

Run-3 offline data processing and analysis at LHCb

N. Skidmore,^{a,*} E. Rodrigues^b and P. Koppenburg^c

^a*Department of Physics and Astronomy, University of Manchester, Manchester, United Kingdom*

^b*Oliver Lodge Laboratory, University of Liverpool, Liverpool, United Kingdom*

^c*Nikhef National Institute for Subatomic Physics, Amsterdam, Netherlands*

E-mail: nicola.skidmore@cern.ch, Eduardo.Rodrigues@cern.ch,

patrick.koppenburg@cern.ch

The LHCb detector is undergoing a comprehensive upgrade for data taking in the LHC's Run 3, which is scheduled to begin in 2022. The increased data rate in Run 3 poses significant data-processing and handling challenges for the LHCb experiment. The offline computing and dataflow model is consequently also being upgraded to cope with the factor 30 increase in data volume and associated demands of user-data samples of ever-increasing size. Coordinating these efforts is the charge of the newly created Data Processing & Analysis (DPA) project. The DPA project is responsible for ensuring the LHCb experiment can efficiently exploit the Run 3 data, processing the data from the online system with central skimming/slimming (a process known as "Sprucing") and subsequently producing analyst-level ntuples with a centrally managed production system (known as "Analysis Productions") utilising improved analysis tools and infrastructure for continuous integration and validation. It is a multi-disciplinary project involving collaboration between computing experts, trigger experts and physics analysis experts.

*** *The European Physical Society Conference on High Energy Physics (EPS-HEP2021), ****

*** *26-30 July 2021 ****

*** *Online conference, jointly organized by Universität Hamburg and the research center DESY ****

*Speaker

1. Introduction to the DPA project

The LHCb experiment is one of the four main experiments collecting data from proton-proton collisions at the Large Hadron Collider (LHC) [1]. It is a forward arm spectrometer specialising in the decays of beauty and charm hadrons. During the LHC's first and second data-taking runs LHCb collected data corresponding to an integrated luminosity of 9 fb^{-1} , equating to $> 10^{12} b\bar{b}$ pairs in the LHCb acceptance. In 2022 the LHC will commence its third run — Run 3 — of data-taking. LHCb has undergone a comprehensive upgrade in anticipation of a factor five increase in luminosity, resulting in an increase in the volume of data collected by LHCb by more than a factor 30 compared to previously (taking into account a factor three due to the increased average event size and a factor two due to higher trigger efficiencies) [2]. The Upgrade era therefore not only poses data collection (Online) challenges but also significant data processing and handling challenges once the data has been taken. Coordination of the activities and developments required to face these challenges called on the formation of a new software project, namely the Data Processing & Analysis (DPA) project at LHCb; its remit is the processing of the data from when it leaves the fully software trigger [3] to analyst-level data-structures containing derived quantities.

The DPA project comprises six dependent but distinct work packages:

1. Sprucing – Centralised offline data selection, streaming and formatting for data leaving the trigger.
2. Analysis productions – nTuple production defined by the analyst and run centrally using the DIRAC [4] transformation system with maximum automation.
3. Offline analysis tools – Offline, thread safe, analysis application looking to create consistency between tools used online and offline.
4. Innovative analysis techniques – R&D work package for innovative analysis techniques, aiming for successful proofs of concept that may become mainstream in LHCb.
5. Legacy software and data – Maintenance of, and support for, legacy Runs 1 & 2 software and data samples.
6. Analysis preservation and open data – Guidelines and tools for analysis preservation in LHCb, management and release of LHCb data to the CERN Open Data portal.

The LHCb Upgrade offline dataflow is shown in Figure 1, where the remit of DPA is identified as "Offline processing". The work packages and their activities will be expanded upon in the following sections.

2. WP1 – Sprucing

In Run 3 and beyond the default model for LHCb is to persist only a custom set of physics objects and their reconstruction to the Offline system for further analysis, the so-called Turbo Selective Persistence (Turbo SP) stream [6]. The rest of the event is discarded, reducing the event size to 4–16 kB [7]. No further trimming or skimming of this data is required and only a

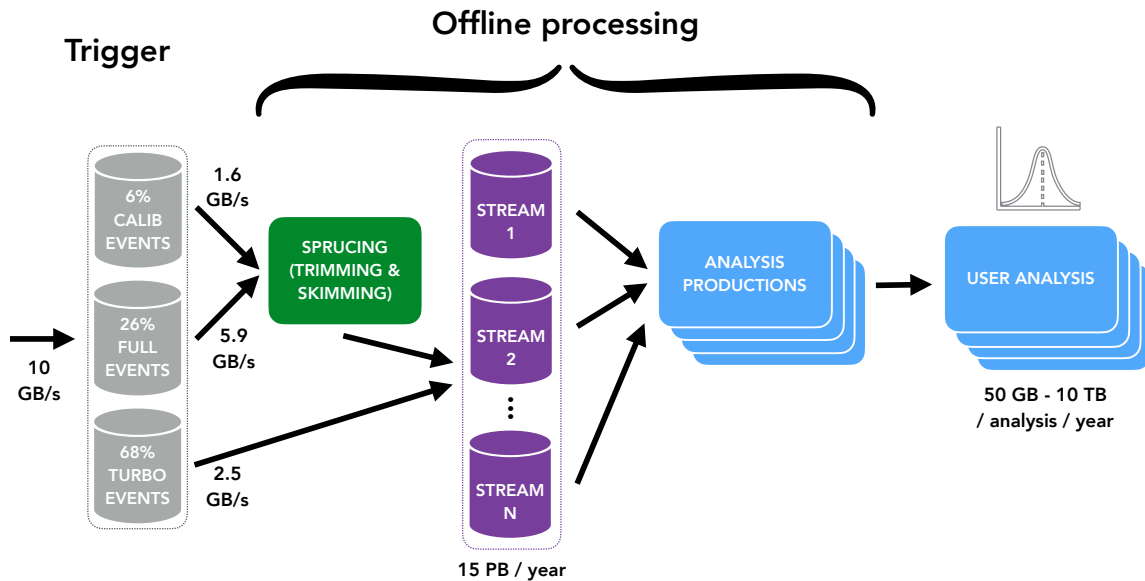


Figure 1: LHCb dataflow showing remit of DPA project. Figure taken from Ref. [5].

simple reformatting is performed before these events are saved to disk where the data is available to analysts. There are numerous physics cases, however, where (almost) the whole event needs to be accessible for later analysis. This Full data stream is saved to tape whereby centrally-managed offline re-processings and selections on full physics events can be periodically run and where the data is slimmed and skimmed so that it can be saved to disk. In Run 3 this is known as the Sprucing processing stage, as seen in Figure 1.

The Sprucing shares the same software framework and application as the trigger, namely the MOORE application [8]. Furthermore, the same algorithms and tools are shared between the trigger, Sprucing and the offline analysis software project DAVINCI [9], namely the THOR [10] based selection and combinatorics algorithms. Trigger lines and Sprucing lines are, by design, interchangeable.

Similar to the Stripping of Runs 1 and 2 [11], the Sprucing will run concurrently with data taking. Re-Sprucing campaigns will typically take place in end-of-year LHC shutdowns whereby the data will be staged from tape and selections – also, for instance, calibrations – can be re-run.

3. WP2 – Analysis Productions

The Turbo (SP) data and the Full data post Sprucing are split into multiple streams (files), which are directly accessible to analysts. In the legacy data model these would be processed by submitting user jobs to LHCbDIRAC [12] 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 $O(10^3)$. While this model works well for smaller datasets, 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.

Analysis Productions are an extension of the LHCbDIRAC transformation system, which has been primarily used for the centralised processing of LHCb data and simulation up till now. Analysis Productions are submitted declaratively via YAML files by providing the job configuration and bookkeeping query for the input data, the latter enabling 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 system, enabling high quality analysis preservation and additional safety checks to be performed. Interactive analysis work can directly leverage 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 pre-submission 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 user-prepared configurations are correct, extensive tests are run on the GitLab Continuous Integration platform prior to getting them approved by physics working group liaisons. The productions and all details of the validation tests are summarised on a dedicated website.

4. WP3 – Offline analysis tools

Whilst in Runs 1 and 2 the trigger software was largely based on the previously designed offline analysis framework, the opposite is true for the upgraded LHCb software stack. As much as possible, the offline analysis framework is based on the software developments made for the trigger, which constitute the MOORE application stack. This approach avoids duplication of effort, profits from these huge developments towards modern and resources-efficient code, and guarantees a similar look-and-feel of the software to be used by analysts. Most importantly, all basic building blocks are shared between the trigger, the Sprucing and the offline analysis, thus guaranteeing that the same software is consistently used to compute the same quantities online and offline.

Following this approach, the DAVINCI analysis software is descoped (compared to its legacy version) to only take care of the production of output tuples from input data (including simulation). Its core is a new tupling algorithm that fills measured quantities using *functors* provided by the MOORE THOR [10] and LoKi frameworks [13]. These are the same functors called in the trigger and Sprucing algorithms, thus ensuring a one-to-one correspondence between applied selection requirements and observables used offline. It should be emphasised that since no offline reconstruction is foreseen in the Run-3 computing model, the input data objects (tracks, clusters, *etc.*) are identical online and offline [2, 14].

5. WP4 – Innovative analysis techniques

Work package 4 of the DPA project is a think-tank for innovative analysis techniques and exploitation of new analysis facilities with heterogeneous computing resources. In Run 3 the first stage of LHCb's trigger system will run on GPUs [15]. This means that during detector down-time

LHCb has the opportunity to exploit a GPU farm for payloads such as detector simulation and computationally intensive analysis tasks such as amplitude analysis fitting studies.

6. WP5 – Legacy software and data

The Runs 1 and 2 datasets will continue to be of significant interest to LHCb analysts and the Particle Physics community for many years to come. DPA's WP5 coordinates regular reprocessing of these so-called legacy data samples and maintains the software and tools to analyse them. WP5 collaborates closely with WP6 to meet LHCb's commitments to open data access.

7. WP6 – Analysis preservation and open data

To facilitate flexibility and creativity for LHCb analysts only minimal constraints are placed on analysis code. While commonly used libraries, such as RooFit [16], Scikit-HEP project packages [17], etc., are available through CVMFS [18] or can be managed using the conda [19] package manager, analysts write custom scripts and oftentimes 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 at gitlab.cern.ch.

The centrally produced data is preserved on the distributed computing infrastructure and catalogued in a dedicated LHCb bookkeeping system. In 2020 LHCb ratified the CERN Open Data policy [20, 21]. In accordance with the access policies outlined in the documents, LHCb will make the output of the Turbo-lines as well as the output of the Sprucing available to the public through the CERN Open Data portal. Broadly speaking, the LHC experiments agreed to make available 50% of their data samples after 5 years, rising to 100% after 10 years from the end of the LHC running period. The release of the Run-1 data and related documentation to the Open Data portal is presently under preparation.

The software necessary to read the data is open-source and preserved as CVMFS releases. However, due to the lack of an analysis-object-only data format, only a small fraction of LHCb data can be stored on the Open Data Portal. This motivates providing secure access to replicas of the data stored on the grid via an ntupling service for third parties. A web interface is under development, which will allow users to configure analysis production jobs with minimal knowledge of the LHCb software. This "Ntuple-Wizard" will also allow fine-grained control over the level of access given to the public, as well as act as a firewall protecting against the injection of nefarious code into the production system.

8. Summary

The Data Processing & Analysis (DPA) project at LHCb is coordinating offline infrastructure developments to ensure full and efficient exploitation of data in Run 3 and beyond. DPA builds on the two main ideas of centralised skimming and trimming of trigger outputs and centralised analysis productions for physics working groups and analysts.

References

- [1] A. A. Alves Jr. et al. “The LHCb detector at the LHC”. In: *JINST* 3 (2008), S08005. DOI: [10.1088/1748-0221/3/08/S08005](https://doi.org/10.1088/1748-0221/3/08/S08005).
- [2] LHCb Collaboration. *Computing Model of the Upgrade LHCb experiment*. Tech. rep. CERN-LHCC-2018-014, LHCb-TDR-018. Geneva: CERN, May 2018. URL: <https://cds.cern.ch/record/2319756>.
- [3] R. Aaij et al. “A comprehensive real-time analysis model at the LHCb experiment”. In: *JINST* 14 (Mar. 2019). arXiv: [1903.01360](https://arxiv.org/abs/1903.01360). URL: <http://cds.cern.ch/record/2665946>.
- [4] A. Tsaregorodtsev et al. “DIRAC3: The new generation of the LHCb grid software”. In: *J. Phys. Conf. Ser.* 219 (2010), p. 062029. DOI: [10.1088/1742-6596/219/6/062029](https://doi.org/10.1088/1742-6596/219/6/062029).
- [5] LHCb Collaboration. *RTA and DPA dataflow diagrams for Run 1, Run 2, and the upgraded LHCb detector*. Tech. rep. LHCb-FIGURE-2020-016. Sept. 2020. URL: <https://cds.cern.ch/record/2730181>.
- [6] R. Aaij et al. “Tesla : an application for real-time data analysis in High Energy Physics”. In: *Comput. Phys. Commun.* 208 (Apr. 2016). arXiv: [1604.05596](https://arxiv.org/abs/1604.05596). URL: <https://cds.cern.ch/record/2147693>.
- [7] Alejandro Alfonso Albero et al. *Upgrade trigger selection studies*. Tech. rep. Geneva: CERN, Sept. 2019. URL: <https://cds.cern.ch/record/2688423>.
- [8] LHCb Collaboration. *Moore application GitLab repository*. <https://gitlab.cern.ch/lhcb/Moore>.
- [9] LHCb Collaboration. *DaVinci application GitLab repository*. <https://gitlab.cern.ch/lhcb/DaVinci/>.
- [10] LHCb Collaboration. *ThOr functors*. https://lhcbdoc.web.cern.ch/lhcbdoc/moore/master/selection/thor_functors.html.
- [11] LHCb Collaboration. *LHCb computing: Technical Design Report*. Tech. rep. CERN-LHCC-2005-019, LHCb-TDR-11. Geneva, 2005. URL: <http://cds.cern.ch/record/835156>.
- [12] *LHCbDIRAC Documentation*. <https://lhcb-dirac.readthedocs.io>.
- [13] I Belyaev et al. “PYTHON-based Physics Analysis Environment for LHCb”. In: LHCb-PROC-2004-021 (Oct. 2004), 4 p. DOI: [10.5170/CERN-2005-002.377](https://doi.org/10.5170/CERN-2005-002.377). URL: <https://cds.cern.ch/record/1490626>.
- [14] LHCb collaboration. *LHCb Trigger and Online Upgrade Technical Design Report*. Tech. rep. CERN-LHCC-2014-016. Geneva: CERN, 2014.
- [15] R. Aaij et al. “Allen: A High-Level Trigger on GPUs for LHCb”. In: *Computing and Software for Big Science* 4.1 (Apr. 2020). ISSN: 2510-2044. DOI: [10.1007/s41781-020-00039-7](https://doi.org/10.1007/s41781-020-00039-7). URL: <http://dx.doi.org/10.1007/s41781-020-00039-7>.
- [16] Wouter Verkerke and David P. Kirkby. “The RooFit toolkit for data modeling”. In: *eConf C0303241* (2003). Ed. by L. Lyons and Muge Karagoz, MOLT007. arXiv: [physics/0306116](https://arxiv.org/abs/physics/0306116).

- [17] Eduardo Rodrigues et al. “The Scikit HEP Project overview and prospects”. In: *EPJ Web of Conferences* 245 (2020). Ed. by C. Doglioni et al., p. 06028. ISSN: 2100-014X. DOI: 10.1051/epjconf/202024506028. URL: <http://dx.doi.org/10.1051/epjconf/202024506028>.
- [18] *CernVM-FS Documentation*. <https://cvmfs.readthedocs.io>.
- [19] *Conda Documentation*. <https://docs.conda.io/>.
- [20] *CERN Open Data Policy for the LHC Experiments*. Tech. rep. Geneva: CERN, Nov. 2020. URL: <https://cds.cern.ch/record/2745133>.
- [21] CERN Scientific Information Policy Board. “CERN Open Data Policy for LHC Experiments: implementation plan”. In: (Nov. 2020). URL: <https://cds.cern.ch/record/2745081>.