# Offline data processing and analysis at LHCb in the 2020s

## International Conference of High Energy Physics 2022
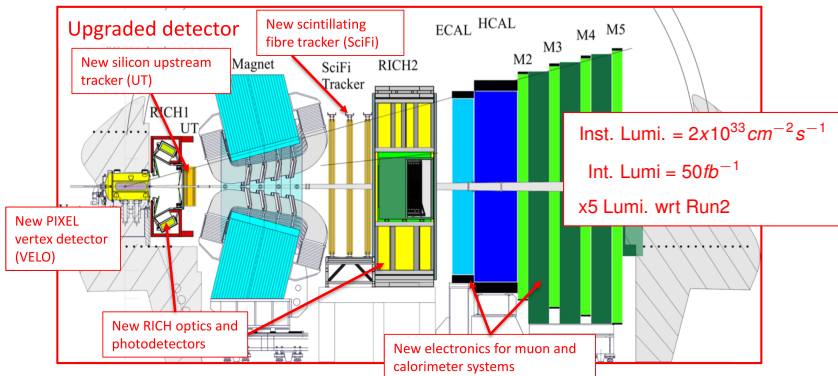
**Davide Fazzini**
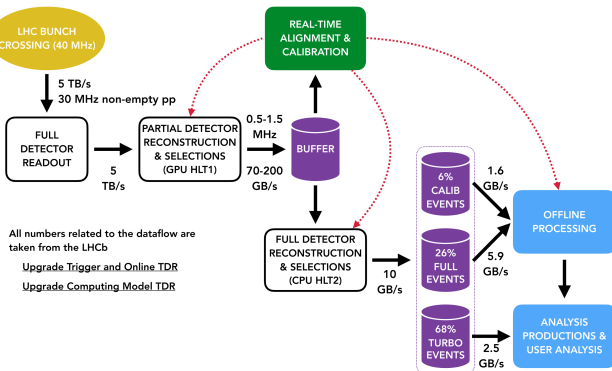on behalf of the LHCb Collaboration

July 6-13 2022, Bologna, Italy

- LHCb is a single-arm forward spectrometer redesigned for Run3:
  - about 95% of the sub-detectors is completely new
  - high-precision tracking and & vertexing
  - excellent PID & hadron separation
- Physics programme in practice far more general
  - Electroweak, Exotica, LLPs, Fixed-target, heavy ions
- **Think of it as a more general purpose detector!**



Upgraded detector

New scintillating fibre tracker (SciFi)

New silicon upstream tracker (UT)

New PIXEL vertex detector (VELO)

New RICH optics and photodetectors

New electronics for muon and calorimeter systems

Inst. Lumi. $= 2x10^{33} cm^{-2} s^{-1}$

Int. Lumi $= 50 fb^{-1}$
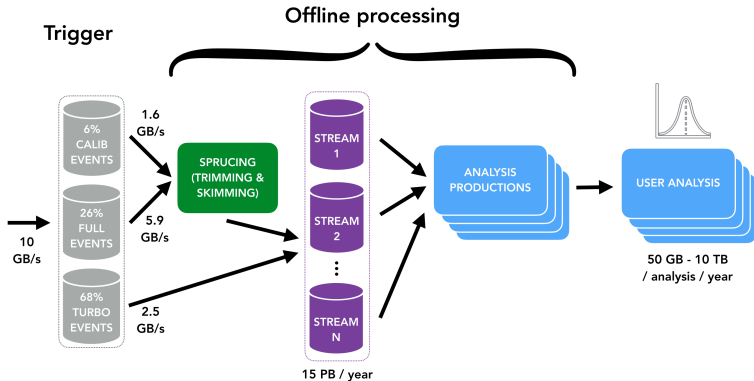
x5 Lumi. wrt Run2

- **Fully-software** *High Level Trigger* (HLT), 30 MHz event reconstruction
  (Daniel Cervenkov's talk "*First performance of the real-time reconstruction at LHCb*")
- Consists of two stages:
  1. GPU-based trigger $\Longrightarrow$ calibrate in real-time (Cristina Agapopoulou's talk "*LHCb HLT1: Tracking and vertexing at 30MHz with GPUs*")
  2. Full reconstruction in CPU-based trigger $\Longrightarrow$ 10 GB/s output (Miroslav Saur' talk "*LHCb HLT2: Real-time alignment, calibration, and software quality-assurance*")



[LHCB-FIGURE-2020-016]

- Increased data rate in Run3 poses significant Offline data processing challenges
- Coordination of these activities by DPA project:
  - software project on same level as detector ones
- DPA consists of 6 work-packages focused on specific tasks



[LHCB-FIGURE-2020-016]

**DPA work-packages**

## WP1: Sprucing

- Offline, central data skimming and slimming
- Sharing of HLT2 framework
- Ensemble of "Sprucing selections" from physics WGs

## WP2: Analysis productions

- Centralised nTupling via DIRAC production system
- Maximal automation
- Inbuilt testing/validation & analysis preservation

## WP3: Offline analysis tool

- Offline analysis application sharing HLT2/Sprucing tools
- Modern design used by Analysis Productions
- Thread safe application

## WP4: Innovative analysis

- R&D for innovative analysis techniques to be adopted in the future
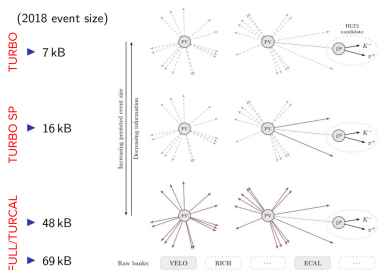- Quantum computing

## WP5: Legacy data & software

- Continued re-stripping of legacy Run1+2 data
- Maintenance of legacy software stacks fro Run1+2 data
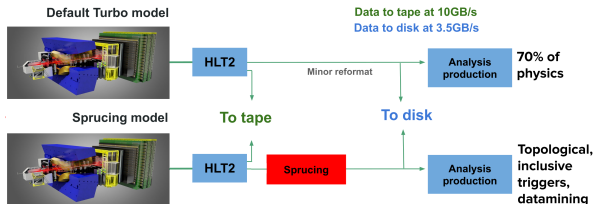
## WP6: Analysis preservation & Open data

- Release LHCb data to CERN Open Data
- Guidelines and tools for analysis preservation
- LHCb use of CERN's CAP and REANA

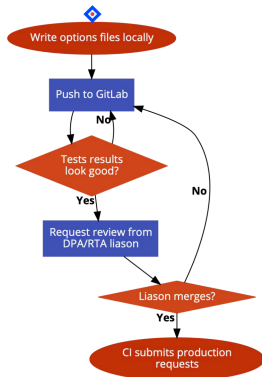● In Run3 event persistency is **customisable** depending on the Physics involved



(2018 event size)

TURBO
▸ 7 kB

TURBO SP
▸ 16 kB

FULL/TURCAL
▸ 48 kB
▸ 69 kB

| Event Size | Persisted objects | Saved to disk |
|---|---|---|
| 4-16 kB | Only the signal candidate is saved, discarding the rest of the event | Yes. No further trimming of data is needed. |
| $\sim$ 16 kB | The signal candidate is saved together with a custom set of other physics objects. | Yes. No further trimming of data is needed. |
| 48-69 kB | The whole event information is retained. | No, saved to tape, not accessible to users. Move to disk only after a sprucing selection. |



**Default Turbo model**

Data to tape at 10GB/s
Data to disk at 3.5GB/s

HLT2 → Minor reformat → Analysis production → **70% of physics**

**Sprucing model**

To tape

To disk

HLT2 → Sprucing → Analysis production → **Topological, inclusive triggers, datamining**

● Sprucing: further offline reduction/selection between tape and disk storage
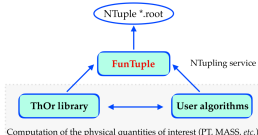  ● same HLT2 selection framework

- After Sprucing/Turbo(SP) stage, data is saved to disk in **Streams** grouped according to the physics of interest

- In Run1+2 analysts created nTuples individually from data on disk using Ganga (a gateway to the Grid):
  - does not scale well for Run3
  - 1000s of faulty jobs can be submitted instantly
  - failed jobs re-submitted manually by user
  - time consuming, O(weeks) for Run1+2 tuples

- Analysis productions submit nTupling jobs centrally using the DIRAC transformation System
- Ensuring a fully automatic:
  - management of files grouping and job failures
  - test of job options
  - preservation of job details in LHCb bookkeeping/EOS
  - error interpretation/advice
  - visualization of the results on a web page



Flowchart for Run3

**WP3: Offline analysis tools**

- In Run1+2 nTuples used "TupleTools" from DaVinci (user analysis) application
  - Pro: easy to implement building blocks("TupleTools") to save variable branches cases
  - Cons: adds lots of redundant branches (often 500+ variables)
    $\implies$ 500GB-10TB of data for only a single Run1+2 analysis!

- The whole DaVinci framework has been **re-optimised** for Run3:
  - *modern thread-safe* (ThOr) functors as in Sprucing used to create light-weight nTuples
  - *consistency* between Online and Offline selections/tools/algorithms
  - analyst has *full control* over which variables for which particles are persisted in nTuple
  - *transparent* configuration in dedicated YAML/JSON and python files
  - unit-testing routines for debugging and CI within the GitLab platform
  - AnaProds will run analysts' DaVinci options



```
"""
Tuple observables relatex to Jpsi -> mu+ mu-
"""
# FunTuple make fields to tuple
fields_jpsi = {
    'Jpsi': "[J/psi(1S) -> mu+ mu-]CC",
    'mup': "[J/psi(1S) -> ^mu+ mu-]CC",
    'mum': "[J/psi(1S) -> mu+ ^mu-]CC"
}
# FunTuple make collection of functors for Jpsi
variables_jpsi = FunctorCollection({
    'THOR_P': F.P,
    'THOR_PT': F.PT,
    'THOR_mom_PT': F.CHILD(1, F.PT),
    'THOR_num_PT': F.CHILD(2, F.PT)
})
```

```
Upgrade_Bd2KstarMuMu_ldst:
    filenames:
        - '/path_to_inputfile_1/file1.dst'
        - '/path_to_inputfile_2/file2.dst'
    qualifiers:
        data_type: Upgrade
        input_type: DST
        simulation: true
```

- Focus on exploitation of new analysis facilities with heterogeneous computing resources (GPU/CPU/FPGA)

- Worldwide LHC Computing Grid consists of $\sim$ 1M CPU cores over 170 sites
  - main LHCb activities based on CPUs
  - supporting High Performance Computing centers providing large GPU resources
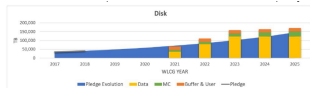  - potential to utilise LHCb's HLT1 GPU farm during detector downtime
- Development to run LHCb payloads on GPUs
  - use advanced algorithms, as Generative Adversarial Networks (GANs), to train models describing LHCb sub-detector $\Longrightarrow$ GPUs speed up GAN training, *Ultra-fast-simulation*
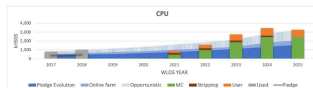  - users using GPUs for complex amplitude analysis models with large statistics

- In Run3 LHCb will produce $\sim$ 15PB of data on disk per year



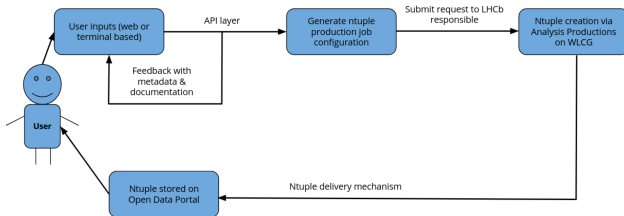- Simulation will require 90% of total offline CPU resources



LHCb-TDR-18

- First investigation into use of Quantum Machine Learning for jet tagging (Davide Zuliani's talk "Quantum Machine Learning for b-jet identification")

**WP6: Analysis preservation and open data**
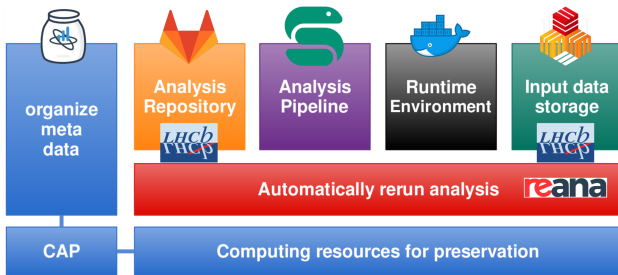
- In accordance with CERN Open data policy, part of the LHCb dataset is made available to general public
- Run3 LHCb data to be released on CERN Open Data portal
- Development of Open Data nTuple Wizard (not yet released and in production):
  - auto-generates options from intuitive user input
  - no prior knowledge of LHCb software is required
  - launches AnaProds & returns nTuple to user on CERN Open Data Portal

  (Ryunosuke O'Neil's poster "An NTuple production service for accessing LHCb Open Data")
- Much smaller storage and bandwidth requirements on Open Data Portal
- Much easier accessing LHCb data



LHCb Analysis Preservation

- Fully compatibility with Snakemake to preserve workflow and provenance tracking
- Capture all dependencies from CERN/LHCb docker containers
- Additional software environments configured from CVMFS
- Deploy analysis to REANA via GitLab to demonstrate preservation status
- Possibility to produce plots and final results documenting their provenance



LHCb Analysis Preservation

- LHCb for Run3 has a *brand new detector & software*
- LHCb will have to process data offline an order of magnitude larger than in Run2
- LHCb is progressing well to meet the Offline demands that Run3 will bring
- Most analyses will be based on turbo candidates
- Analysis workflow will change to cope with the high rates