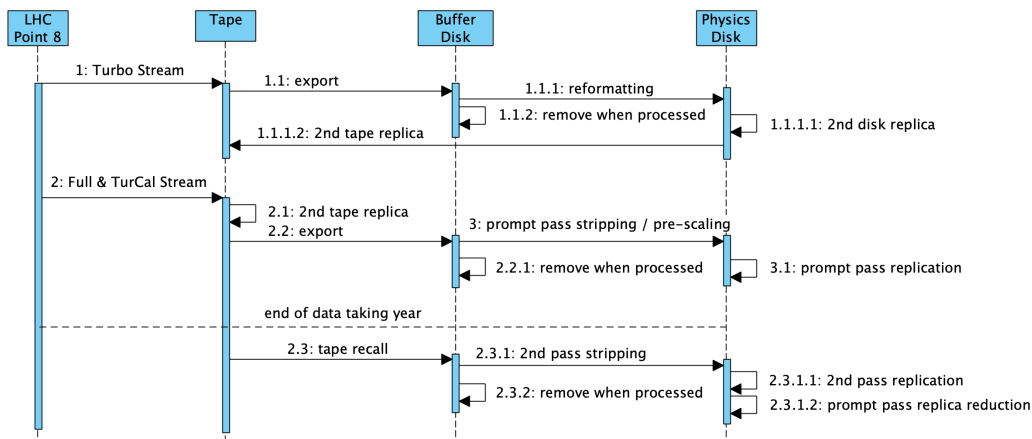
**Figure 122.** The LHCb offline data processing workflow. Reproduced from [194]. CC BY 4.0.

generation and propagation through the detector, and of the digitisation, respectively. Various Monte Carlo simulation techniques that are faster than the standard, detailed GEANT4-based simulation, are available, as described in section 12.5.

The simulation workflow in the LHCb Upgrade is very similar to the one used for Run 1 and Run 2. A number of steps are run in sequence. The intermediate files created at the end of each step are transient and deleted when no longer necessary. The only notable exception to this workflow is the fully parametric simulation where the data are saved directly in the form of high-level objects that are ready to be used in physics analysis. In this case, the digitisation, trigger emulation and event reconstruction steps are skipped.

### 12.8.2 Resource provisioning

The provisioning of resources for the computing infrastructure provided by the Worldwide LHC Computing Grid (WLCG) follows a pledging scheme where computing sites provide a dedicated amount of resources to the experiments. The pledged resources are based on requests submitted by the experiments for the forthcoming years and accompanying resource usage reports. Both documents are provided twice a year to the relevant funding bodies. The WLCG infrastructure is setup in tier levels. The Tiers used by LHCb are the Tier 0 at CERN, major Tier 1 sites in several countries and approximately ninety additional Tier 2 sites both in countries with Tier 1 centres and in other countries. CPU resources are provided on all Tier levels. Tape storage is only provided at Tier 0 and Tier 1 sites. Disk storage is provided at Tier 0, Tier 1 and a limited number of Tier 2 sites. Limiting the storage resources (tape and disk) to a restricted number of sites has proven to be a successful operational model.

Using the *mesh processing* paradigm [269], the so called Tier 2 *helper sites* are attached to one or more storage sites. Each helper site receives a payload, downloads the input data from the remote storage site, processes the files locally and subsequently uploads the output data again to the same storage site from which the input data were downloaded. This concept is used for data reprocessing campaigns to increase the throughput but also during prompt processing in case the Tier 0 and Tier 1 sites do not have enough resources to cope with the load. The *GAUDI federation* concept [270] is used to read input data from other than the initially foreseen storage sites. Within this scheme, an application is deployed with additional information on the location of all replicas of all needed input

data files. If the first priority copy of a file is not readable (for example because the file is corrupted or the disk storage is not available), the application searches over the network for remote replicas of the input file across the federation, and reads data from there. This concept is especially useful for user analysis files where multiple replicas can be available.

Data handling and data replication follow a *democratic* principle where data are replicated over all possible storage sites depending on the available space and the capacity of the corresponding site. For raw detector data, the smallest block to be replicated is represented by the files corresponding to one detector run. Derived data sets are also kept on the same storage site. In case of data replication, all descendant files from one run are also replicated. Intermediate files within a simulation job (typically executed on a Tier 2 site) are stored on a topologically close Tier 1 site and deleted after they have been processed by the corresponding application. The final simulated event files, ready for user analysis, are also uploaded to a topologically close Tier 1 site and then replicated following the democratic data replication policy. This principle has proven to be successful, as it has made the data set handling operations easier and allowed to optimise the load of applications using the distributed data.

In addition to the CPU resources pledged via the WLCG, LHCb also uses several additional computing resources in an opportunistic and/or ad hoc way. They come from two different sources: the LHCb online farm and opportunistic resources not owned by and not under the control of the experiment. The online farm is used by LHCb for offline data processes outside data-taking periods. It is possible in principle to run concurrently online and offline applications, thanks to the fact that the same hardware infrastructure is used in both cases, without any intermediate virtualisation layer, and to the implementation of a *fast stop* mechanism that allows to switch between offline and online usage within a time of about 1-2 hours. Opportunistic resources not owned by LHCb include HPC centres and resources hosted on WLCG sites that are not pledged to LHCb. In both cases, the main usage is for simulation, as it does not require input data. The amount of tasks that can be performed with these CPU resources is unpredictable as the priority given to the LHCb applications compared to other users is lower, and LHCb jobs are essentially used to fill otherwise unexploited CPU time.

### 12.8.3  Resource requirements

The production of simulated events dominates the offline CPU computing needs. The number of events to be simulated is estimated to be $\sim 5 \times 10^9 \, \text{fb}^{-1}$ of integrated luminosity per calendar year. The production of events simulated according to a given data-taking year extends normally up to six years afterwards. The amount of corresponding resource requirement is mitigated by exploiting faster simulation options. The storage needs are instead dominated by data and crucially depend on the HLT output bandwidth of 10 GB per live second of LHC. While the associated tape needs are not compressible, a mitigation is achieved for disk. A fraction of about 70% of triggered events are saved in the `Turbo` format. However, the majority (6.5 GB/s out of 10 GB/s) of the bandwidth is taken by the remaining 30% of events in the `FULL` and `TurCal` streams, where the entire event is persisted. The events in these two streams are therefore slimmed and filtered offline by the `sprucing` process, such that the total (logical) bandwidth to be saved on disk is only $\sim 30\%$ of the original. The extrapolated throughput to tape and disk for the three data streams are reported in table 13.

The impact of simulated events on storage requests is small, since all data produced during the intermediate production steps are not saved and the simulation output is persisted in a compressed format, thus achieving a size reduction per event of a factor up to twenty. In addition, analysis-dependent

**Table 13.** Extrapolated throughput to tape and disk for the FULL, Turbo and TurCal streams.

| Stream | Rate fraction | TAPE throughput (GB/s) | TAPE bandwidth fraction | DISK throughput (GB/s) | DISK bandwidth fraction |
|--------|---------------|------------------------|-------------------------|------------------------|-------------------------|
| FULL   | 26%           | 5.9                    | 59%                     | 0.8                    | 22%                     |
| Turbo  | 68%           | 2.5                    | 25%                     | 2.5                    | 72%                     |
| TurCal | 6%            | 1.6                    | 16%                     | 0.2                    | 6%                      |
| Total  | 100%          | 10.0                   | 100%                    | 3.5                    | 100%                    |

**Table 14.** Summary of the LHCb upgrade computing model requirements. Top section: main assumptions of the model. Bottom section: indicative resource requirements.

| Model assumptions | | | |
|---|---|---|---|
| $\mathcal{L}$ [cm$^{-2}$s$^{-1}$] | $2 \times 10^{33}$ | | |
| Pileup | 6 | | |
| Running time [s] | $5 \times 10^6$ ($2.5 \times 10^6$ in 2021) | | |
| Output bandwidth (GB/s) | 10 | | |
| Fraction of Turbo events | 73% | | |
| Ratio Turbo/FULL event size | 16.7% | | |
| Ratio full/fast/param. simulations | 40:40:20 | | |
| Data replicas on tape | 2 | | |
| Data replicas on disk | 2 (Turbo); 3 (FULL, TurCal) | | |
| Resource requirements | | | |
| WLCG Year | Disk (PB) | Tape (PB) | CPU (kHS06) |
| 2021 | 66  | 142 | 863   |
| 2022 | 111 | 243 | 1.579 |
| 2023 | 159 | 345 | 2.753 |
| 2024 | 165 | 348 | 3.467 |
| 2025 | 171 | 351 | 3.267 |

filtering criteria are also applied to reduce the number of events written on storage. The offline reprocessing of the FULL and TurCal streams requires a significant reading throughput from tape. The needed throughput is estimated by considering that reprocessing is done over a two-months period, with a provisional buffering space corresponding to two weeks of data staging. Half of data is staged at Tier 1 sites, the other half at CERN, 50% of which is then transferred at the Tier 1 sites for processing. For safety reasons, in general two copies of all data that are impossible to be regenerated are saved on tape. Therefore, two copies of the primary data sets, i.e. those originating from the online system, are stored on tape. An archive copy of the offline processed data for each of the three streams is also saved on tape. After the offline processing, two copies of the Turbo stream and up to three replicas of the FULL and TurCal streams will be saved on disk. A single copy of all the simulated events is kept on tape, while two copies of the most used simulated data sets ($\sim 30\%$ of the total) is stored on disk.

A summary of the computing model parameters and an indicative estimation of computing resource requirements is given in table 14.

## 13  Performance

The expected performance of the upgraded LHCb detector has been extensively studied with detailed simulations. Dedicated samples of specific physics channels have been generated and processed through the full LHCb detector simulation to study in particular the HLT1 reconstruction and selection efficiency. In the LHCb upgrade full-software trigger scheme, the physics selections are largely performed at HLT2 level; their efficiency will only be available when the specific analyses will be carried out and finalised during data taking.

### 13.1  Computational performance

The physics performance of the upgraded LHCb detector critically depends on the computational performance of its real-time software, which enables the relevant algorithms to be deployed in HLT1 and HLT2 as described in section 11. The performance is measured in Hertz, giving the number of events which can be processed each second on a representative processing unit. The results obtained with the HLT1 application running on a range of currently available GPU cards are shown in figure 123. The GPU card used by LHCb for the first data taking run is the A5000.[90] The figure also shows the performance of the HLT1 code on a modern CPU server which is relevant for running HLT1 in simulation on the grid where GPU resources are not necessarily available. The computational performance of HLT2 is shown in figure 124 for a reference CPU server which will be used for Run 3 data taking.[91] In both cases the computational performance, albeit evaluated on simulation, is comfortably adequate for the foreseen nominal Run 3 data taking conditions.



**Figure 123.** Throughput of the HLT1 application on a selected subset of current generation GPU cards and a representative modern CPU server. Reproduced with permission from [189].

---

[90]NVIDIA RTX A5000 graphics card.

[91]About 3500 HLT2 nodes are installed in the LHCb data centre.

**Figure 124.** Throughput of the HLT2 application and fraction of HLT2 resources used by different parts of the reconstruction and selections, measured on a representative HLT2 server. A total of 1111 selection algorithms were executed as part of this test. Reproduced with permission from [271].

## 13.2 Reconstruction performance

In the upgraded LHCb the event reconstruction is performed in real-time at the trigger level. The reconstruction efficiencies in HLT1 and HLT2 have been studied in great detail while designing the trigger algorithms. The main results are shown in the next sections.

### 13.2.1 Tracking performance

Performance figures are produced with simulated event samples of $B_s^0 \rightarrow J/\psi\phi$, $B_s^0 \rightarrow \phi\phi$, $B^0 \rightarrow K^{*0}e^+e^-$ and $D^+ \rightarrow K_S^0\pi^+$ (5000 events per magnet polarity), including the effect of VELO and SciFi Tracker radiation damage after $5\,\text{fb}^{-1}$, as well as a sample of $D^+ \rightarrow K_S^0\pi^+$ decays enriched in tracks originating outside the vertex detector.

The tracking performance in the baseline HLT2 reconstruction is shown in terms of reconstruction efficiency and corresponding fake-track reconstruction rate (*ghost rate*) for long and downstream tracks as a function of momentum $p$, transverse momentum $p_{\mathrm{T}}$, pseudo-rapidity $\eta$, and number of primary vertices. Figures 125 and 126 show the track reconstruction efficiency for long tracks originating from *B*-meson decays while the corresponding fake-track rate is shown in figure 127. The downstream-track reconstruction efficiency for particles originating in strange and *B* or *D* decays is shown in Figs 128 and 129, respectively; the corresponding fake-track rate is reported in figure 130. Finally the seeding reconstruction efficiency for particles originating in *B* decays and reconstructible[92] as long tracks is shown in figures 131 and 132, and the corresponding fake-track rate is displayed in figure 133.

---

[92]In the context of this section *reconstructible* means an object (track, vertex) or an event with visible activity in the detector that satisfies some minimum requirements to allow its reconstruction by software algorithms.

**Figure 125.** Long track reconstruction efficiency versus momentum $p$, transverse momentum $p_{\mathrm{T}}$, pseudo-rapidity $\eta$, and number of primary vertices for long reconstructible electrons (blue squares) and non-electron (black dots) particles from $B$ decays within $2 < \eta < 5$. Shaded histograms show the distributions of reconstructible particles. Reproduced with permission from [271].
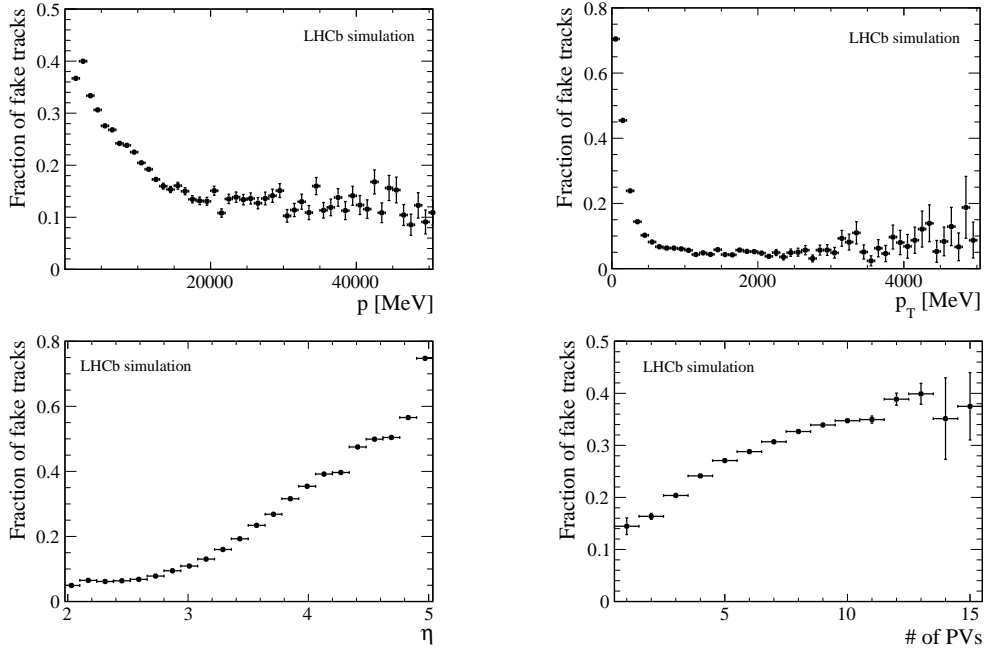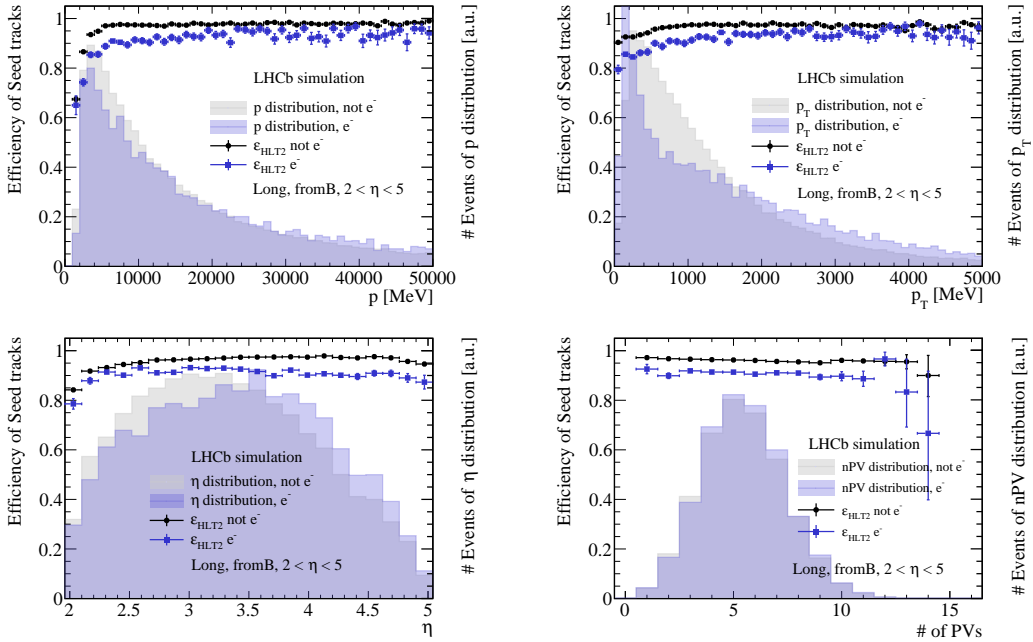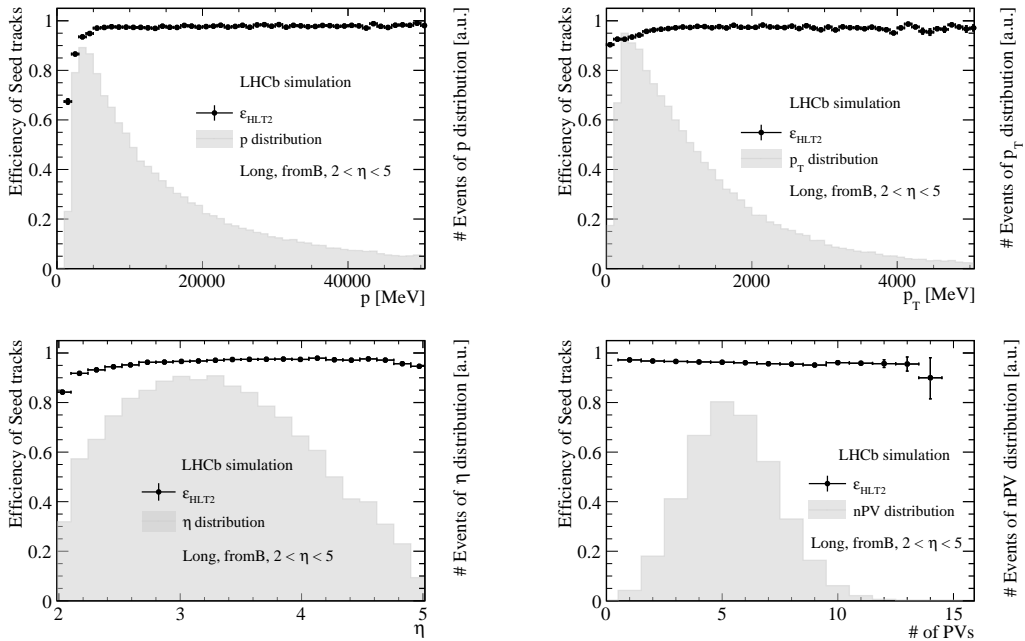


**Figure 126.** Long track reconstruction efficiency versus momentum $p$, transverse momentum $p_{\mathrm{T}}$, pseudo-rapidity $\eta$, and number of primary vertices for long reconstructible particles from $B$ decays within $2 < \eta < 5$. Shaded histograms show the distributions of reconstructible particles. Reproduced with permission from [271].

**Figure 127.** Ghost rate of long tracks reconstructed by the forward and match tracking algorithms as a function of momentum $p$, transverse momentum $p_T$, pseudo-rapidity $\eta$, and number of primary vertices. Reproduced with permission from [271].



**Figure 128.** Downstream track reconstruction efficiency versus momentum $p$, transverse momentum $p_T$, pseudo-rapidity $\eta$, and number of primary vertices for reconstructible particles from long-lived particle (marked as *strange* in the legend) decays within $2 < \eta < 5$ that have no hits in the VELO. Shaded histograms show the distributions of reconstructible particles. Reproduced with permission from [271].

**Figure 129.** Downstream track reconstruction efficiency versus momentum $p$, transverse momentum $p_T$, pseudo-rapidity $\eta$, and number of primary vertices for reconstructible particles from *B/D* decays within $2 < \eta < 5$ that have no hits in VELO. Shaded histograms show the distributions of reconstructible particles. Reproduced with permission from [271].



**Figure 130.** Ghost rate of downstream tracks reconstructed by the forward- and match-tracking algorithms as a function of momentum $p$, transverse momentum $p_T$, pseudo-rapidity $\eta$ and number of primary vertices. Reproduced with permission from [271].

**Figure 131.** Seeding track-reconstruction efficiency versus momentum $p$, transverse momentum $p_{\mathrm{T}}$, pseudo-rapidity $\eta$, and number of primary vertices for long reconstructible electrons (blue squares) and non-electron (black dots) particles within $2 < \eta < 5$. Shaded histograms show the distributions of reconstructible particles. Reproduced with permission from [271].



**Figure 132.** Seeding track reconstruction efficiency versus momentum $p$, transverse momentum $p_{\mathrm{T}}$, pseudo-rapidity $\eta$, and number of primary vertices for long reconstructible particles from $B$ decays within $2 < \eta < 5$. Shaded histograms show the distributions of reconstructible particles. Reproduced with permission from [271].

Figures 134 and 135 show the momentum and impact parameter (IP) resolution for tracks after the Kalman fit. Finally, the tracking efficiency is measured as a function of detector occupancy in lead-lead collisions, in order to understand the suitability of the upgraded LHCb detector for heavy ion data taking. This is shown in figure 136. A gradual degradation of performance is observed with occupancy but, encouragingly, there is no sharp edge where performance collapses.

A simultaneous data taking of beam-beam and beam-gas collisions can be envisaged due to the new SMOG design, with well separated gas-target and beam-crossing regions, and since the expected rate of beam-gas collisions is one order of magnitude smaller than the beam-beam rate. Due to the tight constraints of the online reconstruction framework discussed in section 11, the simultaneous data-taking must not spoil the track reconstruction performance in $pp$ collisions while ensuring good tracking efficiency for beam-gas events. Results of the studies carried out on simulated samples of stand-alone proton-helium ($p$He) collisions and with overlapping $pp$ and $p$He or $p$Ar are shown here. The upgrade luminosity conditions and one proton-gas collision per event confined in the SMOG cell region $z \in [-500, -300]$ mm according to the expected gas pressure into the storage cell are assumed. The track reconstruction efficiency for stand-alone $p$He and $pp$ collisions, and for overlapping $pp$ and $p$He or $p$Ar are shown in figure 137 as a function of the $z$ position of the primary vertex ($PV_z$), along with the corresponding fake track rate as a function of the track momentum $p$. The results are obtained with long tracks with $p > 3$ GeV and $p_T > 0.5$ GeV. The primary vertex reconstruction efficiency is also investigated with the same simulated event samples. The results are shown in figure 138. In the top plot the $v$ reconstruction efficiency is reported as a function of the $z$ coordinate. The distribution of the reconstructible $v$s (in arbitrary units) is also shown. In the bottom plot, the $v$ $z$ resolution is shown as a function of $z$. The time needed by the HLT1 application to process overlapping beam-gas and $pp$ events increases by about 1–3% with respect to $pp$ collision processing time. These studies demonstrate the feasibility of simultaneous running in $pp$ and fixed-target mode.

### 13.2.2 Calorimeter performance

In figure 139 the efficiency of reconstructing ECAL clusters from the energy deposited by photons is shown, while the position resolution of the reconstructed clusters is given in figure 140. These figures are produced with $B^0 \to K^{*0}\gamma$ simulation samples. Figure 141 shows the ECAL cluster position resolution for merged $\pi^0 \to \gamma\gamma$ decays from $B$ decays. The resolutions are shown separately for the three ECAL regions, which are characterised by differing sizes of the ECAL cells which affect in particular the position resolution.

The power of the ECAL variables to separate electrons from other charged particles is shown in figure 142 by the distribution of the ratio between ECAL energy EcalE and track momentum $p$, where the variable EcalE is the sum of energies of the ECAL cells intersecting the track extrapolation and those compatible with potential bremsstrahlung emissions. These bremsstrahlung emissions are determined by projecting the track direction before bending in the magnetic field to energy deposited in the ECAL. These plots are produced with $B^0 \to K^{*0}e^+e^-$ simulation samples.

### 13.3 Selection performance

To cover the current LHCb physics programme, $O(100)$ selections are deployed in HLT1 and $O(1000)$ in HLT2, with their number expected to increase significantly with the evolution of physics studies. The precise balance between efficiency and rate for each physics signature will only be established
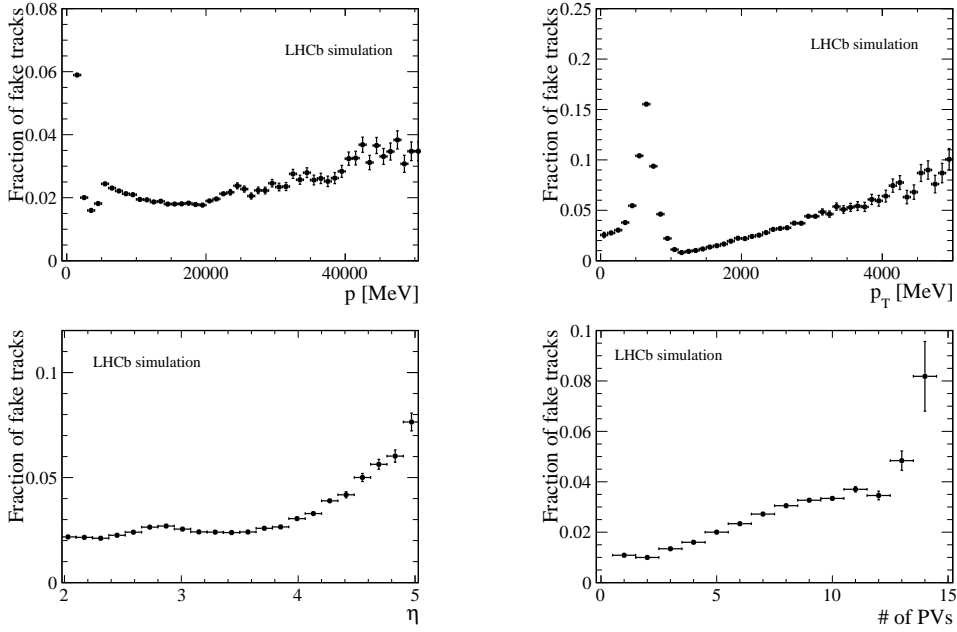
**Figure 133.** Ghost rate of standalone seeding tracks as a function of momentum $p$, transverse momentum $p_{\mathrm{T}}$, pseudo-rapidity $\eta$, and number of primary vertices. The prominent peak in the ghost rate at low transverse momentum (top-right panel) results from a combination of geometric and kinematic effects. Reproduced with permission from [271].
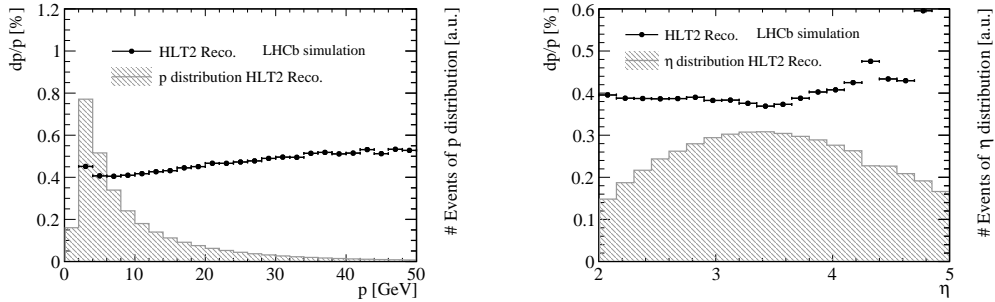


**Figure 134.** Relative resolution of the momentum of reconstructed tracks as a function of momentum $p$, and pseudo-rapidity $\eta$. Reproduced with permission from [271].



**Figure 135.** Resolution of the $x$ projection of the impact parameter, $\sigma_{\mathrm{IPx}}$ (left) and $\sigma_{\mathrm{IPy}}$ (right) as a function of the inverse of transverse momentum $1/p_{\mathrm{T}}$. A minimum bias sample is used for the IP resolution study. Reproduced with permission from [271].

**Figure 136.** Reconstruction efficiency of (left) VELO and (right) long tracks as a function of the occupancy in the vertex detector and SciFi Tracker, respectively. Reproduced with permission from [271].



**Figure 137.** Top: track reconstruction efficiency as a function of the $z$ position of the primary vertex ($PV_z$), for simulated samples with stand-alone (blue) $p$He and (green) $pp$, and overlapping (red) $pp + p$He and (orange) $pp + p$Ar collisions. The distribution of $PV_z$ for reconstructible $vs$ is also shown (shaded histogram, arbitrary units). Bottom: corresponding rate of fake reconstructed tracks as a function of track momentum $p$. Reproduced with permission from [271].

**Figure 138.** Primary vertex reconstruction (top) efficiency and (bottom) resolution as a function of the $z$ coordinate measured on simulated samples with stand-alone (green) $pp$, (blue) $p$He and overlapping (red) $pp + p$He and (orange) $pp + p$Ar collisions. Reproduced with permission from [271].

once the detector is commissioned and the backgrounds observed in data during this new high-pileup regime are understood. Nevertheless, the general selection performance is evaluated for certain archetypal HLT1 and HLT2 selections on a subset of representative signal topologies. In all cases the efficiency is calculated on events in which the signal topology of interest can be fully reconstructed inside the LHCb detector acceptance (i.e. factorising out the geometrical acceptance), but no other offline selection criteria are applied. This rather loose normalisation is used to illustrate the work which the collaboration will have to do in order to fully benefit from the removal of the first-level hardware trigger, which biased all signals to have large transverse momentum. In the case of HLT2 selections, the efficiencies are calculated relative to signal events passing HLT1 requirements, in order to factorise the performance of the first- and second-level triggers.

The performance of HLT1 inclusive selections — a single displaced track trigger and a displaced vertex trigger — is shown in figure 143. The performance of HLT1 inclusive muon selections is shown in figure 144. The performance of the HLT2 inclusive triggers is shown in figure 145. Finally, the performance of two representative exclusive HLT2 triggers is shown in figure 146. While the

**Figure 139.** ECAL cluster reconstruction efficiency versus energy $E$, transverse energy $E_{\mathrm{T}}$ and $x$ and $y$ position in the ECAL for reconstructible photons from $B^0 \to K^{*0}\gamma$ decays. Reproduced with permission from [271].



**Figure 140.** ECAL-cluster (left) $x$ position and (right) $y$ position resolution versus energy for reconstructible photons from $B^0 \to K^{*0}\gamma$ decays. Reproduced with permission from [271].



**Figure 141.** Merged $\pi^0$ (left) $x$ position and (right) $y$ position resolution versus energy for $\pi^0 \to \gamma\gamma$ from $B$ decays. Reproduced with permission from [271].

**Figure 142.** Main electron PID variables for the ECAL: distributions for signal and background separately for the variables (left) EcalE/$p$ and (right) matching $\chi^2$ of a bremsstrahlung cluster candidate to a track. The distributions of the bremsstrahlung matching $\chi^2$ are conditional on having a cluster candidate in a $3 \times 3$ cell grid around the bremsstrahlung track extrapolation. Reproduced with permission from [271].
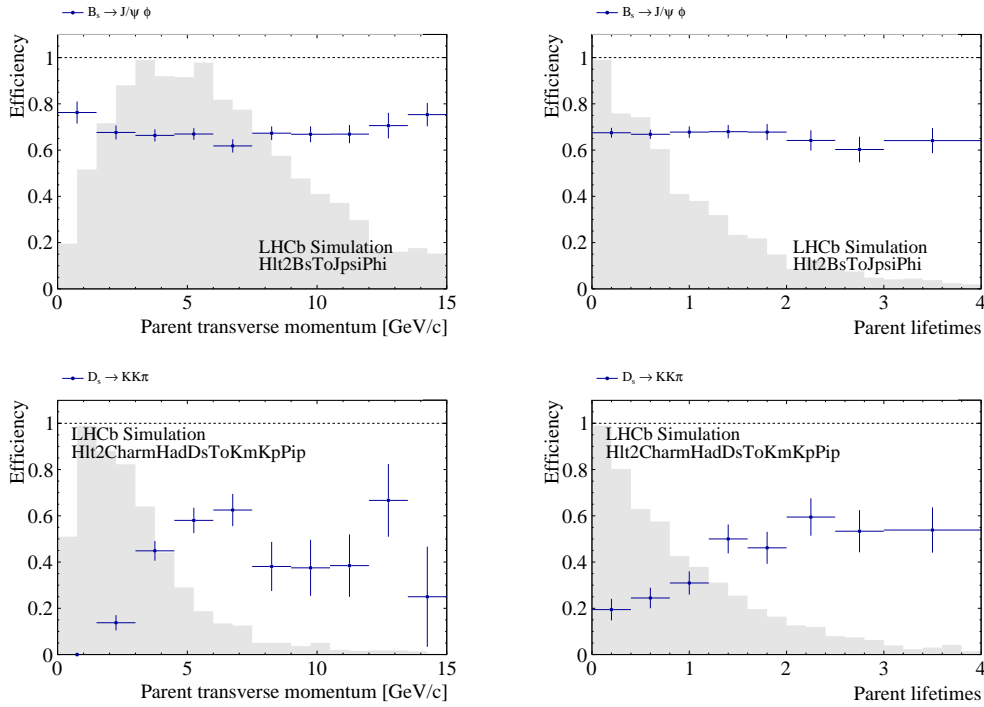


**Figure 143.** Performance of the HLT1 inclusive selections as a function of (left) parent-particle transverse momentum and (right) parent-particle decay time. The top row plots are the single-track selections, while the bottom row plots are the two-track displaced vertex selections. The signal topologies are indicated in the legend above each plot. The decay time plots are drawn such that the $x$ axis is binned in units of the lifetime for each hadron in its rest frame. Reproduced with permission from [271].

plotted efficiencies may appear low in many cases, this is because of the lack of offline criteria in the denominator. For the same reason the HLT2 efficiencies, where the denominator are events passing HLT1 conditions, are higher. Nevertheless, the examples of HLT2 single high-$p_T$ muon triggers and the exclusive $B_s^0 \rightarrow J/\psi \phi$ triggers show that exclusive triggers targeting specific well-defined decay chains can achieve excellent efficiencies even with respect to reconstructible decays.

**Figure 144.** Performance of the HLT1 muon selections. The signal topologies are indicated in the legend above each plot. In the top row the performance of the dimuon selections is plotted as a function of (left) parent-particle transverse momentum and (right) parent-particle decay time. The decay time plot is drawn such that the $x$ axis is binned in units of the lifetime for each hadron in its rest frame. In the bottom row the performance of the single high-$p_T$ muon selection is plotted as a function of parent transverse momentum. The shaded histograms indicate the distribution of the parent particle prior to any trigger selection. Reproduced with permission from [271].



**Figure 145.** Performance of the HLT2 inclusive selections. The signal topologies are indicated in the legend above each plot. In the top row the performance of the inclusive displaced-vertex selections is plotted as a function of (left) parent-particle transverse momentum and (right) parent-particle decay time. The decay-time plot is drawn such that the $x$ axis is binned in units of the lifetime for each hadron in its rest frame. In the bottom row the performance of the single high-$p_T$ muon selection is plotted as a function of parent transverse momentum. The shaded histograms indicate the distribution of the parent particle prior to any trigger selection. Reproduced with permission from [271].

**Figure 146.** Performance of example HLT2 exclusive selections as a function of (left) parent-particle transverse momentum and (right) parent-particle decay time. The signal topologies are indicated in the legend above each plot. The decay-time plots are drawn such that the *x* axis is binned in units of the lifetime for each hadron in its rest frame. The shaded histograms indicate the distribution of the parent particle prior to any trigger selection. Reproduced with permission from [271].

# 14  Summary

The LHCb upgraded experiment has been described including the detector, the online system, the all-software trigger and the software and computing infrastructure.

The upgrade of LHCb consists of: a tracking system including a new silicon-pixel vertex detector, a new silicon-strip tracker upstream of the dipole magnet and a new scintillating-fibre tracker downstream of the dipole magnet; a complete rebuild of the photon detection system of the Cherenkov detectors using multianode photomultipliers tubes; and redesigned and updated electronics of the calorimeters and the muon detector. A novel all-software trigger running on GPUs and on a dedicated computing farm has been deployed, and a completely renewed online system installed. To match the new trigger scheme, the software code base and computing model have been fully redesigned and reimplemented.

The performance of the new detector systems, as studied in the laboratory and with test-beam measurements, has been discussed, along with the expected overall experiment performance, estimated through Monte Carlo simulations, and is found to be as good as the previous experiment, if not better, while facing much higher luminosity and pileup running conditions.

The upgraded experiment will significantly extend the physics programme of LHCb, providing substantially larger statistics for precision studies and new physics searches, and opening new fields of investigation not only in the flavour physics domain but also, as a general purpose detector, in heavy ion, electroweak and fixed target physics.

## Acknowledgments

## Acronyms

**GPIO** general purpose input/output. 83, 107

**GPU** graphics processing unit. 128, 132, 133, 140, 146, 160, 162, 178, 179

**GWP** global warming potential. 7

**GWT** gigabit wireline transmitter. 17, 32, 129

**HCAL** hadronic calorimeter. 3, 101–106, 109, 111, 112, 124, 125, 152, 156

**HDR** high dynamic range. 132

**HEH** high-energy hadrons. 86, 87

**HLT** high level trigger. 32, 105, 127, 136, 137, 143, 148, 166, 169, 174, 176, 198

**HLT1** high level trigger first stage. 32, 128, 132, 133, 139, 140, 143–150, 153, 156, 157, 159, 160, 178, 180, 191–193

**HLT2** high level trigger second stage. 32, 76, 128, 133, 139, 140, 143, 144, 146, 148, 150, 152–154, 156–158, 160, 169, 170, 178, 180, 191, 194, 195

**HPC** high performance computing. 173, 176

**HPD** hybrid photon detector. 77

**HV** high voltage. 5, 46, 79, 81, 85, 94–96, 111, 112, 125, 126

**I2C** inter-integrated cicuit. 27, 45, 46, 70, 83, 84, 107, 116, 120, 121, 138

**IP** impact parameter. 179, 186

**JCOP** joint control project. 136–138

**JTAG** Joint Test Action Group industry standard. 83, 138

**L0** level-0 trigger. 1, 3, 9, 105, 109, 115, 116, 143, 156

**LC** Lucent connector. 86

**LEDTSB** LED trigger signal board. 111, 112

**LET** linear energy transfer. 86

**LHC** Large Hadron Collider. ii, 1, 2, 4–6, 9, 11–15, 30–35, 44, 50, 67, 72, 75, 105, 112, 113, 120, 125–127, 129, 133–136, 140, 145

**LLT** low level trigger. 105–107, 109, 110, 150

**LS1** LHC long shutdown 1. 8, 126

**LS2** LHC long shutdown 2. 6, 12, 125, 157

**LS3** LHC long shutdown 3. 104, 124

**LV** low voltage. 5, 50, 51, 84, 94, 95

**LVCMOS** low-voltage complementary metal-oxide semiconductor. 84

**MaPMT** multi-anode photomultiplier tube. 77–82, 84–86, 88, 90–92, 95–101

**MPO** multifibre push on connector. 86

**MSS** magnet safety system. 6

**MWPC** multi-wire proportional chambers. v, 115, 116, 124–127

**nCB** new custom backplane. 116, 121

**NEG** non-evaporable getter. 14, 31

**NIEL** non-ionising energy loss. 59, 60

**nODE** new off-detector electronics. v, 116–118, 125

**nPDM** new pulse distribution module. 116, 121

**nSB** new service board. 116, 121

**nSBS** new service board system. 116, 121, 122, 124

**nSYNC** new SYNC ASIC. 116–120, 123

**ODE** off-detector electronics. 116, 117, 119, 124

**OPB** opto- and power board. 25, 27, 28, 30, 33

**PA** pitch adapter. 43, 44, 53

**PCIe** peripheral component interconnect express. 28, 96, 123, 129, 130, 132, 135

**PCIe40** PCIe generic back-end board. v, 6, 32, 75, 129, 130, 132, 133

**PDM** photon detector module. 83, 84

**PDMDB** photon detector module digital board. 81–87, 95, 96, 98

**PID** particle identification. 3, 76–78, 98–100, 153, 167, 168, 190

**PLL** phase-locked loop. 45, 46, 50, 75, 120

**PLUME** probe for luminosity measurement. ii, 9–12

**PMT** photomultiplier tube. 9–11, 101–106, 108, 111, 112, 114, 156

**PV** primary vertex. 144, 150, 155, 180, 187

**QA** quality assurance. 79

**QE** quantum efficiency. 79

**RICH** ring imaging Cherenkov detector. 3, 7, 156, 157

**RMS** radiation monitoring system. ii, 9, 10, 12

**SALT** Silicon ASIC for LHCb Tracking. iii, 44–47, 49–51, 54

**SB** service board. 116

**SCA** slow control adapter. 5, 95, 121

**SCADA** supervisory control and data acquisition. 137, 138

**SciFi** scintillating fibre. 57, 58, 66, 70, 71

**SciFi Tracker** scintillating fibre tracker. 3, 6, 7, 38, 57–60, 63, 64, 67, 70–73, 75, 115, 147, 148, 150, 155, 156, 179, 186

**SEE** single event effect. 5, 86, 96

**SEL** single event latchup. 86, 87, 110

**SEU** single event upset. 17, 18, 45, 86, 87, 120, 121

**SEY** secondary electron yield. 14, 36

**SiPM** silicon photomultiplier. 6, 7, 57–61, 63–75

**SMOG** system for measuring the overlap with gas. 33–37

**SOL40** PCIe40 board for controls. 6, 25, 26, 32, 95, 116, 121, 130, 135, 136, 138

**SPI** serial peripheral interface. 51, 83, 107–109, 138

**SPP** super-pixel packet. 16, 17, 25, 32

**SPS** Super Proton Synchrotron. 75, 113

**STP** standard temperature and pressure. 77

**TCM** timing and control module. 83, 84, 95

**TELL40** PCIe40 board for data acquisition. 6, 11, 25, 26, 32, 55, 70, 87, 95, 96, 105, 116, 117, 121–123, 128, 130–132, 135, 141, 148

**TFC** timing and fast signal control. v, 5, 6, 11, 45, 46, 50, 67, 72, 84, 94, 95, 110, 116–122, 127, 128, 132–135

**UPS** uninterruptible power supply. 6

**UT** upstream tracker. 3, 7, 38–44, 46, 47, 49–57, 87, 147, 148, 150, 155, 156

**VELO** vertex locator. 3, 5, 7, 12–16, 22, 23, 25, 26, 28–38, 52, 87, 89, 129, 144, 145, 147, 148, 150, 151, 155–157, 179, 182, 183, 186

## References

[1] LHCb collaboration, *The LHCb Detector at the LHC*, 2008 *JINST* **3** S08005.

[2] LHCb collaboration, *LHCb Detector Performance*, *Int. J. Mod. Phys. A* **30** (2015) 1530022 [arXiv:1412.6352].

[3] LHCb collaboration, *Implications of LHCb measurements and future prospects*, *Eur. Phys. J. C* **73** (2013) 2373 [arXiv:1208.3355].

[4] LHCb collaboration, *LHCb technical proposal: A Large Hadron Collider Beauty Experiment for Precision Measurements of CP Violation and Rare Decays*, Tech. Rep. CERN-LHCC-98-04, CERN-LHCC-98-4, CERN-LHCC-P-4 (1998).

[5] LHCb collaboration, *Letter of Intent for the LHCb Upgrade*, Tech. Rep. CERN-LHCC-2011-001, CERN, Geneva (2011).

[6] LHCb collaboration, *Framework TDR for the LHCb Upgrade: Technical Design Report*, Tech. Rep. CERN-LHCC-2012-007 (2012).

[7] R. Aaij et al., *The LHCb Trigger and its Performance in 2011*, 2013 *JINST* **8** P04022 [arXiv:1211.3055].

[8] C. Fitzpatrick and V.V. Gligorov, *Anatomy of an upgrade event in the upgrade era, and implications for the LHCb trigger*, Tech. Rep. LHCb-PUB-2014-027, CERN-LHCb-PUB-2014-027, CERN, Geneva (2014).

[9] LHCb collaboration, *LHCb Trigger and Online Upgrade Technical Design Report*, Tech. Rep. CERN-LHCC-2014-016 (2014).

[10] LHCb collaboration, *LHCb magnet: Technical Design Report*, Tech. Rep. CERN-LHCC-2000-007, CERN, Geneva (2000).

[11] J. Andre et al., *Status of the LHCb magnet system*, *IEEE Trans. Appl. Supercond.* **12** (2002) 366.

[12] J. André et al., *Status of the LHCb dipole magnet*, *IEEE Trans. Appl. Supercond.* **14** (2004) 509.

[13] P. Moreira et al., *The GBT Project*, in the proceedings of the *Topical Workshop on Electronics for Particle Physics*, Paris, France, 21–25 September 2009, pp. 342–346 [DOI:10.5170/CERN-2009-006.342].

[14] K. Wyllie and F. Alessio, *Electronics Architecture of LHCb for Run3 and Future Upgrades*, LHCb-PUB-2022-014, CERN, Geneva (2022).

[15] CMS collaboration, *The versatile link, a common project for super-LHC*, 2009 *JINST* **4** P12003.

[16] A. Caratelli et al., *The GBT-SCA, a radiation tolerant ASIC for detector control and monitoring applications in HEP experiments*, 2015 *JINST* **10** C03034.

[17] F. Faccio et al., *Development of custom radiation-tolerant DCDC converter ASICs*, 2010 *JINST* **5** C11016.

[18] A. Ferrari, P.R. Sala, A. Fasso and J. Ranft, *FLUKA: A multi-particle transport code (Program version 2005)*, Tech. Rep. CERN-2005-010, SLAC-R-773, INFN-TC-05-11, CERN-2005-10 (2005) [DOI:10.2172/877507].

[19] T.T. Böhlen et al., *The FLUKA Code: Developments and Challenges for High Energy and Medical Applications*, *Nucl. Data Sheets* **120** (2014) 211.

[20] P. Petagna, B. Verlaat and A. Francescon, *Two-phase thermal management of silicon detectors for high energy physics*, in *Encyclopedia of two-phase heat transfer and flow III*, pp. 335–412, World Scientific, Singapore (2018).

[21] LHCb collaboration, *LHCb reoptimized detector design and performance: Technical Design Report*, Tech. Rep. CERN-LHCC-2003-030, CERN, Geneva (2003).

[22] L. Leduc, G. Corti and R. Veness, *Design of a highly optimised vacuum chamber support for the LHCb experiment*, Tech. Rep. LHCb-PROC-2011-048, CERN-ATS-2011-262, CERN, Geneva (2011).

[23] R. Alemany-Fernandez, F. Follin and R. Jacobsson, *The LHCb Online Luminosity Control and Monitoring*, in the proceedings of the *4th International Particle Accelerator Conference*, Shanghai, China, 12–17 May 2013, pp. 1346–1348.

[24] S. Barsuk et al., *Probe for LUminosity MEasurement in LHCb*, Tech. Rep. LHCb-PUB-2020-008, CERN, Geneva (2020).

[25] LHCb collaboration, *LHCb PLUME: Probe for LUminosity MEasurement*, Tech. Rep. CERN-LHCC-2021-002, CERN, Geneva (2021) [DOI:10.17181/CERN.WLU0.M37F].

[26] T. Sjöstrand, S. Mrenna and P.Z. Skands, *A Brief Introduction to PYTHIA 8.1*, *Comput. Phys. Commun.* **178** (2008) 852 [arXiv:0710.3820].

[27] GEANT4 collaboration, *GEANT4–a simulation toolkit*, *Nucl. Instrum. Meth. A* **506** (2003) 250.

[28] J. Allison et al., *Geant4 developments and applications*, *IEEE Trans. Nucl. Sci.* **53** (2006) 270.

[29] I. Guz, V. Belyaev, E. Chernov, V. Egorychev, S. Kandybei, T. Kvaratskheliya et al., *Upgrade of the monitoring system of LHCb ECAL*, Tech. Rep. LHCb-PUB-2016-018, CERN, Geneva (2016).

[30] G. Avoni et al., *The new LUCID-2 detector for luminosity measurement and monitoring in ATLAS*, 2018 *JINST* **13** P07017.

[31] C. Ilgner et al., *The Beam Conditions Monitor of the LHCb Experiment*, arXiv:1001.2487.

[32] O. Okhrimenko, V. Iakovenko, V. Pugatch, F. Alessio and G. Corti, *The Radiation Monitoring System for the LHCb Inner Tracker*, in the proceedings of the *13th International Conference on Accelerator and Large Experimental Physics Control Systems*, Grenoble, France, 10–14 October 2011, pp. 1115–1118, https://cds.cern.ch/record/1563821.

[33] LHCb collaboration, *LHCb VELO Upgrade Technical Design Report*, Tech. Rep. CERN-LHCC-2013-021 (2013).

[34] E.A. Papadelis, *Characterisation and Commissioning of the LHCb VELO Detector*, Ph.D. thesis, Vrije University, Amsterdam (2009).

[35] R.B. Appleby, M. Ferro-Luzzi, M. Giovannozzi, B. Holzer and M. Neat, *VELO aperture considerations for the LHCb upgrade*, Tech. Rep. LHCb-PUB-2012-018, CERN-LHCb-PUB-2012-018, CERN-ATS-Note-2012-101, CERN, Geneva (2012).

[36] C. Boscolo Meneguolo, R. Bruce, M. Ferro-Luzzi, M. Giovannozzi and S. Redaelli, *Calculation of the allowed aperture for a gas storage cell in IP8*, Tech. Rep. CERN-PBC-Notes-2018-008, CERN, Geneva (2018).

[37] V. Balagura, *Van der Meer scan luminosity measurement and beam-beam correction*, *Eur. Phys. J. C* **81** (2021) 26 [`arXiv:2012.07752`].

[38] B.K. Popovic and C. Vollinger, *Measurement and Simulation of the Longitudinal Impedance of the LHCb VELO*, Tech. Rep. CERN-ACC-NOTE-2019-0048 CERN, Geneva (2019).

[39] B.K. Popovic, *VELO with SMOG2 impedance-based heating localization analysis*, Tech. Rep. CERN-PBC-Notes-2019-003, CERN, Geneva (2019).

[40] R. Wanzenberg and O. Zagorodnova, *Calculation of Wakefields for the New Design of the LHCb Vertex Locator*, Tech. Rep. CERN-ACC-NOTE-2017-0034 CERN, Geneva (2016).

[41] P. Chiggiato and P.C. Pinto, *Ti–Zr–V non-evaporable getter films: From development to large scale production for the Large Hadron Collider*, *Thin Solid Films* **515** (2006) 382.

[42] T.M. Bird, *Flavour studies with LHCb: b-meson mixing, lepton-flavour violation and the VELO upgrade*, Ph.D. thesis, University of Manchester, Manchester, U.K. (2015).

[43] M. van Beuzekom et al., *VeloPix ASIC development for LHCb VELO upgrade*, *Nucl. Instrum. Meth. A* **731** (2013) 92.

[44] V. Gromov et al., *Development and applications of the Timepix3 readout chip*, *PoS* **VERTEX2011** (2011) 046.

[45] T. Poikela et al., *VeloPix: the pixel ASIC for the LHCb upgrade*, 2015 *JINST* **10** C01057.

[46] V. Gromov et al., *Development of a low power 5.12 Gbps data serializer and wireline transmitter circuit for the VeloPix chip*, 2015 *JINST* **10** C01054.

[47] B. van der Heijden et al., *SPIDR, a general-purpose readout system for pixel ASICs*, 2017 *JINST* **12** C02040.

[48] LHCb collaboration, *Microchannel Cooling for the LHCb VELO Upgrade I*, *Nucl. Instrum. Meth. A* **1039** (2022) 166874 [`arXiv:2112.12763`].

[49] R. Bates et al., *High speed electrical transmission line design and characterization*, 2017 *JINST* **12** C02002.

[50] L. Eklund et al., *The VELO optical and power board*, Tech. Rep. LHCb-PUB-2021-012, CERN-LHCb-PUB-2021-012, CERN, Geneva (2021) [`DOI:10.17181/CERN.TIOS.BTPB`].

[51] F. Faccio et al., *FEAST2: A Radiation and Magnetic Field Tolerant Point-of-Load Buck DC/DC Converter*, in the proceedings of the *Radiation Effects Data Workshop*, Paris, France, 14–18 July 2014 [`DOI:10.1109/REDW.2014.7004569`].

[52] P. Moreira, A. Marchioro and Kloukinas, *The GBT: A proposed architecure for multi-Gb/s data transmission in high energy physics*, in the proceedings of the *Topical Workshop on Electronics for Particle Physics*, Prague, Czech Republic, 3–7 September 2007, pp. 332–336 [`DOI:10.5170/CERN-2007-007.332`].

[53] P. Moreira et al., *The GBT-SerDes ASIC prototype*, 2010 *JINST* **5** C11022.

[54] G. Mazza et al., *A radiation tolerant 5 Gb/s Laser Driver in 130-nm CMOS technology*, 2012 *JINST* **7** C01052.

[55] G. Mazza et al., *High-speed, radiation-tolerant laser drivers in 0.13 μm CMOS technology for HEP applications*, *IEEE Trans. Nucl. Sci.* **61** (2014) 3653.

[56] K. Hennessy et al., *Readout Firmware of the Vertex Locator for LHCb Run 3 and Beyond*, *IEEE Trans. Nucl. Sci.* **68** (2021) 2472.

[57] G. Bassi et al., *A real-time FPGA-based cluster finding algorithm for LHCb silicon pixel detector*, *EPJ Web Conf.* **251** (2021) 04016.

[58] G. Bassi et al., *A FPGA-Based Architecture for Real-Time Cluster Finding in the LHCb Silicon Pixel Detector*, *IEEE Trans. Nucl. Sci.* **70** (2023) 1189 [`arXiv:2302.03972`].

[59] C. Barschel, *Precision luminosity measurement at LHCb with beam-gas imaging*, Ph.D. thesis, RWTH Aachen University, Aachen, Germany (2014).

[60] LHCb collaboration, *Precision luminosity measurements at LHCb*, 2014 *JINST* **9** P12005 [`arXiv:1410.0149`].

[61] LHCb collaboration, *LHCb SMOG Upgrade*, Tech. Rep. CERN-LHCC-2019-005, CERN, Geneva (2019) [`DOI:10.17181/CERN.SAQC.EOWH`].

[62] E. Steffens and W. Haeberli, *Polarized gas targets*, *Rept. Prog. Phys.* **66** (2003) R02.

[63] C. Yin Vallgren et al., *Amorphous carbon coatings for the mitigation of electron cloud in the CERN Super Proton Synchrotron*, *Phys. Rev. ST Accel. Beams* **14** (2011) 071001.

[64] R. Kersevan and M. Ady, *Recent developments of Monte-Carlo codes Molflow+ and Synrad+*, in the proceedings of the *10th International Particle Accelerator Conference*, Melbourne, Australia, 19–24 May 2019, pp. 1327–1330 [`DOI:10.18429/JACoW-IPAC2019-TUPMP037`].

[65] R. Mountain, *LHCb Upstream Tracker (UT)*, https://cds.cern.ch/record/2807067 (2021).

[66] LHCb collaboration, *Upstream Tracker - The silicon strip tracking detector for the LHCb Upgrade*, *PoS* **ICHEP2020** (2021) 724.

[67] M. Firlej et al., *A fast, low-power, 6-bit SAR ADC for readout of strip detectors in the LHCb Upgrade experiment*, 2014 *JINST* **9** P07006.

[68] K. Wyllie, F. Alessio, C. Gaspar, R. Jacobsson, R. Le Gac, N. Neufeld et al., *Electronics architecture of the LHCb upgrade*, Tech. Rep. LHCb-PUB-2011-011, CERN-LHCb-PUB-2011-011, CERN, Geneva (2013).

[69] F. Alessio and R. Jacobsson, *Readout control specifications for the front-end and back-end of the LHCb upgrade*, Tech. Rep. LHCb-PUB-2012-017, CERN-LHCb-PUB-2012-017. LHCb-INT-2012-018, CERN, Geneva (2014).

[70] C. Abellan Beteta et al., *The SALT—Readout ASIC for Silicon Strip Sensors of Upstream Tracker in the Upgraded LHCb Experiment*, *Sensors* **22** (2021) 107.

[71] G. Papotti, *An Error-Correcting Line Encoding ASIC for a HEP Rad-Hard Multi-GigaBit Optical Link*, in the proceedings of the *2006 Ph.D. Research in Microelectronics and Electronics*, Otranto, Italy, 12–15 June 2006, pp. 225–228 [`DOI:10.1109/rme.2006.1689937`].

[72] S. Coelli, *Development and test of the $CO_2$ evaporative cooling system for the LHCb UT Tracker Upgrade*, 2017 *JINST* **12** C03087.

[73] A. Abba et al., *Testbeam studies of pre-prototype silicon strip sensors for the LHCb UT upgrade project*, *Nucl. Instrum. Meth. A* **806** (2016) 244 [`arXiv:1506.00229`].

[74] A. Abba et al., *Study of prototype sensors for the Upstream Tracker upgrade*, Tech. Rep. LHCb-PUB-2016-007, CERN, Geneva (2016).

[75] *Beetle manual*, http://www.kip.uni-heidelberg.de/lhcb/Publications/BeetleRefMan_v1_3.pdf (accessed 2022-12-22).

[76] M. Artuso et al., *Signal coupling to embedded pitch adapters in silicon sensors*, *Nucl. Instrum. Meth. A* **877** (2018) 252 [`arXiv:1708.03371`].

[77] M. Artuso et al., *First beam test of UT sensors with the SALT 3.0 readout ASIC* , Tech. Rep. LHCb-PUB-2019-009, CERN, Geneva (2019).

[78] D. van Eijk et al., *Radiation hardness of the LHCb Outer Tracker*, *Nucl. Instrum. Meth. A* **685** (2012) 62.

[79] LHCb collaboration, *LHCb Tracker Upgrade Technical Design Report*, Tech. Rep. CERN-LHCC-2014-001, CERN, Geneva (2014).

[80] C. Joram, A.B. Rodriguez Cavalcante, L. Gavardi, F. Ravotti and T. Schneider, *Irradiation test of mirror samples for the LHCb SciFi tracker*, Tech. Rep. LHCb-PUB-2016-006, CERN-LHCb-PUB-2016-006, CERN, Geneva (2016).

[81] *Kuraray plastic scintillating fibres*, http://kuraraypsf.jp/psf/sf.html.

[82] O. Borshchev et al., *Development of a New Class of Scintillating Fibres with Very Short Decay Time and High Light Yield*, 2017 *JINST* **12** P05013.

[83] O. Shinji, *Oxygen diffusion in a polystyrene optical fibre surrounded by an oxygen gas barrier*, EP-Tech-Note-2019-001, CERN, Geneva (2019).

[84] O. Shinji, *Estimation of the optical aging rate for 1.00 mm plastic scintillating fibre of type SCSF-78M*, Tech. Rep. EP-Tech-Note-2019-002, CERN, Geneva (2019).

[85] A.B.R. Cavalcante et al., *Refining and testing 12,000 km of scintillating plastic fibre for the LHCb SciFi tracker*, 2018 *JINST* **13** P10025.

[86] C. Alfieri, A. Cavalcante, C. Joram and M. Kenzie, *An experimental set-up to measure light yield of scintillating fibres*, Tech. Rep. LHCb-PUB-2015-012, CERN-LHCb-PUB-2015-012, CERN, Geneva (2015).

[87] C. Joram, U. Uwer, B.D. Leverington, T. Kirn, S. Bachmann, R.J. Ekelhof et al., *LHCb scintillating fibre tracker Engineering Design Review report: Fibres, mats and modules*, Tech. Rep. LHCb-PUB-2015-008, CERN-LHCb-PUB-2015-008, CERN, Geneva (2015).

[88] Particle Data Group collaboration, *Review of Particle Physics*, *PTEP* **2020** (2020) 083C01.

[89] O. Girard et al., *Characterisation of silicon photomultipliers based on statistical analysis of pulse-shape and time distributions*, arXiv:1808.05775.

[90] C. Joram and U. Uwer, *LHCb scintillating fibre tracker Engineering Design Review: Front-end electronics*, Tech. Rep. LHCb-PUB-2016-012, CERN-LHCb-PUB-2016-012, CERN, Geneva (2016).

[91] L. Gruber, *LHCb SciFi — Upgrading LHCb with a scintillating fibre tracker*, *Nucl. Instrum. Meth. A* **958** (2020) 162025.

[92] B.D. Leverington, *LHCb Upgrade — The Scintillating Fibre Tracker*, *PoS* **EPS-HEP2015** (2015) 254.

[93] A. Comerma and J. Mazorra, *PACIFICr5 data sheet*, Tech. Rep. EDMS 1841761, Heidelberg Universität (2020).

[94] S. Vinogradov, *Analytical models of probability distribution and excess noise factor of Solid State Photomultiplier signals with crosstalk*, *Nucl. Instrum. Meth. A* **695** (2012) 247 [arXiv:1109.2014].

[95] L. Witola, *Calibration and performance studies of the readout ASIC for the LHCb SciFi tracker*, MSc thesis, Heidelberg Universität (2019).

[96] F. Vasey, *Versatile Link specification part 2.1: Front-end versatile transceiver and twin transmitter*, Tech. Rep. EDMS 1140665, CERN, Geneva (2013).

[97] P. Sainvitu, A. Zemanek, K. Nikolitsas and N. Witold, *LHCb-SCIFI assembly survey and photogrammetry measurement*, Tech. Rep. EDMS 2652340, CERN, Geneva (2021).

[98] K.S. Hashemi and J. Bensinger, *The BCAM Camera*, Tech. Rep. ATL-MUON-2000-024, CERN, Geneva (2000).

[99] L. Del Buono, H. Chanal, O. Le Dortz, A. Pellegrino and W. Vink, *LHCb upgrade SciFi tracker front-end data format*, Tech. Rep. EDMS 1898940, CERN, Geneva (2021).

[100] L. Del Buono, O. Le Dortz, A. Pellegrino and W. Vink, *LHCb upgrade SciFi tracker TELL40 data processing*, Tech. Rep. EDMS 1904563, CERN, Geneva (2020).

[101] M. Clemencic et al., *The LHCb simulation application, Gauss: Design, evolution and experience*, *J. Phys. Conf. Ser.* **331** (2011) 032023.

[102] S. Beranek, M.S. Bieker, M. Demmer, R.J. Ekelhof, C.O. Gerber, M.P. Whitehead et al., *Simulation of the light yield attenuation maps for the LHCb SciFi tracker upgrade*, Tech. Rep. LHCb-PUB-2019-007, CERN-LHCb-PUB-2019-007, CERN, Geneva (2019).

[103] J.P. Cachemiche et al., *The PCIe-based readout system for the LHCb experiment*, 2016 *JINST* **11** P02013.

[104] V. Blobel and C. Kleinwort, *A New method for the high precision alignment of track detectors*, in the proceedings of the *Conference on Advanced Statistical Techniques in Particle Physics*, Durham, U.K., 18–22 March 2002, pp. 268–277 [hep-ex/0208021].

[105] M. Adinolfi et al., *Performance of the LHCb RICH detector at the LHC*, *Eur. Phys. J. C* **73** (2013) 2431 [arXiv:1211.6759].

[106] R. Calabrese et al., *Performance of the LHCb RICH detectors during LHC Run 2*, 2022 *JINST* **17** P07013 [arXiv:2205.13400].

[107] C. D'Ambrosio, S. Easo, C. Frei and A. Petrolini, *RICH2019: a proposal for the LHCb RICH upgrade*, Tech. Rep. LHCb-PUB-2013-011, CERN-LHCb-PUB-2013-011, CERN, Geneva (2013).

[108] LHCb collaboration, *LHCb PID Upgrade Technical Design Report*, Tech. Rep. CERN-LHCC-2013-022 (2013).

[109] S. Okamura, *Commissioning of the upgraded RICH system at the LHCb experiment*, *JINST* **17** (2022) C11006.

[110] M. Alemi et al., *First operation of a hybrid photon detector prototype with electrostatic cross-focussing and integrated silicon pixel readout*, *Nucl. Instrum. Meth. A* **449** (2000) 48.

[111] C. Giugliano, *Quality assurance for the LHCb RICH upgrade Photon-Detection chain*, *Nucl. Instrum. Meth. A* **1055** (2023) 168436.

[112] M. Andreotti et al., *A Fast and Radiation-Hard Single-Photon Counting ASIC for the Upgrade of the LHCb RICH Detector at CERN*, in the proceedings of the *2017 IEEE Radiation Effects Data Workshop (REDW)*, New Orleans, LA, U.S.A., 17–21 July 2017 [DOI:10.1109/nsrec.2017.8115435].

[113] M. Baszczyk et al., *CLARO: an ASIC for high rate single photon counting with multi-anode photomultipliers*, 2017 *JINST* **12** P08019.

[114] C. Gotti, *An ASIC for fast single photon counting in the LHCb RICH upgrade*, 2017 *JINST* **12** C03016.

[115] M. Fiorini et al., *Radiation hardness tests and characterization of the CLARO-CMOS, a low power and fast single-photon counting ASIC in 0.35 micron CMOS technology*, *Nucl. Instrum. Meth. A* **766** (2014) 228.

[116] M. Andreotti et al., *Irradiation of the CLARO-CMOS chip, a fast ASIC for single-photon counting*, *Nucl. Instrum. Meth. A* **787** (2015) 234.

[117] J. Ramos-Martos et al., *Radiation characterization of the austriamicrosystems 0.35 µm CMOS technology*, in the proceedings of the *12th European Conference on Radiation and Its Effects on Components and Systems*, Seville, Spain, 19–23 September 2011, pp. 806–811 [`DOI:10.1109/radecs.2011.6131335`].

[118] J. Ramos et al., *SEE characterization of the AMS 0.35 µm CMOS technology*, in the proceedings of the *14th European Conference on Radiation and Its Effects on Components and Systems (RADECS)*, Oxford, U.K., 23–27 September 2013 [`DOI:10.1109/radecs.2013.6937402`].

[119] A. Thornton, Tech. Rep. *CHARM Facility Test Area Radiation Field Description*, CERN-ACC-NOTE-2016-0041 CERN, Geneva (2016).

[120] V.M. Placinta, *Complex integrated circuits in the radiation environment at the LHCb high energy physics experiment and extrapolation to the case of space-based experiments*, Ph.D. thesis, Polytechnic Institute, Bucharest, Romania (2020).

[121] H. Boterenbrood and B.I. Hallgren, *The Development of Embedded Local Monitor Board (ELMB)*, Tech. Rep. ATL-DAQ-2003-053, CERN, Geneva (2003) [`DOI:10.5170/CERN-2003-006.331`].

[122] LHCb collaboration, *Characterisation of signal-induced noise in Hamamatsu R11265 Multianode Photomultiplier Tubes*, 2021 *JINST* **16** P11030 [`arXiv:2110.00831`].

[123] T. Blake et al., *Quenching the scintillation in CF$_4$ Cherenkov gas radiator*, *Nucl. Instrum. Meth. A* **791** (2015) 27.

[124] S. Easo, *Overview of LHCb-RICH upgrade*, *Nucl. Instrum. Meth. A* **876** (2017) 160.

[125] LHCb collaboration, *LHCb calorimeters: Technical Design Report*, Tech. Rep. CERN-LHCC-2000-036, CERN, Geneva (2000).

[126] S. Barsuk et al., *Design and construction of electromagnetic calorimeter for LHCb experiment*, Tech. Rep. LHCb-2000-043 CERN, Geneva (2000).

[127] R. Djeliadine, O. Iouchtchenko and V.F. Obraztsov, *LHCb hadron trigger and HCAL cell size and length optimization*, Tech. Rep. LHCb-99-035, CERN, Geneva (1999).

[128] C. Abellán Beteta et al., *Calibration and performance of the LHCb calorimeters in Run 1 and 2 at the LHC*, Tech. Rep. LHCb-DP-2020-001, CERN, Geneva (2020) [`arXiv:2008.11556`].

[129] E. Picatoste et al., *Low noise front end ICECAL ASIC for the upgrade of the LHCb calorimeter*, 2012 *JINST* **7** C01080.

[130] C. Beigbeder-Beau et al., *The front-end electronics for LHCb calorimeters*, Tech. Rep. LHCb-2000-028, CERN, Geneva (2000).

[131] E. Picatoste et al., *Low noise 4-channel front end ASIC with on-chip DLL for the upgrade of the LHCb Calorimeter*, 2015 *JINST* **10** C04017.

[132] Y. Gilitsky et al., *LHCb calorimeters high voltage system*, *Nucl. Instrum. Meth. A* **571** (2007) 294.

[133] A. Konoplyannikov, *Electronics of LHCb calorimeter monitoring system*, in the proceedings of the *Topical Workshop on Electronics for Particle Physics*, Naxos, Greece, 15–19 September 2008, pp. 392–396, http://cds.cern.ch/record/1217584 [`DOI:10.5170/CERN-2008-008.392`].

[134] Y. Guz, R. Dzhelyadin, A. Konoplyannikov, V. Matveev, V. Novikov and M. Soldatov, *Design and integration of HV, LED monitoring and radioactive-source system for HCAL*, Tech. Rep. LHCb-2003-005, CERN, Geneva (2003).

[135] D. Breton and D. Charlet, *Using the SPECS in LHCb*, Tech. Rep. LHCb-2003-005, CERN, Geneva (2003).

[136] LHCb collaboration, *LHCb muon system: Technical Design Report*, Tech. Rep. CERN-LHCC-2001-010, CERN, Geneva (2001).

[137] LHCb collaboration, *LHCb: Addendum to the Muon System Technical Design Report*, Tech. Rep. CERN-LHCC-2003-002, CERN, Geneva (2003).

[138] LHCb collaboration, Tech. Rep. *LHCb: Second Addendum to the Muon System Technical Design Report*, CERN-LHCC-2005-012, CERN, Geneva (2005).

[139] A.A. Alves Jr. et al., *Performance of the LHCb muon system*, 2013 *JINST* **8** P02022 [arXiv:1211.1346].

[140] F. Archilli et al., *Performance of the Muon Identification at LHCb*, 2013 *JINST* **8** P10020 [arXiv:1306.0249].

[141] W. Bonivento, D. Marras and G. Auriemma, *Production of the front-end boards of the LHCb muon system*, Tech. Rep. CERN-LHCB-2007-150, CERN, Geneva (2008).

[142] W. Bonivento et al., *Development of the CARIOCA front-end chip for the LHCb muon detector*, *Nucl. Instrum. Meth. A* **491** (2002) 233.

[143] S. Cadeddu, A. Lai and C. Deplano, *The DIALOG chip in the front-end electronics of the LHCb muon detector*, *IEEE Trans. Nucl. Sci.* **52** (2005) 2726.

[144] P. Fresch, G. Chiodi, F. Iacoangeli and V. Bocci, *First Prototype of the Muon Frontend Control Electronics for the LHCb Upgrade: Hardware Realization and Test*, *Springer Proc. Phys.* **212** (2018) 173.

[145] A. Cardini, *The LHCb Muon Upgrade*, 2014 *JINST* **9** C02014.

[146] S. Bonacini, K. Kloukinas and P. Moreira, *E-link: A Radiation-Hard Low-Power Electrical Link for Chip-to-Chip Communication*, in the proceedings of the *Topical Workshop on Electronics for Particle Physics*, Paris, France, 21–25 September 2009, pp. 422–425 [DOI:10.5170/CERN-2009-006.422].

[147] S. Cadeddu et al., *The nSYNC ASIC for the new readout electronics of the LHCb Muon Detector Upgrade*, *Nucl. Instrum. Meth. A* **936** (2019) 378.

[148] S. Cadeddu, V. De Leo, C. Deplano and A. Lai, *The SYNC chip in the electronics architecture of the LHCb muon detector*, *IEEE Trans. Nucl. Sci.* **57** (2010) 2790.

[149] R. Giordano et al., *High-Resolution Synthesizable Digitally-Controlled Delay Lines*, *IEEE Trans. Nucl. Sci.* **62** (2015) 3163.

[150] D. Brundu, *Radiation hardness of the upgraded LHCb muon detector electronics and prospects for a full angular analysis in multi-body rare charm decays*, Ph.D. thesis, Università Di Cagliari, Cagliari, Italy (2020).

[151] D. Brundu, S. Cadeddu, A. Cardini and L. Casu, *Radiation hardness test of the UMC 130 nm nSYNC ASIC with a 60 MeV proton beam and X-Rays*, in the proceedings of the *2018 IEEE Nuclear Science Symposium and Medical Imaging Conference*, Sydney, NSW, Australia, 10–17 November 2018 [DOI:10.1109/NSSMIC.2018.8824458].

[152] D. Brundu, A. Cardini, S. Cadeddu, K. Wyllie and P. Ciambrone, *An X-Ray facility to perform irradiation tests and TID studies on electronics and detectors*, Tech. Rep. LHCb-PUB-2018-007. CERN-LHCb-PUB-2018-007, CERN, Geneva (2018).

[153] F. Alessio, J.-P. Cachemiche, P. Durante, N. Neufeld, R. Schwemmer, S. T'Jampens et al., *Readout protocol for TELL40*, Tech. Rep. EDMS 1606939, CERN, Geneva (2016).

[154] D. Pinci, *Performance of the muon MWPC in high luminosity runs*, Tech. Rep. LHCb-PUB-2013-005, CERN-LHCb-PUB-2013-005, CERN, Geneva (2013).

[155] N. Bondar, D. Ilin and O. Maev, *Proposal on application of the multi-wire proportional chambers of the LHCb MUON Detector at very high rates for the future upgrades*, 2020 *JINST* **15** C07001 [arXiv:2007.03058].

[156] F.P. Albicocco et al., *Long-term operation of the multi-wire-proportional-chambers of the LHCb muon system*, 2019 *JINST* **14** P11031 [arXiv:1908.02178].

[157] L. Malter, *Thin film field emission*, *Phys. Rev.* **50** (1936) 48.

[158] LHCb collaboration, *LHCb online system, data acquisition and experiment control: Technical Design Report*, Tech. Rep. CERN-LHCC-2001-040, CERN, Geneva (2001).

[159] L. Calefice et al., *Effect of the high-level trigger for detecting long-lived particles at LHCb*, *Front. Big Data* **5** (2022) 1008737.

[160] *X.200: Data networks and open system communications open systems interconnection — Model and notation*, http://www.itu.int/rec/T-REC-X.200-199407-I/en (Accessed 2022-12-23).

[161] F. Alessio and R. Jacobsson, *System-level specifications of the timing and fast control system for the LHCb upgrade*, Tech. Rep. LHCb-PUB-2012-001, CERN, Geneva (2012).

[162] F. Alessio et al., *A generic firmware core to drive the Front-End GBT-SCAs for the LHCb upgrade*, 2015 *JINST* **10** C02013.

[163] J. Barbosa, F. Alessio and C. Gaspar, *The new version of the LHCb SOL40-SCA core to drive front-end GBT-SCAs for the LHCb upgrade*, *PoS* **TWEPP-17** (2018) 078.

[164] F. Alessio et al., *Clock and timing distribution in the LHCb upgraded detector and readout system*, 2015 *JINST* **10** C02033.

[165] F. Alessio, P. Durante and G. Vouters, *The readout supervisor firmware for controlling the upgraded LHCb detector and readout system*, arXiv:1806.08626.

[166] F. Alessio, Z. Guzik and R. Jacobsson, *LHCb Global Timing and Monitoring of the LHC Filling Scheme*, Tech. Rep. LHCb-PUB-2011-004, CERN-LHCb-PUB-2011-004, CERN, Geneva (2011).

[167] D.R. Myers et al., *The LHC experiments joint controls project, JCOP*, in the proceedings of the 7th Biennial International Conference on Accelerator and Large Experimental Physics Control Systems, Trieste, Italy, 4–8 October 1999, 633.

[168] B. Franek and C. Gaspar, *SMI++ object oriented framework for designing and implementing distributed control systems*, *IEEE Trans. Nucl. Sci.* **45** (1998) 1946.

[169] C. Gaspar, *The LHCb experiment control system: On the path to full automation*, in the proceedings of the 13th International Conference on Accelerator and Large Experimental Physics Control Systems, Grenoble, France, 10–14 October 2011, pp. 20–23.

[170] M. Adinolfi et al., *LHCb data quality monitoring*, *J. Phys. Conf. Ser.* **898** (2017) 092027.

[171] C. Gaspar, M. Dönszelmann and P. Charpentier, *DIM, a portable, light weight package for information publishing, data transfer and inter-process communication*, *Comput. Phys. Commun.* **140** (2001) 102.

[172] *Foreman*, https://www.theforeman.org/.

[173] Puppet Inc., https://puppet.com/.

[174] Netapp Inc., https://www.netapp.com/.

[175] Kubernetes, https://kubernetes.io/.

[176] R. Aaij et al., *A comprehensive real-time analysis model at the LHCb experiment*, 2019 *JINST* **14** P04006 [arXiv:1903.01360].

[177] 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].

[178] V.V. Gligorov, *Conceptualization, implementation, and commissioning of real-time analysis in the High Level Trigger of the LHCb experiment*, Ph.D. thesis, Université Paris VI-VII, Paris, France (2018) [arXiv:1806.10912].

[179] C. Fitzpatrick, J.M. Williams, S. Meloni, T.J. Boettcher, M.P. Whitehead, A. Dziurda et al., *Upgrade trigger: Bandwidth strategy proposal*, Tech. Rep. LHCb-PUB-2017-006, CERN-LHCb-PUB-2017-006, CERN, Geneva (2017).

[180] LHCb collaboration, *RTA and DPA dataflow diagrams for Run 1, Run 2, and the upgraded LHCb detector*, Tech. Rep. LHCb-FIGURE-2020-016, CERN, Geneva (2020).

[181] LHCb collaboration, *Design and performance of the LHCb trigger and full real-time reconstruction in Run 2 of the LHC*, 2019 *JINST* **14** P04013 [arXiv:1812.10790].

[182] LHCb collaboration, *A Comparison of CPU and GPU Implementations for the LHCb Experiment Run 3 Trigger*, *Comput. Softw. Big Sci.* **6** (2022) 1 [arXiv:2105.04031].

[183] G. Barrand et al., *GAUDI — A software architecture and framework for building HEP data processing applications*, *Comput. Phys. Commun.* **140** (2001) 45.

[184] P. Mato, *GAUDI-Architecture design document*, Tech. Rep. LHCb-98-064 CERN, Geneva (1998).

[185] R. Aaij et al., *Allen: A high level trigger on GPUs for LHCb*, *Comput. Softw. Big Sci.* **4** (2020) 7 [arXiv:1912.09161].

[186] LHCb collaboration, *LHCb Upgrade GPU High Level Trigger Technical Design Report*, Tech. Rep. CERN-LHCC-2020-006, CERN, Geneva (2020) [DOI:10.17181/CERN.QDVA.5PIR].

[187] R.E. Kalman, *A New Approach to Linear Filtering and Prediction Problems*, *J. Basic Eng.* **82** (1960) 35.

[188] R. Fruhwirth, *Application of Kalman filtering to track and vertex fitting*, *Nucl. Instrum. Meth. A* **262** (1987) 444.

[189] LHCb collaboration, *Performance of the GPU HLT1 (Allen)*, Tech. Rep. LHCB-FIGURE-2020-014, CERN, Geneva (2020).

[190] V.V. Gligorov, *A single track HLT1 trigger*, Tech. Rep. LHCb-PUB-2011-003, CERN-LHCb-PUB-2011-003, LHCb-INT-2010-053, CERN, Geneva (2011).

[191] A. Hennequin et al., *A fast and efficient SIMD track reconstruction algorithm for the LHCb Upgrade 1 VELO-PIX detector*, 2020 *JINST* **15** P06018 [arXiv:1912.09901].

[192] P. Billoir, M. De Cian, P.A. Günther and S. Stemmle, *A parametrized Kalman filter for fast track fitting at LHCb*, *Comput. Phys. Commun.* **265** (2021) 108026 [arXiv:2101.12040].

[193] L. Anderlini et al., *Muon identification for LHCb Run 3*, 2020 *JINST* **15** T12005 [arXiv:2008.01579].

[194] LHCb collaboration, *Computing Model of the Upgrade LHCb experiment*, Tech. Rep. CERN-LHCC-2018-014, CERN, Geneva (2018) [DOI:10.17181/CERN.Q0P4.57ON].

[195] W. Hulsbergen, *The Global covariance matrix of tracks fitted with a Kalman filter and an application in detector alignment*, *Nucl. Instrum. Meth. A* **600** (2009) 471 [arXiv:0810.2241].

[196] J. Amoraal et al., *Application of vertex and mass constraints in track-based alignment*, *Nucl. Instrum. Meth. A* **712** (2013) 48 [arXiv:1207.4756].

[197] LHCb collaboration, *Novel real-time alignment and calibration of the LHCb detector and its performance*, *Nucl. Instrum. Meth. A* **845** (2017) 560.

[198] LHCb collaboration, *2018 Alignment stability plots*, Tech. Rep. LHCB-FIGURE-2019-015, CERN, Geneva (2019).

[199] R. Aaij et al., *Selection and processing of calibration samples to measure the particle identification performance of the LHCb experiment in Run 2*, *EPJ Tech. Instrum.* **6** (2019) 1 [arXiv:1803.00824].

[200] LHCb collaboration, *Measurement of the track reconstruction efficiency at LHCb*, 2015 *JINST* **10** P02007 [arXiv:1408.1251].

[201] LHCb collaboration, *Measurement of the electron reconstruction efficiency at LHCb*, 2019 *JINST* **14** P11023 [arXiv:1909.02957].

[202] LHCb collaboration, *LHCb computing: Technical Design Report*, Tech. Rep. CERN-LHCC-2005-019, CERN, Geneva (2005).

[203] LHCb et al. collaborations, *CORAL and COOL during the LHC Long Shutdown.*, *J. Phys. Conf. Ser.* **513** (2014) 042045.

[204] CERN VM File System (CVM-FS), https://cernvm.cern.ch/portal/filesystem.

[205] P. Buncic et al., *CernVM: A virtual software appliance for LHC applications*, *J. Phys. Conf. Ser.* **219** (2010) 042003.

[206] J. Blomer, C. Aguado Sanchez, P. Buncic and A. Harutyunyan, *Distributing LHC application software and conditions databases using the CernVM file system*, *J. Phys. Conf. Ser.* **331** (2011) 042003.

[207] M. Frank, F. Gaede, M. Petric and A. Sailer, *AIDASoft/DD4hep: v01-25-01*, https://zenodo.org/record/592244 [DOI:10.5281/ZENODO.592244].

[208] S. Ponce, P. Mato Vila, A. Valassi and I. Belyaev, *Detector description framework in LHCb*, *eConf* **C0303241** (2003) THJT007 [physics/0306089].

[209] LCG Release Area, https://ep-dep-sft.web.cern.ch/document/lcg-releases.

[210] LCGCmake, https://gitlab.cern.ch/sft/lcgcmake.

[211] git distributed version control system, https://git-scm.com/.

[212] Jira Software, https://www.atlassian.com/software/jira.

[213] A. Trisovic, B. Couturier, V. Gibson and C. Jones, *Recording the LHCb data and software dependencies*, *J. Phys. Conf. Ser.* **898** (2017) 102010.

[214] neo4j graph database, https://neo4j.com/.

[215] M. Clemencic and B. Couturier, *LHCb Build and Deployment Infrastructure for Run 2*, *J. Phys. Conf. Ser.* **664** (2015) 062008.

[216] Jenkins, https://jenkins.io/.

[217] OpenStack, https://www.openstack.org/.

[218] M. Clemencic, B. Couturier and S. Kyriazi, *Improvements to the User Interface for LHCb's Software continuous integration system.*, *J. Phys. Conf. Ser.* **664** (2015) 062025.

[219] A. Mazurov, B. Couturier, D. Popov and N. Farley, *Microservices for systematic profiling and monitoring of the refactoring process at the LHCb experiment*, *J. Phys. Conf. Ser.* **898** (2017) 072037.

[220] RPM package manager, http://rpm.org/.

[221] D. Müller, *Adopting new technologies in the LHCb Gauss simulation framework*, *EPJ Web Conf.* **214** (2019) 02004.

[222] B.G. Siddi and D. Müller, *Gaussino — a Gaudi-Based Core Simulation Framework*, in the proceedings of the *IEEE Nuclear Science Symposium and Medical Imaging Conference*, Manchester, U.K., 26 October–2 November 2019 [DOI:10.1109/NSS/MIC42101.2019.9060074].

[223] LHCb collaboration, *Handling of the generation of primary events in Gauss, the LHCb simulation framework*, *J. Phys. Conf. Ser.* **331** (2011) 032047.

[224] D.J. Lange, *The EvtGen particle decay simulation package*, *Nucl. Instrum. Meth. A* **462** (2001) 152.

[225] M. Mazurek, M. Clemencic and G. Corti, *Gauss and Gaussino: the LHCb simulation software and its new experiment agnostic core framework*, *PoS* **ICHEP2022** (2022) 225.

[226] M. Dobbs and J.B. Hansen, *The HepMC C++ Monte Carlo Event Record for High Energy Physics*, Tech. Rep. ATL-SOFT-2000-001, CERN, Geneva (2000).

[227] LHCb collaboration, *Performance of a multithreaded prototype for the future LHCb simulation framework (Gauss-on-Gaussino)*, Tech. Rep. LHCB-FIGURE-2019-012, CERN, Geneva (2019).

[228] C. Bozzi, *LHCb computing resource usage in 2021*, Tech. Rep. LHCb-PUB-2022-011, CERN-LHCb-PUB-2022-011, CERN, Geneva (2022).

[229] G. Corti et al., *How the Monte Carlo production of a wide variety of different samples is centrally handled in the LHCb experiment*, *J. Phys. Conf. Ser.* **664** (2015) 072014.

[230] M.P. Whitehead, *A palette of fast simulations in LHCb*, *PoS* **ICHEP2018** (2019) 271.

[231] LHCb collaboration, *Amplitude analysis of the $\Lambda_c^+ \to pK^-\pi^+$ and $\Lambda_c^+$ baryon polarization measurement in semileptonic beauty hadron decays*, *Phys. Rev. D* **108** (2023) 012023 [arXiv:2208.03262].

[232] LHCb collaboration, *Evidence for modification of b quark hadronization in high-multiplicity $pp$ collisions at $\sqrt{s} = 13\,TeV$*, arXiv:2204.13042.

[233] LHCb collaboration, *Search for the lepton-flavour violating decays $B^0 \to K^{*0}\mu^\pm e^\mp$ and $B_s^0 \to \phi\mu^\pm e^\mp$*, *JHEP* **06** (2023) 073 [arXiv:2207.04005].

[234] LHCb collaboration, *Search for the radiative $\Xi_b^- \to \Xi^-\gamma$ decay*, *JHEP* **01** (2022) 069 [arXiv:2108.07678].

[235] D. Müller, M. Clemencic, G. Corti and M. Gersabeck, *ReDecay: A novel approach to speed up the simulation at LHCb*, *Eur. Phys. J. C* **78** (2018) 1009 [arXiv:1810.10362].

[236] LHCb collaboration, *Observation of a $\Lambda_b^0 - \overline{\Lambda}_b^0$ production asymmetry in proton-proton collisions at $\sqrt{s} = 7$ and $8\,TeV$*, *JHEP* **10** (2021) 060 [arXiv:2107.09593].

[237] C. Bozzi, *LHCb Computing Resource usage in 2020*, Tech. Rep. LHCb-PUB-2021-003, CERN-LHCb-PUB-2021-003, CERN, Geneva (2021).

[238] M. Rama and G. Vitali, *Calorimeter fast simulation based on hit libraries LHCb Gauss framework*, *EPJ Web Conf.* **214** (2019) 02040.

[239] V. Chekalina et al., *Generative Models for Fast Calorimeter Simulation: the LHCb case*, *EPJ Web Conf.* **214** (2019) 02034 [arXiv:1812.01319].

[240] L. Anderlini, *Machine Learning for the LHCb Simulation*, arXiv:2110.07925.

[241] H. Voss, A. Hocker, J. Stelzer and F. Tegenfeldt, *TMVA, the Toolkit for Multivariate Data Analysis with ROOT*, *PoS* **ACAT** (2007) 040.

[242] A. Hocker et al., *TMVA — Toolkit for Multivariate Data Analysis with ROOT: Users guide*, Tech. Rep. physics/0703039, CERN, Geneva (2007).

[243] F. Pedregosa et al., *Scikit-learn: Machine Learning in Python*, *J. Machine Learning Res.* **12** (2011) 2825 [arXiv:1201.0490].

[244] F. Chollet et al., *Keras*, https://keras.io, 2015.

[245] L. Anderlini and M. Barbetti, *scikinC: a tool for deploying machine learning as binaries*, *PoS* **CompTools2021** (2022) 034.

[246] A. Maevskiy, D. Derkach, N. Kazeev, A. Ustyuzhanin, M. Artemev and L. Anderlini, *Fast Data-Driven Simulation of Cherenkov Detectors Using Generative Adversarial Networks*, *J. Phys. Conf. Ser.* **1525** (2020) 012097 [arXiv:1905.11825].

[247] M. Barbetti, *Techniques for parametric simulation with deep neural networks and implementation for the LHCb experiment at CERN and its future upgrades*, MSc thesis, University of Florence, Florence, Italy (2020).

[248] DaVinci project, http://lhcbdoc.web.cern.ch/lhcbdoc/davinci/.

[249] N. Skidmore, E. Rodrigues and P. Koppenburg, *Run-3 offline data processing and analysis at LHCb*, *PoS* **EPS-HEP2021** (2022) 792.

[250] M. Ferrillo, *New generation offline software for the LHCb upgrade I*, Tech. Rep. LHCb-PROC-2022-004, CERN-LHCb-PROC-2022-004, CERN, Geneva (2022).

[251] *CERN Open Data Policy for the LHC Experiments*, Tech. Rep. CERN-OPEN-2020-013, CERN, Geneva (2020) [DOI:10.17181/CERN.QXNK.8L2G].

[252] E. Rodrigues et al., *The Scikit HEP Project — overview and prospects*, *EPJ Web Conf.* **245** (2020) 06028 [arXiv:2007.03577].

[253] DIRAC collaboration, *DIRAC in Large Particle Physics Experiments*, *J. Phys. Conf. Ser.* **898** (2017) 092020.

[254] R.G. Aparicio and I.C. Coz, *Database on Demand: insight how to build your own DBaaS*, *J. Phys. Conf. Ser.* **664** (2015) 042021.

[255] Z. Mathe, A.C. Ramo, N. Lazovsky and F. Stagni, *The DIRAC Web Portal 2.0*, *J. Phys. Conf. Ser.* **664** (2015) 062039.

[256] A. Casajus Ramo, R. Graciani Diaz and A. Tsaregorodtsev, *DIRAC RESTful API*, *J. Phys. Conf. Ser.* **396** (2012) 052019.

[257] V. Méndez Muñoz et al., *The integration of CloudStack and OCCI/OpenNebula with DIRAC*, *J. Phys. Conf. Ser.* **396** (2012) 032075.

[258] DIRAC collaboration, *DIRAC universal pilots*, *J. Phys. Conf. Ser.* **898** (2017) 092024.

[259] F. Stagni et al., *LHCbDirac: Distributed computing in LHCb*, *J. Phys. Conf. Ser.* **396** (2012) 032104.

[260] Z. Mathe and P. Charpentier, *Optimising query execution time in LHCb Bookkeeping System using partition pruning and Partition-Wise joins*, *J. Phys. Conf. Ser.* **513** (2014) 042032.

[261] LHCb collaboration, *The LHCb DIRAC-based production and data management operations systems*, *J. Phys. Conf. Ser.* **368** (2012) 012010.

[262] F. Stagni, A. Tsaregorodtsev, A. McNab and C. Luzzi, *Pilots 2.0: DIRAC pilots for all the skies*, *J. Phys. Conf. Ser.* **664** (2015) 062061.

[263] A. Filipčič, D. Cameron and J.K. Nilsen, *Dynamic Resource Allocation with the arcControlTower*, *J. Phys. Conf. Ser.* **664** (2015) 062015.

[264] B. Bockelman et al., *Commissioning the HTCondor-CE for the Open Science Grid*, *J. Phys. Conf. Ser.* **664** (2015) 062003.

[265] A. McNab, *The vacuum platform*, *J. Phys. Conf. Ser.* **898** (2017) 052028.

[266] R. Currie et al., *Recent developments in user-job management with Ganga*, *J. Phys. Conf. Ser.* **664** (2015) 062010.

[267] R. Currie et al., *Expanding the user base beyond HEP for the Ganga distributed analysis user interface*, *J. Phys. Conf. Ser.* **898** (2017) 052032.

[268] Boole, http://lhcbdoc.web.cern.ch/lhcbdoc/boole.

[269] L. Arrabito et al., *Major changes to the LHCb Grid computing model in year 2 of LHC data*, *J. Phys. Conf. Ser.* **396** (2012) 032092.

[270] C. Haen, P. Charpentier, M. Frank and A. Tsaregorodtsev, *The DIRAC Data Management System and the Gaudi dataset federation*, *J. Phys. Conf. Ser.* **664** (2015) 042025.

[271] LHCb collaboration, *Selected HLT2 reconstruction performance for the LHCb upgrade*, LHCB-FIGURE-2021-003, CERN, Geneva, Switzerland (2021).

## The LHCb collaboration

R. Aaij [32], A.S.W. Abdelmotteleb [50], C. Abellan Beteta [44], F. Abudinén [50], C. Achard [9], T. Ackernley [54],
B. Adeva [40], M. Adinolfi [48], P. Adlarson [77], H. Afsharnia [9], C. Agapopoulou [13], C.A. Aidala [78],
Z. Ajaltouni [9], S. Akar [59], K. Akiba [32], P. Albicocco [23], J. Albrecht [15], F. Alessio [42], M. Alexander [53],
A. Alfonso Albero [39], Z. Aliouche [56], P. Alvarez Cartelle [49], R. Amalric [13], S. Amato [2],
J.L. Amey [48], Y. Amhis [11,42], L. An [42], L. Anderlini [22], M. Andersson [44], A. Andreani [25,m],
A. Andreianov [38], M. Andreotti [21], D. Andreou [62], J.E. Andrews [60], M. Anelli [23], A. Anjam [17], D. Ao [6],
F. Archilli [31,u], K. Arnaud [10], A. Artamonov [38], M. Artuso [62], J. Ashby [53], E. Aslanides [10],
M. Atzeni [44], B. Audurier [12], D. Ayres Rocha [1], I.B Bachiller Perea [8], S. Bachmann [17],
M. Bachmayer [43], J.J. Back [50], A. Bailly-reyre [13], P. Baladron Rodriguez [40], V. Balagura [12], G. Balbi [20],
W. Baldini [21,42], A. Balla [23], M. Baltazar [11], H. Band [32], J. Baptista de Souza Leite [1], M. Barbetti [22,k],
P. Barclay [51], R.J. Barlow [56], S. Barsuk [11], W. Barter [52], M. Bartolini [49], F. Baryshnikov [38],
J.M. Basels [14], G. Bassi [29,r], M. Baszczyk [35,w], J.C. Batista Lopes [42], B. Batsukh [4], A. Battig [15],
A. Bay [43], A. Beck [50], M. Becker [15], F. Bedeschi [29], I.B. Bediaga [1], C. Beigbeder-Beau [11],
A. Beiter [62], S. Belin [40], V. Bellee [44], K. Belous [38], I. Belov [38], I. Belyaev [38], G. Benane [10],
G. Bencivenni [23], M. Benettoni [28], E. Ben-Haim [13], A. Berezhnoy [38], F. Bernard [43], R. Bernet [44],
S. Bernet Andres [76], D. Berninghoff [17], H.C. Bernstein [62], C. Bertella [56], A. Bertolin [28], C. Betancourt [44],
F. Betti [42], Ia. Bezshyiko [44], O. Bezshyyko [80], S. Bhasin [48], J. Bhom [35], L. Bian [68], M.S. Bieker [15],
N.V. Biesuz [21], P. Billoir [13], A. Biolchini [32], M. Birch [55], F.C.R. Bishop [49], A. Bitadze [56],
A. Bizzeti [], M.P. Blago [49], T. Blake [50], F. Blanc [43], J.E. Blank [15], S. Blusk [62], D. Bobulska [53],
B. Bochin [38], J.A. Boelhauve [15], O. Boente Garcia [12], T. Boettcher [59], G. Bogdanova [38],
I. Boiaryntseva [79], A. Boldyrev [38], C.S. Bolognani [74], R. Bolzonella [21,j], N. Bondar [38,42],
M.J. Booth [57,51], F. Borgato [28], S. Borghi [56], M. Borsato [17], J.T. Borsuk [35], H. Boterenbrood [32],
S.A. Bouchiba [43], T.J.V. Bowcock [54], A. Boyaryntsev [79], A. Boyer [42], C. Bozzi [21], M.J. Bradley [55],
S. Braun [60], A. Brea Rodriguez [40], G. Bregliozzi [42], K. Bridges [54], M.M.J. Briere [11], M. Brock [57],
M. Brodski [42], J. Brodzicka [35], A. Brossa Gonzalo [40], C. Brown [62], J. Brown [54], A.J. Brummitt [51],
D. Brundu [27], L. Brunetti [8], L. Buda [62], A. Buonaura [44], L. Buonincontri [28], A.T. Burke [56],
L. Burmistrov [11], C. Burr [42], A. Bursche [66], A. Butkevich [38], J.S. Butter [32], J. Buytaert [42],
W. Byczynski [42], J.P. Cachemiche [10], S. Cadeddu [27], H. Cai [68], A. Caillet [42], R. Calabrese [21,j],
L. Calefice [15], D. Calegari [42], S. Cali [23], M. Calvi [26,n], M. Calvo Gomez [76], P. Campana [23],
D.H. Campora Perez [74], A.F. Campoverde Quezada [6], S. Canfer [51], S. Capelli [26,n], L. Capriotti [20],
V. Carassiti [21], A. Carbone [20,h], A. Carbone [25,m], R. Cardinale [24,l], A. Cardini [27], M. Carletti [23],
P. Carniti [26,n], J. Carroll [54], L. Carus [14], A. Casais Vidal [40], R. Caspary [17], G. Casse [54], M. Cattaneo [42],
G. Cavallero [55,42], V. Cavallini [21,j], L. Ceelie [32], S. Celani [43], J. Cerasoli [10], D. Cervenkov [57],
S. Cesare [25,m], B. Chadaj [42], A.J. Chadwick [54], I. Chahrour [78], H. Chanal [9], M.G. Chapman [48],
M. Charles [13], Ph. Charpentier [42], V.J. Chaumat [11], C.A. Chavez Barajas [54], M. Chefdeville [8],
C. Chen [10], S. Chen [4], A. Chernov [35], E. Chernov [38], S. Chernyshenko [46], S. Chiozzi [21],
V. Chobanova [40], S. Cholak [43], M. Chrzaszcz [35], A. Chubykin [38], V. Chulikov [38], P. Ciambrone [23],
M.F. Cicala [50], X. Cid Vidal [40], G. Ciezarek [42], P. Cifra [42], M. Citterio [25], G. Ciullo [j,21], K. Clark [48],
P.E.L. Clarke [52], M. Clemencic [42], H.V. Cliff [49], J. Closier [42], J.L. Cobbledick [56], V. Coco [42],
S. Coelli [25], J. Cogan [10], E. Cogneras [9], L. Cojocariu [37], P. Collins [42], T. Colombo [42],
L. Congedo [19], N. Conti [25], A. Contu [27], N. Cooke [47], I. Corredoira [40], G. Corti [42],

A. Cotta Ramusino [21], B. Couturier [42], G.A. Cowan [52], D.C. Craik [44], M. Cruz Torres [1,f],
R. Currie [52], C.L. Da Silva [61], S. Dadabaev [38], L. Dai [65], X. Dai [5], E. Dall'Occo [15], J. Dalseno [40],
C. D'Ambrosio [42], A. Damen [32], J. Daniel [9], A. Danilina [38], P. d'Argent [19], F. Daudon [9], J.E. Davies [56],
A. Davis [56], J. Davis [60], O. De Aguiar Francisco [56], F. De Benedetti [25,42], J. de Boer [42], K. De Bruyn [73],
S. De Capua [56], M. De Cian [43], U. De Freitas Carneiro Da Graca [1], E. De Lucia [23], J.M. De Miranda [1],
R. de Oliveira [42], L. De Paula [2], K. De Roo [32], M. De Serio [19,g], D. De Simone [44], P. De Simone [23],
F. De Vellis [15], J.A. de Vries [74], E. De Wit [32], C.T. Dean [61], F. Debernardis [19,g], D. Decamp [8],
M. Deckenhoff [15], V. Dedu [10], L. Del Buono [13], B. Delaney [58], H.-P. Dembinski [15], C. Denis [42],
V. Denysenko [44], O. Deschamps [9], F. Dettori [27,i], B. Dey [71], D. Di Bari [23], P. Di Nezza [23],
I. Diachkov [38], S. Didenko [38], L. Dieste Maronas [40], H. Dijkstra [42], S. Ding [62], V. Dobishuk [46],
M. Doets [32], F. Doherty [53], A. Dolmatov [38], M. Domke [15], C. Dong [3], A.M. Donohoe [18], F. Dordei [27],
P. Dorosz [35,w], A.C. dos Reis [1], L. Douglas [53], A.G. Downes [8], O. Duarte [11], P. Duda [75],
M.W. Dudek [35], L. Dufour [42], V. Duk [72], R. Dumps [42], P. Durante [42], M.M. Duras [75],
J.M. Durham [61], D. Dutta [56], P.Y. Duval [10], M. Dziewiecki [17], A. Dziurda [35], A. Dzyuba [38],
S. Easo [51], U. Egede [63], V. Egorychev [38], C. Eirea Orro [40], S. Eisenhardt [52], E. Ejopu [56], R. Ekelhof [15],
S. Ek-In [43], L. Eklund [77], M.E Elashri [59], J. Ellbracht [15], A. Elvin [56], S. Ely [55], A. Ene [37],
E. Epple [59], S. Escher [14], J. Eschle [44], S. Esen [44], T. Evans [56], F. Fabiano [27,i], L.N. Falcao [1],
Y. Fan [6], B. Fang [11,68], L. Fantini [72,q], M. Faria [43], S. Farry [54], D. Fazzini [26,n], L.F Felkowski [75],
M. Feo [42], P. Fernandez Declara [42], M. Fernandez Gomez [40], A. Fernandez Prieto [40], A.D. Fernez [60],
F. Ferrari [20], R. Ferreira [42], L. Ferreira Lopes [43], F. Ferreira Rodrigues [2], S. Ferreres Sole [32],
M. Ferrillo [44], M. Ferro-Luzzi [42], S. Filippov [38], R.A. Fini [19], M. Fiorini [21,j], M. Firlej [34],
K.M. Fischer [57], D.S. Fitzgerald [78], C. Fitzpatrick [56], T. Fiutowski [34], F. Fleuret [12], L. Flores [53],
M. Fontana [13], F. Fontanelli [24,l], R. Forty [42], D. Foulds-Holt [49], C. Fournier [42], V. Franco Lima [54],
M. Franco Sevilla [60], M. Frank [42], E. Franzoso [21,j], G. Frau [17], J. Freestone [56], C. Frei [42], R. Frei [43],
J. Frelier [62], D.A. Friday [53], L.F Frontini [25], J. Fu [6], Q. Fuehring [15], T. Fulghesu [13], C. Fuzipeg [56],
E. Gabriel [32], G. Galati [19,g], M.D. Galati [32], M. Galka [42], A. Gallas Torreira [40], D. Galli [20,h],
S. Gallorini [28,42], S. Gambetta [52,42], Y. Gan [3], M. Gandelman [2], P. Gandini [25], R. Gao [57], Y. Gao [7],
Y. Gao [5], M. Garau [27,i], L.M. Garcia Martin [50], P. Garcia Moreno [39], J. García Pardiñas [26,n],
B. Garcia Plana [40], F.A. Garcia Rosales [12], L. Garrido [39], N. Garroum [13], P.J. Garsed [49], D. Gascon [39],
C. Gaspar [42], C. Gasq [9], M. Gatta [23], L. Gavardi [15], P.M. Gebolis [42], R.E. Geertsema [32], D. Gerick [17],
L.L. Gerken [15], D. Germann [62], E. Gersabeck [56], M. Gersabeck [56], T. Gershon [50], S.A. Getz [38],
L. Giambastiani [28], V. Gibson [49], H.K. Giemza [36], A.L. Gilman [57], M. Giovannetti [23,u],
A. Gioventù [40], O.G. Girard [43], P. Gironella Gironell [39], C. Giugliano [21,j], M.A. Giza [35],
K. Gizdov [52], E.L. Gkougkousis [42], V.V. Gligorov [13,42], C. Göbel [64], L. Golinka-Bezshyyko [80],
E. Golobardes [76], D. Golubkov [38], A. Golutvin [55,38], A. Gomes [1,a], S. Gomez Fernandez [39],
F. Goncalves Abrantes [57], M. Goncerz [35], G. Gong [3], I.V. Gorelov [38], C. Gotti [26], J.P. Grabowski [70],
T. Grammatico [13], L.A. Granado Cardoso [42], F. Grant [53], E. Graugés [39], E. Graverini [43], G. Graziani [],
A.T. Grecu [37], L.M. Greeven [32], R. Greim [32], N.A. Grieser [59], L. Grillo [53], S. Gromov [38], V. Gromov [32],
N. Grub [42], B.R. Gruberg Cazon [57], B. Grynyov [79], C. Gu [3], M. Guarise [21,j], S. Guerin [62],
M. Guittiere [11], P.A. Günther [17], E. Gushchin [38], A. Guth [14], Y. Guz [38], T. Gys [42], F. Hachon [10],
T. Hadavizadeh [63], C. Hadjivasiliou [60], G. Haefeli [43], C. Haen [42], J. Haimberger [42], S.C. Haines [49],
T. Halewood-leagas [54], M.M. Halvorsen [42], P.M. Hamilton [60], J. Hammerich [54], S. Hamrat [9], Q. Han [7],
X. Han [17], E.B. Hansen [56], S. Hansmann-Menzemer [17], L. Hao [6], N. Harnew [57], T. Harrison [54],

C. Hasse [42], M. Hatch [42], J. He [6,c], K. Heijhoff [32], F.H Hemmer [42], C. Henderson [59],
R.D.L. Henderson [63,50], A.M. Hennequin [58], K. Hennessy [54], L. Henry [42], J. Herd [55], T. Herold [17],
J. Heuel [14], A. Hicheur [2], D. Hill [43], M. Hilton [56], G.T. Hoft [32], S.E. Hollitt [15], P.H. Hopchev [43],
O. Hornberger [17], J. Horswill [56], R. Hou [7], Y. Hou [8], J. Hu [17], J. Hu [66], W. Hu [5], X. Hu [3],
W. Huang [6], X. Huang [68], W. Hulsbergen [32], S. Hummel [17], R.J. Hunter [50], M. Hushchyn [38],
O.E. Hutanu [37], D. Hutchcroft [54], D. Hynds [32], P. Ibis [15], M. Idzik [34], D. Ilin [38], P. Ilten [59],
A. Inglessi [38], A. Iniukhin [38], C. Insa [9], A. Ishteev [38], K. Ivshin [38], R. Jacobsson [42], H. Jage [14],
S.J. Jaimes Elles [41], S. Jakobsen [42], O. Jamet [42], E. Jans [32], B.K. Jashal [41], M. Jaspers [32],
A. Jawahery [60], M. Jevaud [10,†], V. Jevtic [15], E. Jiang [60], X. Jiang [4,6], Y. Jiang [6], D. John [32],
M. John [57], D. Johnson [58], C.R. Jones [49], T.P. Jones [50], B. Jost [42], N. Jurik [42], I. Juszczak [35],
S. Kandybei [45], Y. Kang [3], M. Karacson [42], J.M. Kariuki [57], D. Karpenkov [38], W. Karpinski [14],
M. Karpov [38], K. Kaufmann [42], J.W. Kautz [59], F. Kayzel [32], F. Keizer [42], D.M. Keller [62], M. Kenzie [50],
T. Ketel [32], B. Khanji [15], A. Kharisova [38], S. Kholodenko [38], G. Khreich [11], T. Kirn [14],
V.S. Kirsebom [43], O. Kitouni [58], S. Klaver [33], N. Kleijne [29,r], K. Klimaszewski [36], M.R. Kmiec [36],
H. Kok [32], S. Koliiev [46], L. Kolk [15], A. Kondybayeva [38], A. Konoplyannikov [38], P. Kopciewicz [34],
R. Kopecna [17], P. Koppenburg [32], M. Korolev [38], J. Kos [32], I. Kostiuk [32], O. Kot [46], S. Kotriakhova [],
A. Kozachuk [38], V.S. Kozlov [38], M. Kraan [32], P. Kravchenko [38], L. Kravchuk [38], R.D. Krawczyk [42],
M. Kreps [50], S. Kretzschmar [14], P. Krokovny [38], W. Krupa [34], W. Krzemien [36], J. Kubat [17],
S. Kubis [75], W. Kucewicz [35,w], M. Kucharczyk [35], V. Kudryavtsev [38], A. Kuhlman [62], W.C. Kuilman [32],
E.K Kulikova [38], A.K. Kuonen [43], N. Kupfer [17], A. Kupsc [77], T. Kvaratskheliya [38], D. Lacarrere [42],
G. Lafferty [56], A. Lai [27], A. Lampis [27,i], D. Lancierini [44], C. Landesa Gomez [40], J.J. Lane [56],
R. Lane [48], C. Langenbruch [14], J. Langer [15], M. Langstaff [56], O. Lantwin [38], T. Latham [50],
F. Lazzari [29,s], M. Lazzaroni [25,m], O. Le Dortz [13], R. Le Gac [10], S.H. Lee [78], R. Lefèvre [9],
A. Leflat [38], S. Legotin [38], F. Lemaitre [42], E. Lemos Cid [40], P. Lenisa [j,21], O. Leroy [10], T. Lesiak [35],
B. Leverington [17], A. Li [3], H. Li [66], K. Li [7], P. Li [42], P.-R. Li [67], S. Li [7], T. Li [4], T. Li [66],
Y. Li [4], Z. Li [62], X. Liang [62], B. Lieunard [8], C. Lin [6], T. Lin [51], R. Lindner [42], V. Lisovskyi [15],
R. Litvinov [27,i], G. Liu [66], H. Liu [6], Q. Liu [6], S. Liu [4,6], A. Lobo Salvia [39], A. Loi [27],
R. Lollini [72], J. Lomba Castro [40], I. Longstaff [53], J.H. Lopes [2], A. Lopez Huertas [39], S. López Soliño [40],
D. Louis [14], G.H. Lovell [49], P. Loveridge [51], A.D. Lowe [57], Y. Lu [4,b], C. Lucarelli [22,k],
D. Lucchesi [28,p], S. Luchuk [38], M. Lucio Martinez [74], V. Lukashenko [32,46], A. Lukianov [38], H. Luo [52],
Y. Luo [3], A. Lupato [56], E. Luppi [21,j], O. Lupton [50], A. Lusiani [29,r], L.F. Lutz [60], K. Lynch [18],
X.-R. Lyu [6], R. Ma [6], S. Maccolini [15], F. Machefert [11], F. Maciuc [37], I. Mackay [57], V. Macko [43],
P. Mackowiak [15], S. Maddrell-Mander [48], L.R. Madhan Mohan [48], A. Maevskiy [38], M. Magne [9],
D. Maisuzenko [38], M.W. Majewski [34], R. Malaguti [21], J.J. Malczewski [35], S. Malde [57], B. Malecki [35,42],
A. Malinin [38], K. Malkinski [11], T. Maltsev [38], G. Manca [27,i], G. Mancinelli [10], C. Mancuso [11,25,m],
R. Manera Escalero [39], D. Manuzzi [20], C.A. Manzari [44], D. Marangotto [25,m], J.F. Marchand [8],
U. Marconi [20], S. Mariani [22,k], C. Marin Benito [39], J. Marks [17], A.M. Marshall [48], P.J. Marshall [54],
G. Martelli [72,q], G. Martellotti [30], L. Martinazzoli [42,n], M. Martinelli [26,n], D. Martinez Santos [40],
F. Martinez Vidal [41], B. Masic [52], A. Massafferri [1], M. Materok [14], R. Matev [42], A. Mathad [44],
Z. Mathe [42], V. Matiunin [38], C. Matteuzzi [26], K.R. Mattioli [12], A. Mauri [32], E. Maurice [12],
J. Mauricio [39], J. Mazorra de Cos [41], M. Mazurek [42], M. McCann [55], L. Mcconnell [18],
T.H. McGrath [56], N.T. McHugh [53], A. McNab [56], R. McNulty [18], J.V. Mead [54], B. Meadows [59],
G. Meier [15], L. Meier-villardita [44], D. Melnychuk [36], S. Meloni [26,n], M. Merk [32,74], A. Merli [25,m],

J.L. Meunier [13], L. Meyer Garcia [2], D. Miao [4,6], M. Mikhasenko [70,e], D.A. Milanes [69], E. Millard[50],
G. Miller [56], M. Milovanovic[42], M.-N. Minard[8,†], A. Minotti [26,n], S. Minutoli [24], T. Miralles [9],
S.E. Mitchell [52], B. Mitreska [15], T. Mittelstaedt[17], D.S. Mitzel [15], A. Mödden [15], L. Modenese[28],
A. Mogini[13], R.A. Mohammed [57], R.D. Moise [14], S. Mokhnenko [38], T. Mombächer [40],
M. Monk [50,63], I.A. Monroy [69], S. Monteil [9], M. Monti [25], M. Morandin [28], G. Morello [23],
M.J. Morello [29,r], M.P. Morgenthaler [17], J. Moron [34], A.B. Morris [42], A.G. Morris [50],
R. Mountain [62], H. Mu [3], E. Muhammad [50], F. Muheim [52], M. Mulder [73], S. Muley [17], D. Müller[42,56],
K. Müller [44], B. Munneke[32], C.H. Murphy [57], D. Murray [56], R. Murta [55], P. Muzzetto [27,i], P. Naik [48],
S.A. Naik [53], T. Nakada [43], R. Nandakumar [51], T. Nanut [42], I. Nasteva [2], E. Nazarov[38],
M. Needham [52], I. Neri [21,j], N. Neri [25,m], S. Neubert [70], N. Neufeld [42], P. Neustroev[38],
R. Newcombe[55], T. Nguyen Trung[11], J. Nicolini [15,11], D. Nicotra [74], E.M. Niel [43], S. Nieswand[14],
N. Nikitin[38], N.S. Nolte [58], C. Normand [8,i,27], J. Novoa Fernandez [40], G.N Nowak [59], C. Nunez [78],
T. O'Bannon [60], A. Oblakowska-Mucha [34], V. Obraztsov[38], J. O'Dell[51], T. Oeser [14], S. Okamura [21,j],
R. Oldeman [27,i], F. Oliva [52], P. Olive[10], C.J.G. Onderwater [73], R.H. O'Neil [52], V. Orlov[80],
J.M. Otalora Goicochea [2], T. Ovsiannikova[38], P. Owen [44], A. Oyanguren [41], O. Ozcelik [52],
K.O. Padeken [70], B. Pagare [50], P.R. Pais [42], T. Pajero [57], A. Palano [19], M. Palutan [23], Y. Pan [56],
G. Panshin [38], E. Paoletti[23], L. Paolucci [50], A. Papanestis [51], M. Pappagallo [19,g], L.L. Pappalardo [21,j],
C. Pappenheimer [59], W. Parker [60], C. Parkes [56], L. Pasquali[23,†], B. Passalacqua [21,j], G. Passaleva [22,*],
A. Pastore [19], M. Patel [55], C. Patrignani [20,h], D. Pavlenko[38], C.J. Pawley [74], A. Pazos Alvarez [40],
A. Pearce [42], M.D.P. Peco Regales[34], A. Pellegrino [32], F. Peltier[8], M. Pepe Altarelli [42], S. Perazzini [20],
D. Pereima [38], A. Pereiro Castro [40], E. Perez Trigo [40], P. Perret [9], A. Perro [42], M. Perry [56],
G. Pessina [26], K. Petridis [48], A. Petrolini [24,l], S. Petrucci [52], M. Petruzzo [25], H. Pham [62],
A. Philippov [38], R. Piandani [6], L. Pica [29,r], E. Picatoste Olloqui [39], M. Piccini [72], D. Piedigrossi[42],
B. Pietrzyk [8], G. Pietrzyk [11], M. Pili [57], N. Pillet [9], E.M. Pilorz[34], D. Pinci [30], F. Pisani [42],
M. Pizzichemi [26,n,42], V. Placinta [37], J. Plews [47], M. Plo Casasus [40], F. Polci [13,42], M. Poli Lener [23],
A. Poluektov [10], N. Polukhina [38], I. Polyakov [42], V. Polyakov[38], E. Polycarpo [2], G.J. Pomery[48],
S. Ponce [42], X. Pons[42], K. Poplawski[32], D. Popov [6,42], S. Poslavskii[38], K. Prasanth [35], D. Pratt[62],
L. Promberger [17], C. Prouve [40], V. Pugatch [46], V. Puill [11], G. Punzi [29,s], H.R. Qi [3], W. Qian [6],
N. Qin [3], S. Qu [3], R. Quagliani [43], N.V. Raab [18], B. Rachwal [34], J.H. Rademacker [48],
R. Rajagopalan[62], M. Rama [29], J.J. Ramaherison[9], M. Ramos Pernas [50], M.S. Rangel [2], F. Ratnikov [38],
G. Raven [33,42], M. Rebollo De Miguel [41], F. Redi [42], J. Reich [48], F. Reiss [56], C. Remon Alepuz[41],
Z. Ren [3], P.K. Resmi [57], F. Rethore[10], D. Reynet[11], R. Ribatti [29,r], A.M. Ricci [27], S. Ricciardi [51],
D.S. Richards[51], K. Richardson [58], M. Richardson-Slipper [52], J. Riedinger[17], K. Rinnert [54], P. Robbe [11],
G. Robertson [52], J. Rochet[42], A.B. Rodrigues [43], E. Rodrigues [54], E. Rodriguez Fernandez[40],
J.A. Rodriguez Lopez [69], P. Rodriguez Perez[56,†], E. Rodriguez Rodriguez [40], E. Roeland[32], D.L. Rolf [42],
A. Rollings [57], P. Roloff [42], V. Romanovskiy [38], M. Romero Lamas [40], A. Romero Vidal [40], P. Rosier[11],
J.D. Roth[78,†], M. Rotondo [23], J. Rovekamp[32], L. Roy[42], F. Rudnyckyj[11], M.S. Rudolph [62], T. Ruf [42],
R.A. Ruiz Fernandez [40], J. Ruiz Vidal[41], A. Ryzhikov [38], J. Ryzka [34], J.J. Saborido Silva [40],
N. Sagidova [38], N. Sahoo [47], B. Saitta [27,i], M. Salomoni [42], C. Sanchez Gras [32], F. Sanders[32],
I. Sanderswood [41], R. Santacesaria [30], C. Santamarina Rios [40], M. Santimaria [23], E. Santovetti [31,u],
A. Saputi [23], D. Saranin [38], G. Sarpis [14], M. Sarpis [70], A. Sarti [30], C. Satriano [30,t], A. Satta [31],
M. Saur [15], A. Saussac[11], D. Savrina [38], H. Sazak [9], F. Sborzacchi[23,42], L.G. Scantlebury Smead [57],
A. Scarabotto [13], S. Schael [14], S. Scherl [54], M. Schiller [53], A. Schimmel[32], H. Schindler [42],

J.D. Schipper[32], R. Schmeitz[32], M. Schmelling [16], B. Schmidt [42], S. Schmitt [14], O. Schneider [43], T. Schneider[42], A. Schopper [42], M. Schubiger [32], S. Schulte [43], M.H. Schune [11], R. Schwemmer [42], B. Sciascia [23], A. Sciuccati [42], S. Sellam [40], A. Semennikov [38], M. Senghi Soares [33], A. Sergi [24,l], N. Serra [44], J. Sestak[42], L. Sestini [28], A. Seuthe [15], P. Seyfert[42], Y. Shang [5], D.M. Shangase [78], M. Shapkin [38], I. Shchemerov [38], L. Shchutska [43], T. Shears [54], L. Shekhtman [38], Z. Shen [5], S. Sheng [4,6], M.s Sherman [62], V. Shevchenko [38], B. Shi [6], E.B. Shields [26,n], Y. Shimizu [11], E. Shmanin [38], R. Shorkin [38], J.D. Shupperd [62], B.G. Siddi [21,j], S. Siebig [17], D. Sigmund [17], S. Sigurdsson[49], R. Silva Coutinho [62], G. Simi [28], S. Simone [19,g], M. Singla [63], N. Skidmore [56], R. Skuza [17], T. Skwarnicki [62], M.W. Slater [47], K. Slattery[59], I. Slazyk [21,j], J.C. Smallwood [57], J.G. Smeaton [49], E. Smith [44], K. Smith [61], M. Smith [55], N.A. Smith [54], A. Snoch [32], L. Soares Lavra [9], J-L. Socha[11], M.D. Sokoloff [59], F.J.P. Soler [53], A. Solomin [38,48], A. Solovev [38], I. Solovyev [38], R. Song [63], F.L. Souza De Almeida [2], B. Souza De Paula [2], B. Spaan[15,†], E. Spadaro Norella [25,m], E. Spedicato [20], E. Spiridenkov[38], P. Spradlin [53], S. Squerzanti[21], V. Sriskaran [42], F. Stagni [42], M. Stahl [42], S. Stahl [42], S. Stanislaus [57], E. Steffens [d], E.N. Stein [42], O. Steinkamp [44], O. Stenyakin[38], H. Stevens [15], S. Stone [62,†], M.E. Stramaglia [43], D. Strekalina [38], Y.S Su [6], F. Suljik [57], J. Sun [27], L. Sun [68], Y. Sun [60], P. Svihra [56], P.N. Swallow [47], K. Swientek [34], S. Swientek[15], A. Szabelski [36], T. Szumlak [34], M. Szymanski [42], G Tagliente [19], Y. Tan [3], S. Taneja [56], M.D. Tat [57], M. Taurigna Quere[11], A. Terentev [44], D.F. Terront[13], F. Teubert [42], E. Thomas [42], D.J.D. Thompson [47], K.A. Thomson [54], H. Tilquin [55], V. Tisserand [9], S. T'Jampens [8], M. Tobin [4], L. Tomassetti [21,j], G. Tonani [25,m], X. Tong [5], S. Topp-Joergensen [57], D. Torres Machado [1], D.Y. Tou [3], S.M. Trilov [48], C. Trippl [43], G. Tuci [6], N. Tuning [32], A. Ukleja [36], D.J. Unverzagt [17], A. Usachov [33], A. Ustyuzhanin [38], U. Uwer [17], A. Vagner[38], V. Vagnoni [20], A. Valassi [42], S. Valat[42], G. Valenti [20], N. Valls Canudas [76], M. van Beuzekom [32], P.W. Van De Kraats [32], B. van der Heijden[32], M. Van Dijk [43], J. van Dongen[33], H. Van Hecke [61], E. van Herwijnen [55], C.B. Van Hulse [40,x], L. Van Nieuwland[32], M. van Overbeek[32], M. Van Stenis[42], M. van Veghel [32], R. Vandaele[9], R. Vazquez Gomez [39], P. Vazquez Regueiro [40], C. Vázquez Sierra [42], S. Vecchi [21], L. Veldt[32], J.J. Velthuis [48], M. Veltri [22,v], A. Venkateswaran [43], H. Verkooijnen[32], M. Veronesi [32], M. Vesterinen [50], J.V. Viana Barbosa[42], D. Vieira [59], M. Vieites Diaz [43], K.J. Viel[42], X. Vilasis-Cardona [76], E. Vilella Figueras [54], A. Villa [20], P. Vincent [13], W. Vink[32], A. Vitkovskiy[32], V. Volkov [38], F.C. Volle [11], D. vom Bruch [10], B. Voneki[42], O. Vorbach[17], A. Vorobyev [38], V. Vorobyev[38], N. Voropaev [38], K. Vos [74], G. Vouters[8], C. Vrahas [52], W. Walet[32], J. Walsh [29], E.J. Walton [63], G. Wan [5], C. Wang [17], G. Wang [7], J. Wang [5], J. Wang [4], J. Wang [3], J. Wang [68], M. Wang [25], R. Wang [48], X. Wang [66], Y. Wang [7], Z. Wang [44], Z. Wang [3], Z. Wang [6], J.A. Ward [50,63], K. Warda[15], N.K. Watson [47], D. Websdale [55], J. Webster [52], Y. Wei [5], B.D.C. Westhenry [48], D.J. White [56], M. Whitehead [53], D. Wieczorek[15], A.R. Wiederhold [50], D. Wiedner [15], G. Wilkinson [57], M.K. Wilkinson [59], I. Williams[49], M. Williams [58], M.R.J. Williams [52], R. Williams [49], F.F. Wilson [51], J. Wimberley [60], B. Windelband[17], W. Wislicki [36], M. Witek [35], L. Witola [17], M. Wlochal [14], C.P. Wong [61], M. Wormald[54], G. Wormser [11], S.A. Wotton [49], K. Wraight [53], H. Wu [62], J. Wu [7], K. Wyllie [42], Z. Xiang [6], Y. Xie [7], A. Xu [5], J. Xu [6], L. Xu [3], L. Xu [3], M. Xu [50], Q. Xu[6], Z. Xu [9], Z. Xu [6], D. Yang [3], S. Yang [6], X. Yang [5], Y. Yang [6], Z. Yang [5], Z. Yang [60], L.E. Yeomans [54], V. Yeroshenko [11], H. Yeung [56], H. Yin [7], J. Yu [65], X. Yuan [62], E. Zaffaroni [43], M. Zavertyaev [16], M. Zdybal [35], O. Zenaiev [42], M. Zeng [3], C. Zhang [5], D. Zhang [7], L. Zhang [3], S. Zhang [65], S. Zhang [5], Y. Zhang [5], Y. Zhang[57], Y. Zhao [17], A. Zharkova [38], A. Zhelezov [17], Y. Zheng [6],

T. Zhou [5], X. Zhou [6], Y. Zhou [6], V. Zhovkovska [11], X. Zhu [3], X. Zhu [7], Z. Zhu [6], V. Zhukov [14,38],
V. Zivkovic[32], Q. Zou [4,6], S. Zucchelli [20,h], D. Zuliani [28], G. Zunica [56], S. Zvyagintsev[38]

1 *Centro Brasileiro de Pesquisas Físicas (CBPF), Rio de Janeiro, Brazil*

2 *Universidade Federal do Rio de Janeiro (UFRJ), Rio de Janeiro, Brazil*

3 *Center for High Energy Physics, Tsinghua University, Beijing, China*

4 *Institute Of High Energy Physics (IHEP), Beijing, China*

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

6 *University of Chinese Academy of Sciences, Beijing, China*

7 *Institute of Particle Physics, Central China Normal University, Wuhan, Hubei, China*

8 *Université Savoie Mont Blanc, CNRS, IN2P3-LAPP, Annecy, France*

9 *Université Clermont Auvergne, CNRS/IN2P3, LPC, Clermont-Ferrand, France*

10 *Aix Marseille Univ, CNRS/IN2P3, CPPM, Marseille, France*

11 *Université Paris-Saclay, CNRS/IN2P3, IJCLab, Orsay, France*

12 *Laboratoire Leprince-Ringuet, CNRS/IN2P3, Ecole Polytechnique, Institut Polytechnique de Paris, Palaiseau, France*

13 *LPNHE, Sorbonne Université, Paris Diderot Sorbonne Paris Cité, CNRS/IN2P3, Paris, France*

14 *I. Physikalisches Institut, RWTH Aachen University, Aachen, Germany*

15 *Fakultät Physik, Technische Universität Dortmund, Dortmund, Germany*

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

17 *Physikalisches Institut, Ruprecht-Karls-Universität Heidelberg, Heidelberg, Germany*

18 *School of Physics, University College Dublin, Dublin, Ireland*

19 *INFN Sezione di Bari, Bari, Italy*

20 *INFN Sezione di Bologna, Bologna, Italy*

21 *INFN Sezione di Ferrara, Ferrara, Italy*

22 *INFN Sezione di Firenze, Firenze, Italy*

23 *INFN Laboratori Nazionali di Frascati, Frascati, Italy*

24 *INFN Sezione di Genova, Genova, Italy*

25 *INFN Sezione di Milano, Milano, Italy*

26 *INFN Sezione di Milano-Bicocca, Milano, Italy*

27 *INFN Sezione di Cagliari, Monserrato, Italy*

28 *Università degli Studi di Padova, Università e INFN, Padova, Padova, Italy*

29 *INFN Sezione di Pisa, Pisa, Italy*

30 *INFN Sezione di Roma La Sapienza, Roma, Italy*

31 *INFN Sezione di Roma Tor Vergata, Roma, Italy*

32 *Nikhef National Institute for Subatomic Physics, Amsterdam, Netherlands*

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

34 *AGH - University of Science and Technology, Faculty of Physics and Applied Computer Science, Kraków, Poland*

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

36 *National Center for Nuclear Research (NCBJ), Warsaw, Poland*

37 *Horia Hulubei National Institute of Physics and Nuclear Engineering, Bucharest-Magurele, Romania*

38 *Affiliated with an institute covered by a cooperation agreement with CERN*

39 *ICCUB, Universitat de Barcelona, Barcelona, Spain*

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

41 *Instituto de Fisica Corpuscular, Centro Mixto Universidad de Valencia - CSIC, Valencia, Spain*

42 *European Organization for Nuclear Research (CERN), Geneva, Switzerland*

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

44 *Physik-Institut, Universität Zürich, Zürich, Switzerland*

45 *NSC Kharkiv Institute of Physics and Technology (NSC KIPT), Kharkiv, Ukraine*

46 *Institute for Nuclear Research of the National Academy of Sciences (KINR), Kyiv, Ukraine*

47 *University of Birmingham, Birmingham, United Kingdom*

48 *H.H. Wills Physics Laboratory, University of Bristol, Bristol, United Kingdom*

49 *Cavendish Laboratory, University of Cambridge, Cambridge, United Kingdom*

50 *Department of Physics, University of Warwick, Coventry, United Kingdom*

51 *STFC Rutherford Appleton Laboratory, Didcot, United Kingdom*

52 *School of Physics and Astronomy, University of Edinburgh, Edinburgh, United Kingdom*

53 *School of Physics and Astronomy, University of Glasgow, Glasgow, United Kingdom*

54 *Oliver Lodge Laboratory, University of Liverpool, Liverpool, United Kingdom*

55 *Imperial College London, London, United Kingdom*

56 *Department of Physics and Astronomy, University of Manchester, Manchester, United Kingdom*

57 *Department of Physics, University of Oxford, Oxford, United Kingdom*

58 *Massachusetts Institute of Technology, Cambridge, MA, United States*

59 *University of Cincinnati, Cincinnati, OH, United States*

60 *University of Maryland, College Park, MD, United States*

61 *Los Alamos National Laboratory (LANL), Los Alamos, NM, United States*

62 *Syracuse University, Syracuse, NY, United States*

63 *School of Physics and Astronomy, Monash University, Melbourne, Australia, associated to* [50]

64 *Pontifícia Universidade Católica do Rio de Janeiro (PUC-Rio), Rio de Janeiro, Brazil, associated to* [2]

65 *Physics and Micro Electronic College, Hunan University, Changsha City, China, associated to* [7]

66 *Guangdong Provincial Key Laboratory of Nuclear Science, Guangdong-Hong Kong Joint Laboratory of Quantum Matter, Institute of Quantum Matter, South China Normal University, Guangzhou, China, associated to* [3]

67 *Lanzhou University, Lanzhou, China, associated to* [4]

68 *School of Physics and Technology, Wuhan University, Wuhan, China, associated to* [3]

69 *Departamento de Fisica, Universidad Nacional de Colombia, Bogota, Colombia, associated to* [13]

70 *Universität Bonn - Helmholtz-Institut für Strahlen und Kernphysik, Bonn, Germany, associated to* [17]

71 *Eotvos Lorand University, Budapest, Hungary, associated to* [42]

72 *INFN Sezione di Perugia, Perugia, Italy, associated to* [21]

73 *Van Swinderen Institute, University of Groningen, Groningen, Netherlands, associated to* [32]

74 *Universiteit Maastricht, Maastricht, Netherlands, associated to* [32]

75 *Tadeusz Kosciuszko Cracow University of Technology, Cracow, Poland, associated to* [35]

76 *DS4DS, La Salle, Universitat Ramon Llull, Barcelona, Spain, associated to* [39]

77 *Department of Physics and Astronomy, Uppsala University, Uppsala, Sweden, associated to* [53]

78 *University of Michigan, Ann Arbor, MI, United States, associated to* [62]

79 *Institute for Scintillation Materials, Kharkiv, Ukraine*

80 *Taras Schevchenko University of Kyiv, Faculty of Physics, Kyiv, Ukraine*

[a] *Universidade de Brasília, Brasília, Brazil*

[b] *Central South U., Changsha, China*

[c] *Hangzhou Institute for Advanced Study, UCAS, Hangzhou, China*

[d] *Friedrich-Alexander-Universitat Erlangen-Nurnberg (FAU), Erlangen-Nurnberg, Germany*

[e] *Excellence Cluster ORIGINS, Munich, Germany*

[f] *Universidad Nacional Autónoma de Honduras, Tegucigalpa, Honduras*

[g] *Università di Bari, Bari, Italy*

[h] *Università di Bologna, Bologna, Italy*

[i] *Università di Cagliari, Cagliari, Italy*

[j] *Università di Ferrara, Ferrara, Italy*

[k] *Università di Firenze, Firenze, Italy*

[l] *Università di Genova, Genova, Italy*

[m] *Università degli Studi di Milano, Milano, Italy*

[n] *Università di Milano Bicocca, Milano, Italy*

[o] *Università di Modena e Reggio Emilia, Modena, Italy*

[p] *Università di Padova, Padova, Italy*

[q] *Università di Perugia, Perugia, Italy*

[r] *Scuola Normale Superiore, Pisa, Italy*

[s] *Università di Pisa, Pisa, Italy*

[t] *Università della Basilicata, Potenza, Italy*

[u] *Università di Roma Tor Vergata, Roma, Italy*

$^v$ *Università di Urbino, Urbino, Italy*

$^w$ *AGH - University of Science and Technology, Faculty of Computer Science, Electronics and Telecommunications, Kraków, Poland*

$^x$ *Universidad de Alcalá, Alcalá de Henares, Spain*

$^\dagger$ *Deceased*

$^*$ *Corresponding author*