

Run-3 offline data processing and analysis at LHCb

N. Skidmore on behalf of the LHCb collaboration

EPS-HEP2021 conference

July 2021



The University of Manchester





LHCb Upgrade I - Hardware and software Upgrade

New scintillating fibre tracker (SciFi) New electronics (CALO, MUON)



New fully software trigger - all sub-detectors must read out at 40MHz HLT1 to run on GPUs

New optics and photodetectors (RICH) In Run 3 the data rate will increase by factor 30 compared to run 2...

See talk by Marianna

New silicon

upstream tracker

(UT)

New pixel vertex detector — (VELO)

Data Processing and Analysis (DPA) in Run 3

- Increased data rate in Run 3 poses significant Offline data processing and analysis challenges
- Coordination of these activities by DPA project
 - Software project on same level as detector projects



DPA project Offline processing/selections/analysis

| WP1 | - Sprucing |
|-----|-------------------------------------|
| WP2 | - Analysis Productions |
| WP3 | - Offline Analysis Tools |
| WP4 | - Innovative Analysis Techniques |
| WP5 | - Legacy Software & Data |
| WP6 | - Analysis Preservation & Open Data |

Six work packages



WP1: Sprucing

In Run 3 event persistency is customisable depending on the Physics involved



| stream | event size | event rate | rate | throughput | bandwidth |
|--------|------------|------------|----------|------------|-----------|
| | (kB) | (kHz) | fraction | (GB/s) | fraction |
| FULL | 70 | 7.0 | 65% | 0.49 | 75% |
| Turbo | 35 | 3.1 | 29% | 0.11 | 17% |
| TurCal | 85 | 0.6 | 6% | 0.05 | 8% |
| total | 61 | 10.8 | 100% | 0.65 | 100% |

In Run 2 - 32% of physics events went to TURBO



WP1: Sprucing

| Output data | Turbo | Physics channels left in FULL |
|-------------|------------------|---|
| volume | physics fraction | |
| 10 GB/s | 73% | EW, high PT, (semi)leptonic and some hadronic B-physics, leptonic charm decays and general LFV searches |
| 7.5 GB/s | 87% | EW, high PT, some leptonic B- physics, some LFV searches and leptonic searches |
| 5 GB/s | 99% | None |

For physics left in FULL stream - Sprucing:

- Utilise cheap tape storage for bulk of bandwidth (FULL stream)
- Perform central offline slimming/skimming



Sprucing

A further offline stage of data reduction/selection between **tape** and **disk** storage when HLT2 line throughput is too large to go straight to disk.

Utilise same selection framework as HLT2

WP2 : Analysis productions (AnaProd)

In Run 1 + 2 analysts created nTuples individually from data on disk using Ganga... does not scale well for Run 3

- 1000s of faulty jobs can be submitted instantly (10% of user jobs fail)
- Time consuming O(weeks) for Run 1 + 2 tuples failed jobs re-submitted manually by user
- No analysis preservation infrastructure

Analysis productions submit nTupling jobs centrally using DIRAC transformation System (AnaProds already in use for legacy data)

- Does not require analyst to babysit grid jobs
- Options tested automatically upon push to GitLab (CI). Final approval must be given from PAWG liaisons
- Job details/configuration/logs automatically preserved in LHCb bookkeeping/EOS
- Automated error interpretation/advice
- Results displayed on webpage



WP2 : Analysis productions (AnaProd)

Automated error interpretation



Options tested on GitLab Cl

| SILIO Presens v | Search or jump to Q | 0.8 1.00 0.00 0 | · • |
|-----------------------|---|----------------------------|---------|
| A AnalysisProductions | LHEb Data Packages > AnalysiaFroductions > Jobs > #10108707 | test | Retry |
| D Project overview | 💿 passed 🛛 Jab #10108707 triggered 6 days ago by 🍘 Chris Burr | | |
| | | Timeout: 1d (from project) | econcis |
| E Hebourpek | B D - + | Dunner (#5453) | |
| Chimmer (6) | 1 Furning with gitlab_runner_api 0.1.dev42+p457db45 | Tars: Day second | |
| | 2 3570:Creating new pipeline for ID 1950460 | and a second second | |
| 11 Merge Requests (8) | Merce Recuests () 3 ALMXYS:Results will be available at https://lhcb-amalysis-productions.web.corn.ch/1555448/ | | |
| | 4 3MPD:Creating production D02HLStarterkit | Commit 55210090 15 In | 110 |
| ≪ ci/co | INFO:Generating configuration options for INERM_Starterkit 2016_MagBown_PromptMC_D02000 (1:1/2) | Minimise repetition | |
| | 3%F0:Generating configuration options for 140000_Starterkit 2416_HegUp_PromptMC_04000K (2.1/2) | | |
| Pipelines | 7 ALWAYS:Submitting jobs for D029H_Starterkit | Pipeline #1968460 for | |
| laba. | 3170:Submitting test for 00204_Starterkit 2016_MagUp_FrompUMC_0020X 12/21 | djubite/starterkit-pre | ectice |
| 3446 | 3xF0:Submitting test for DB2HL_Starterkit 2016_PagDave_PromptHC_DB29K 11/21 | | |
| Schecklos | 10 INFO:Submitted 0DRAC job for D02H0_Starterkit 2016_MagUp_PromptHC_002HX with ID 465432666 | 1410 | |
| | 11 INFO:Subsitted 0D0AC job for D0200_Starterkit 2016_MagBown_PromptMC_0020X with ID 405432667 | | |
| Operations | 12 ALMANS:2 jobs still running at 2020-09-23132:45:07.464027 | → ⊘test | |
| | 13 ALMANS-2 jobs still running at 2020-09-23732:45:55.646284 | | |
| Packages & Registries | 14 ALMANS-2 jobs still running at 2020-09-23732:46:36.042122 | | |
| | 15 ALMANS:2 jobs still running at 2020-09-23732:47:16.577247 | | |
| Le Analytics | 10 ALMANS:2 jobs still running at 2020-09-23732:48:01.112105 | | |
| | 17 #LMKN3:2 jebs still running at 2820-89-23132:48:48.594881 | | |
| A Members | 18 ALMANSI2 jobs still running at 2020-09-23732149123-345684 | | |
| | 19 #14475-3 jobs still running at 2020-06-23712:50:04.024533 | | |
| ♦ Settings | 20 ALMANS:2 jobs still running at 2020-09-23112:50:54.779824 | | |
| | 21 #10#05-2 jobs still running at 2020-09-23122:51:52.300850 | | |
| | 22 ADMINIST 2 1985 STLLL FURNING OT 2020-49-25132152133.424767 | | |
| | 23 ADDATED 3445 STOLL CONTING AT 2020-09-281321532153.346984 | | |
| | 24 ADADA 944 BOLL 100009 BL 200409-201004 BS35558 | | |





Job configuration/logs preserved in bookkeeping

| Menu 🤇 | | | 1 |
|---------------------------------------|---|--|--|
| e 📃 o | Analysis Productions [Unitled 1] \times | | |
| • • • • • • • • • • • • • • • • • • • | FBCD F - Andread Statut G 10 G 11 G 12 Andread SCH Hammel 13 G 14 G 15 G 16 G 17 G 18 G 19 G 10 G 10 G | in a set and the second second is specify insubstance in the second seco | 9220 rosty 92 argust FFNs - Comp - Collect AD - Comp - Collect AD |
| | Control (1) | More Table 2012 Control Tab | Coop: 5000340 Coop: 5000340 2222 Tot? 57 o/Dut Phis Coop: 5000340 Coop: 5000340 |

WP3 : Offline analysis tools



In Run 1+2 nTuples used "TupleTools" from DaVinci (user analysis) application

- "TupleTools" create and save variable branches for typical use cases eg. TupleToolTrackInfo
- Very easy to implement but adds lots of redundant branches can easily save 500+ variables
- 500GB 10TB of data for a single Run 1+2 analysis nTuples tend to be only used for one analysis

For Run 3 complete redesign of DaVinci framework - FUNTUPLE

- Same thread-safe functors as in HLT2/Sprucing used to create light-weight nTuples
 - Consistency between Online and Offline selections/tools/algorithms
- Analyst has full control over exactly which variables for which particles are persisted in nTuple
- AnaProds will run analysts' DaVinci options





WP4 : R&D Innovative analysis techniques

Think tank for innovative analysis techniques and exploitation of new analysis facilities with heterogeneous computing resources (GPU/CPU/FPGA)

Worldwide LHC Computing Grid (WLCG) consists of ~ 1M CPU cores over 170 sites

- Most sites have no GPUs yet push towards High Performance Computing (HPC) centers providing large GPU resources
- Potential to utilise LHCb's HLT1 GPU farm during detector downtime

Need development such that significant LHCb payloads can run on GPUs

- Use advanced algorithms such as Generative Adversarial Networks (GANs) to train models describing LHCb sub-detector response GPUs speed up GAN training *Ultra-fast simulation*
- Users using GPUs for analysis, e.g. TensorFlow for model fitting (Zfit) particularly for complex amplitude analysis models with large statistics

First investigations into use of Quantum Machine Learning for jet tagging (see backup)



In Run 3 LHCb will produce ~ 15PB of data on disk per year

Real data will dominate disk storage but simulation will dominate CPU needs - 90% of total offline CPU resources



WP6: Analysis preservation and open data

Run 1 LHCb data to be released on CERN Open Data portal

- Development of Open Data nTuple wizard
 - Auto-generates options from intuitive user input no knowledge of LHCb software required
 - Launches AnaProd (See slide 8-9)
 - Returns nTuple to user
- Much smaller storage and bandwidth requirements on Open Data Portal
- Will be used by collaboration members to analyse Run 1+2 data



Analysis preservation has been in place at LHCb since 2017 - building on this

...analysis code should be preserved in a long-term archive such as a physics analysis gitlab group, the input ntuples should also be preserved in a long-term archive (EOS) and sufficient documentation to enable a (technically competent) LHCb member to run the code in a standard environment such as lxplus should be included with the code.



- Developing tag-based access to AnaProd output removing need for copying files/hard coding paths
- Support for snakemake (significant adoption in LHCb) in REANA

Summary



- LHCb will have to process data offline an order of magnitude larger than in Run 2 in 2022
- LHCb is progressing well to meet the Offline demands that run 3 will bring, coordinated by the **DPA** project



Backup

WP4 : R&D Innovative analysis techniques

Quantum Machine Learning model for $b vs \overline{b}$ jet tagging

- Quantum Circuit with parameterized gates
- Variables from jet particles encoded in a quantum state
- The state is processed by trainable quantum gates
- Measurements on the final state are mapped to labels (b or \overline{b})
- Parameters (of the gates) are optimized via a Gradient Descent minimization of a cost function (training)





Quantum models can allow the study of correlations among particles inside the jet meaning insight on jet substructure and better identification!

For more information see PyHEP 2021 talk