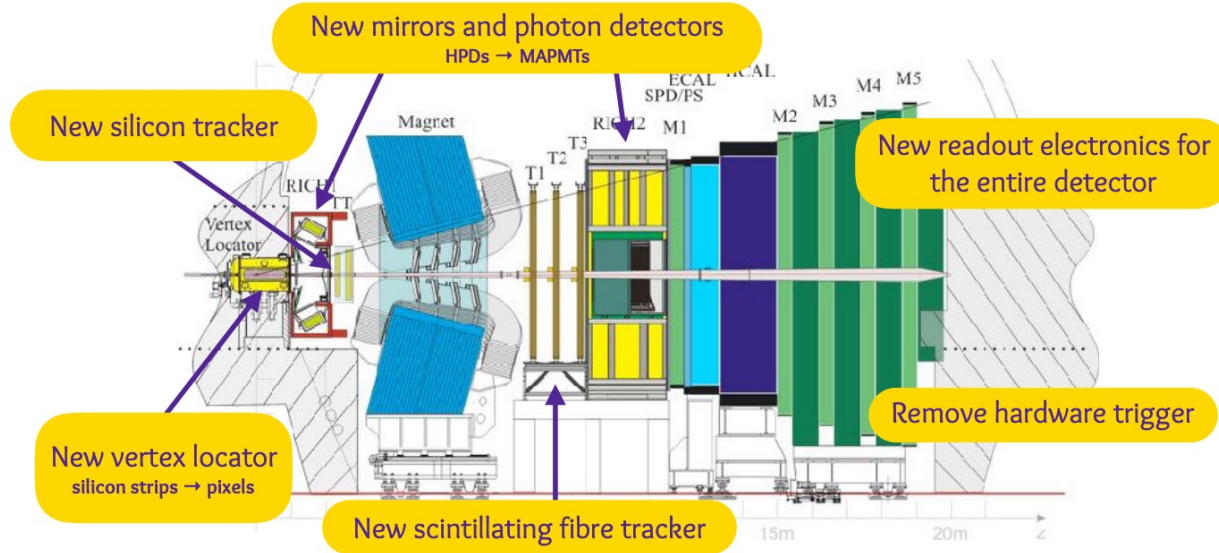# LHCb News

LCG-France Meeting

November 18th 2021

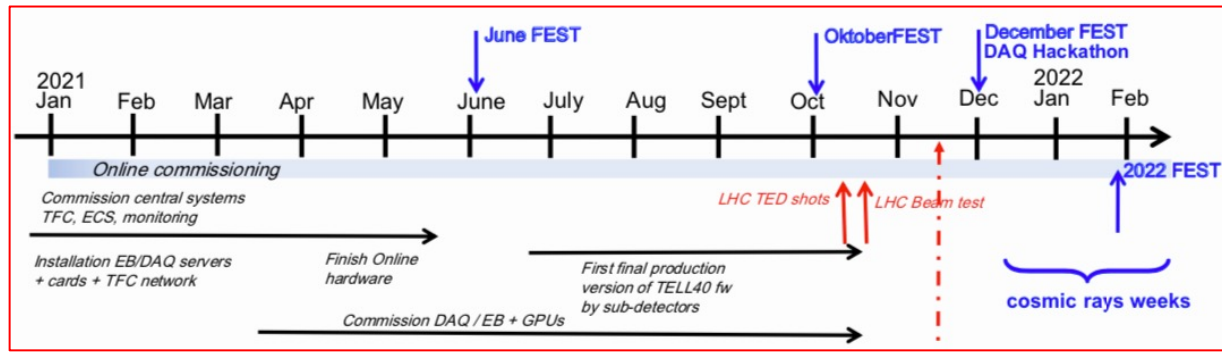Andrei Tsaregorodtsev, CPPM

# Outline

- LHCb Run3 detector upgrade

- Run 3 Computing Model

- Current operations

- Status of requests/pledges 2022-2023 and beyond

- Conclusions

# The upgraded LHCb detector for Run 3



- This is a new detector !
  - Major detector modifications, new tracker
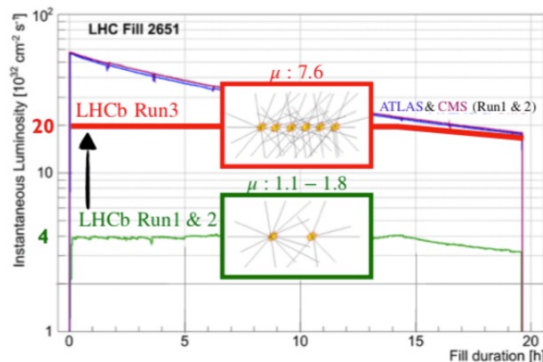  - 100% new RO electronics, DAQ

# Run3 roadmap



**Global activities:**

- Online and sub-detector commissioning during LHC beam test
- Commissioning weeks with cosmic ray tests, integrating newly installed detectors while their stand-alone commissioning progresses
- Full Experiment System Test (FEST)
  - Simulated samples injected in the online system
  - Full dataflow run in commissioning weeks

4

# Run 3 conditions for LHCb

## Luminosity increase: x5

- More interaction vertices per collision of proton bunches, more tracks, more signal
- Beauty and charm signal rates: 1-10MHz
- Almost all events will have a *b* or *c* hadron in Run 3



- Hardware trigger is no more an option
  - No simple local criteria
  - Track reconstruction is needed for event selection
    - Discover event topology as early as possible

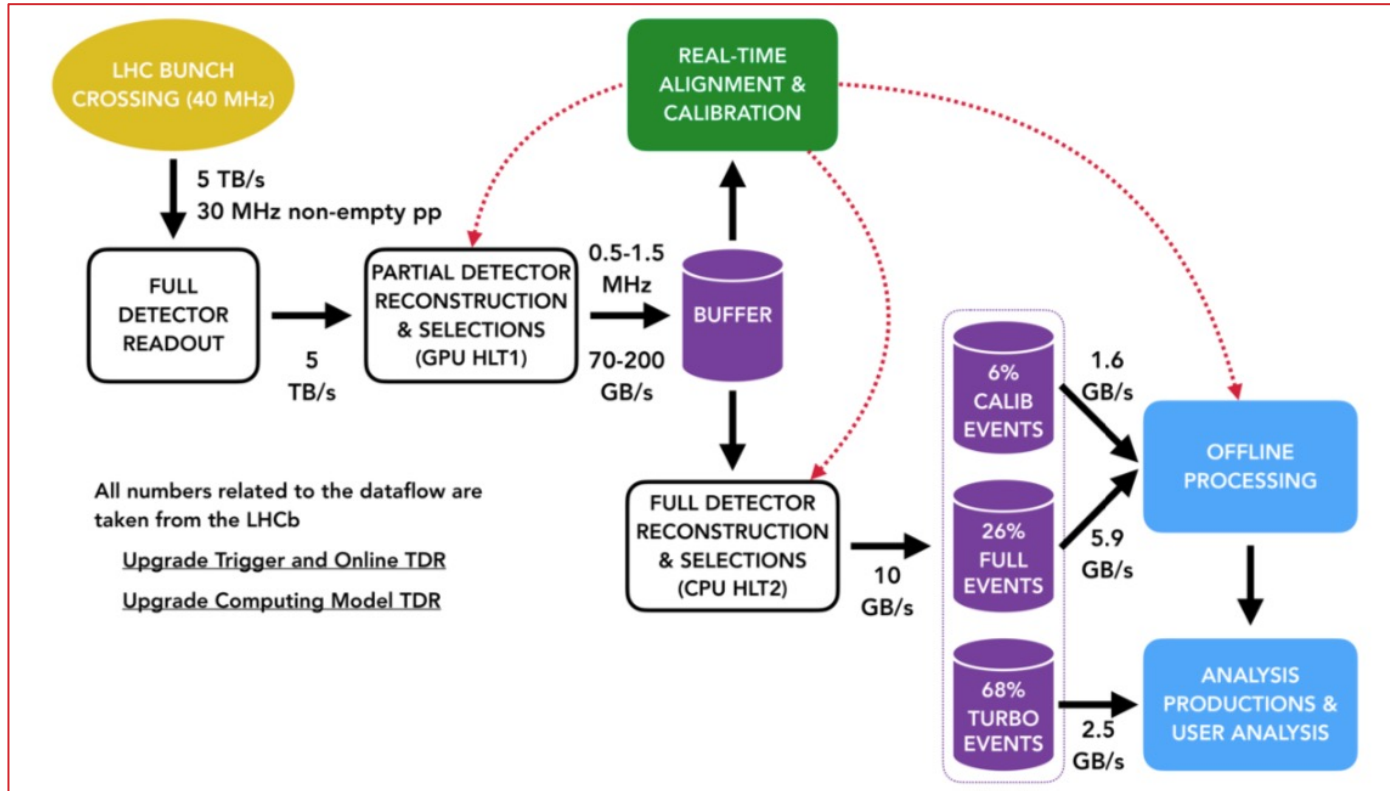- Full software trigger is required

# The MHz signal era



"From a needle in a haystack to an haystack of needles"

# Run3 Computing Model

# LHCb Run 3 Data Flow

# LHCb Run3 HLT practical implementation

200G IB
100GbE
10GbE

173 Event Builder servers

Three TELL40 readout boards per EB server

32 Tb/s

1 Tb/s

1 Tb/s

16 storage servers

40 HLT2 servers   40 HLT2 servers   40 HLT2 servers   40 HLT2 servers

Up to 100 HLT2 sub-farms (4000 servers)

CPU+RAM1   CPU+RAM2

RU   BU   RU   BU

Readout   Readout   200G EB net   Accelerator   10G HLT net   Readout   200G EB net   Accelerator   Accelerator   10G HLT net

GPU-equipped event builder PC, with traffic of all three readout cards.

3 PCIe40 (FPGAs)

2 network connections

1-3 GPUs

PCIe slots

# Persistency model

- Selective persistency: write out only the "interesting" part of the event.



- Turbo stream:
  - Miminum output: only HLT2 signal candidates
  - Optionally: (parts of) pp vertex (e.g. "cone" around candidate for spectroscopy searches)

Limitations: cannot refit tracks and PVs offline, rerun flavour tagging etc.

Advantage: Event size O(10) smaller than RAW

- FULL stream: all reconstructed objects in the event
  - Optionally adding selected RAW banks
- TurCal stream: HLT2 candidates and RAW banks
  - Used for offline calibration and performance measurement

# Output rates

- Moving a larger fraction of
  the physics program to Turbo
  decreases the output bandwidth
  - Turbo events – 16% of Full size events

- Baseline assumes 73% of the
  physicis selections on Turbo
  - Correponds to the output bandwidth of 10GB/s



11

# Data flow evolution Upgrade

Cannot save all HLT output straight to **disk**!
- Utilise cheap **tape storage** for bulk of bandwidth (full stream)
- Rely on central offline slimming/skimming
- Safer option for some physics/allows data mining

**Default model Turbo**



**To tape**

**HLT2**

Minor reformat

**Data analysis**

**To disk**

**Default model - >70% of physics**

**Sprucing model FULL and TURCAL**



**To tape**

**HLT2**

**Sprucing**

**Data analysis**

**To disk**

**Use cases - topological, inclusive triggers, datamining**

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

# Streaming and filtering in Run3

- Can we fit 10 GB/s in a reasonable amount of storage resources ?
  - 10 GB/s to tape
  - Reduce by ~1/6 FULL and Calibration data volume with "sprucing"
    - Selecting events to store
      - $O(10^3)$ selection lines
    - Selecting a subset of reconstructed objects to store

- Save **3.5 GB/s** to disk!

Throughput to tape

| stream | rate fraction | throughput (GB/s) | bandwidth fraction |
|--------|---------------|-------------------|--------------------|
| FULL   | 26%           | 5.9               | 59%                |
| Turbo  | 68%           | 2.5               | 25%                |
| TurCal | 6%            | 1.6               | 16%                |
| total  | 100%          | 10.0              | 100%               |

Throughput to disk

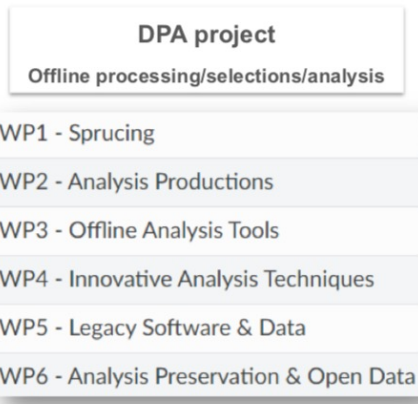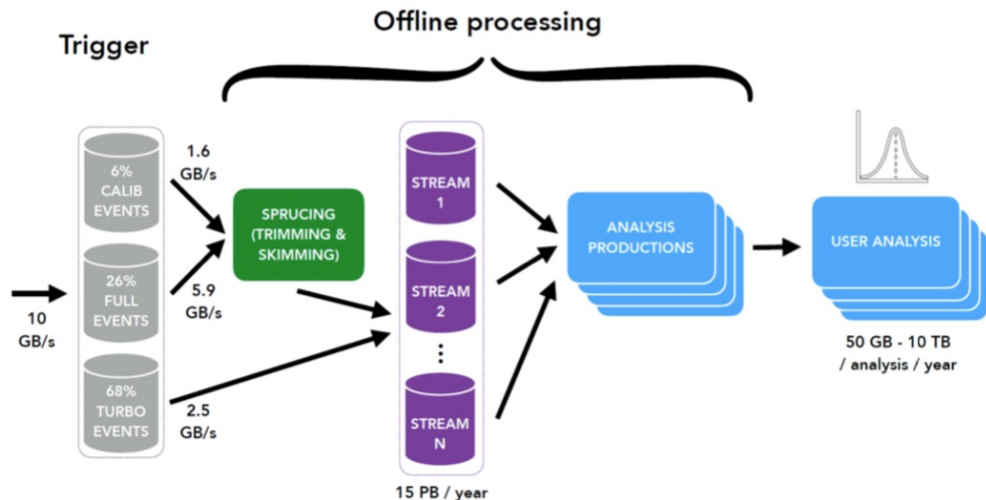| stream | throughput (GB/s) | bandwidth fraction |
|--------|-------------------|--------------------|
| FULL   | 0.8               | 22%                |
| Turbo  | 2.5               | 72%                |
| TurCal | 0.2               | 6%                 |
| total  | 3.5               | 100%               |

# "Data Processing and Analysis" (DPA) project
## An offline workflow for the 2020s

Very large increase in data volume wrt. Run II brings challenges to offline data processing and analysis
DPA built around 2 main ideas:
- Centralised skimming and trimming (aka Sprucing of significant fraction of HLT2 outputs)
- Centralised analysis productions for physics WGs and users

*V.Gligorov*

# The model: what about CPU ?

- CPU is dominated by MC production (~90% of CPU power)

- Expected to be the same at the Upgrade

- Baseline simulation numbers:
  - Event timing:
    - Full/fast/parametric simulation: 120/40/2 seconds
  - Sharing full / fast / parametric: 40/40/20
- Aggressive use of faster simulation techniques:
  - Reduce CPU need
  - No effect on tape
  - No effect on disk
  - May not be feasible, strongly linked to analysis



Running jobs by JobType
11 Weeks from Week 34 of 2019 to Week 45 of 2019

MC production: ~80%

Fast MC production: ~10%

Data processing and analysis ~10%

Max: 134, Min: 3.91, Average: 106, Current: 3.91

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| MCSimulation | 80.1% | MCReconstruction | 3.7% | Merge | 0.1% | MCReprocessing | 0.0% |
| MCFastSimulation | 9.7% | DataStripping | 1.2% | MCMerge | 0.0% | test | 0.0% |
| user | 5.0% | WGProduction | 0.2% | DataReconstruction | 0.0% | unknown | 0.0% |

Generated on 2019-11-17 16:47:39 UTC

# Run 3 Computing model requirements

- Assumptions on simulated event volume
  - N. of  MC events scales with $L_{int}$
  - MC production for a data taking years extends over the following 6 years
  - **MC events saved in MDST format (x40 size reduction!)**
- Assumptions on replicas

| stream | tape | disk |
|---|---|---|
| FULL | $2\times$ RDST + $1\times$ MDST | $3\times$ MDST |
| Turbo | $1\times$ TurboRaw + $1\times$ MDST | $2\times$ MDST |
| TurCal | $2\times$ RDST + $1\times$ MDST | $3\times$ MDST |
| Simulation | $1\times$ MDST | $1\times$ MDST (30% data set only) |

- **All Run 1 + 2 data will be reduced in the end to 1 replica**
- The first year of LHC Run 3 (2021) is considered a "commissioning year" with half the luminosity delivered

16

# WLCG tape challenge

# WLCG tape challenge



target

- Some details [here](here)
- EOS -> T1 write tests
  - Real staging activities in parallel
- DIRAC scaled perfectly
- Met average rate, close to peak rate
  - Issues: FTS settings, number of EOS gridftp gateways, sites configuration
  - the main bottleneck (EOS gridftp gateways) should disappear by the start of Run3
    - Moving to (SRM +) HTTPs
- Not a complete success but
  - Good reminder of the FTS tuning we have to do
  - Highlighted the importance of monitoring
    - efforts required in DIRAC
  - Gave ideas to further optimize the data export from P8

| Site | Expected Speed (GB/s) | Average Speed (GB/s) | Max Speed (GB/s) | Duration (hours) |
|------|------|------|------|------|
| CNAF | 2.24 | 1.07 | 5.16 | 72 |
| IN2P3 | 1.26 | 0.70 | 1.8 | 61 |
| NLT1 | 0.88 | 0.33 | 1.77 | 90 |
| RRC-KI | 0.88 | 0.27 | 1.09 | 112 |
| PIC | 0.58 | 0.24 | 1.15 | 82 |
| RAL | 2.92 | 0.93 | 2.2 | 106 |
| Gridka | 2.24 | 0.35 | 3.41 | 220 |

# WLCG tape challenge: CC/IN2P3



- Immediate start
- Jumps in the throughput
- Target: 1.26 GB/s; average 0.70 GB/s; peak 1.80 GB/s

# Current Operations

# Distributed computing operations

- Computing work dominated by MC production (94%)

- Fast:detailed simulation = 50:50

- Simulating about 180 million events per day

- Incremental stripping of 2018 data recently completed



CPU days used by Site

10 Weeks from Week 35 of 2021 to Week 45 of 2021

LCG.CERN.cern

Online farm 20%

13.5%



Running jobs by JobType
10 Weeks from Week 35 of 2021 to Week 45 of 2021

Fast simulations
Stripping

Max: 167, Min: 5.61, Average: 136, Current: 5.61

| | | | | | | |
|---|---|---|---|---|---|---|
| MCSimulation | 82.5% | DataStripping | 3.2% | Merge | 0.0% | test 0.0% |
| MCFastSimulation | 8.1% | MCReconstruction | 2.8% | unknown | 0.0% | |
| user | 3.3% | WGProduction | 0.1% | MCMerge | 0.0% | |

Generated on 2021-11-12 16:39:28 UTC



Evolution of Mean Number of Simulated Events per Day

Mean 7 days
Mean 30 days
Mean 90 days

1.8E+08
1.4E+08

# Opportunistic resources

- HLT farm 20%
- Non-pledging sites 10%
- HPCs
  - NERSC, CSCS, SDumont now in production
  - Barcelona Supercomputing Center (BSC), still not in production
    - Installation and configuration of ARC CE
  - CINECA/Marconi100
    - GPU + Power9: difficult to use in normal production workflows, no full software build
    - Some user jobs run locally, very limited CPU consumption
    - DIRAC configured for grid-like access, pilots sent but no matching jobs yet
  - O(1000) computing slots in total
    - Not a lot !



Running jobs by Grid
15 Weeks from Week 22 of 2021 to Week 37 of 2021

WLCG

opportunistic

Max: 159, Min: 85.7, Average: 128, Current: 124

| LCG | 68.0% | VAC | 0.5% | ANY | 0.0% | Group | 0.0% |
| DIRAC | 31.4% | Multiple | 0.1% | CLOUD | 0.0% | | |

Generated on 2021-09-17 12:05:11 UTC



Cumulative Pilots by Site
30 Days from 2021-10-13 to 2021-11-12

Max: 539, Min: 0.67, Average: 85.5, Current: 539

DIRAC.CINECAM100.it    538.9

Generated on 2021-11-12 18:08:00 UTC

# French contributions in 2021



CPU days by Country 2021
45 Weeks from Week 00 of 2021 to Week 46 of 2021

| | |
|---|---|
| LHCB | 9648044.3 |
| UK | 6964373.5 |
| CERN | 4779244.7 |
| FR | 3448891.6 |
| IT | 3372455.2 |
| DE | 3299884.2 |
| CH | 1943491.3 |
| PL | 1745838.0 |
| RU | 1355719.2 |
| NL | 1157398.6 |
| US | 682240.1 |
| ES | 452755.1 |
| RO | 422326.3 |
| BR | 374235.0 |
| CN | 301485.3 |
| AU | 60754.8 |
| MULTIPLE | 13417.6 |
| IL | 12142.5 |
| ANY | 2392.4 |
| DIRAC.CLIENT.LOCAL | 7.8 |
| ZONE | 1.6 |

France: 8.5%

Generated on 2021-11-17 21:07:41 UTC



CPU days at French Sites 2021
45 Weeks from Week 00 of 2021 to Week 46 of 2021

| | |
|---|---|
| LCG.IN2P3.fr | 1254343.0 |
| LCG.LAL.fr | 550012.4 |
| LCG.CPPM.fr | 480604.3 |
| LCG.LAPP.fr | 313000.1 |
| LCG.LPNHE.fr | 245418.8 |
| LCG.LPC.fr | 157713.1 |
| LCG.AUVER.fr | 8774.5 |

Generated on 2021-11-17 21:16:22 UTC
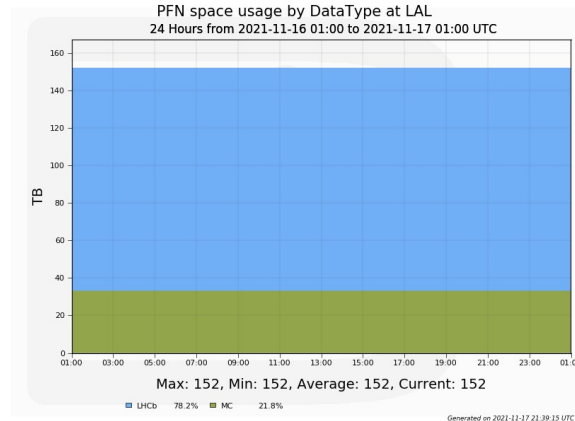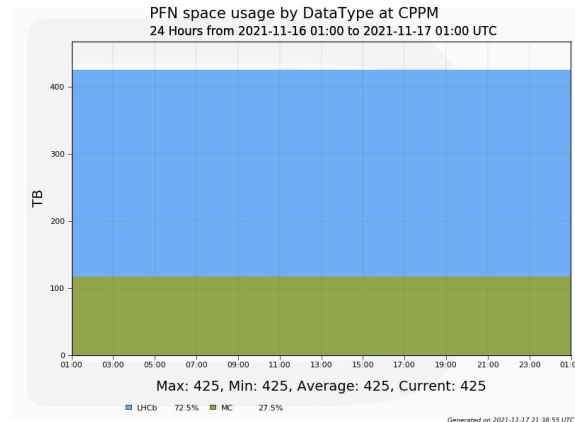
- No particular comments with respect to the French sites functioning
  - Some occasional problems with running pilots at Condor/CC
    – solved by Vanessa

# Disk space usage at T2D's 2021

- CPPM
  - Pledged 600TBs, used 425 TBs
  - Occupancy 71%
- LAL
  - Pledged 383TBs, used 152 TBs
  - Occupancy 40%
- LHCb T2D policy
  - T2D introduced to allow countries without T1's to contribute storage resources
  - No special use of storages at T2's compared to T1 storage - what matters is T1+T2 disk storage
  - But more attention to SEs at T2 sites due to less operational overheads
    - Single person responsible for data management



PFN space usage by DataType at CPPM
24 Hours from 2021-11-16 01:00 to 2021-11-17 01:00 UTC

Max: 425, Min: 425, Average: 425, Current: 425

LHCb 72.5%   MC   27.5%

*Generated on 2021-11-17 21:38:55 UTC*



PFN space usage by DataType at LAL
24 Hours from 2021-11-16 01:00 to 2021-11-17 01:00 UTC

Max: 152, Min: 152, Average: 152, Current: 152

LHCb 78.2%   MC   21.8%

*Generated on 2021-11-17 21:39:15 UTC*

Requests and pledges

# 2022 pledges situation

| Tier | Pledge Type | Year | LHCb Required | LHCb Pledged | LHCb Balance |
|------|-------------|------|---------------|--------------|--------------|
| 0 | Tape | 2022 | 81000 | 81000 | 0 % |
| 0 | Disk | 2022 | 26500 | 26500 | 0 % |
| 0 | CPU | 2022 | 189000 | 189000 | 0 % |
| 1 | Tape | 2022 | 139000 | 116337 | -16 % |
| 1 | Disk | 2022 | 52900 | 47783 | -10 % |
| 1 | CPU | 2022 | 622000 | 514531 | -17 % |
| 2 | CPU | 2022 | 345000 | 332640 | -4 % |
| 2 | Disk | 2022 | 10200 | 6941 | -32 % |
| Tier | Pledge Type | Year | LHCb Required | LHCb Pledged | LHCb Balance |

~10% lower pledges at Tier1s – significantly less disk at Tier2s
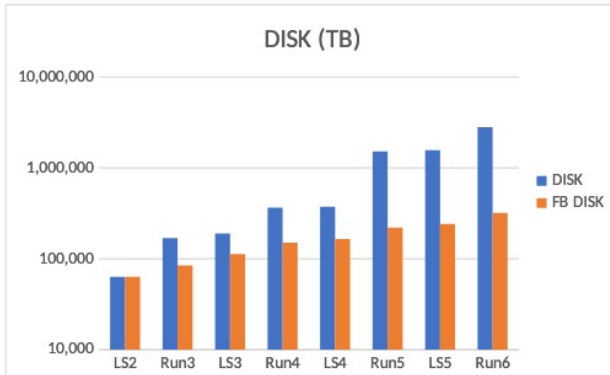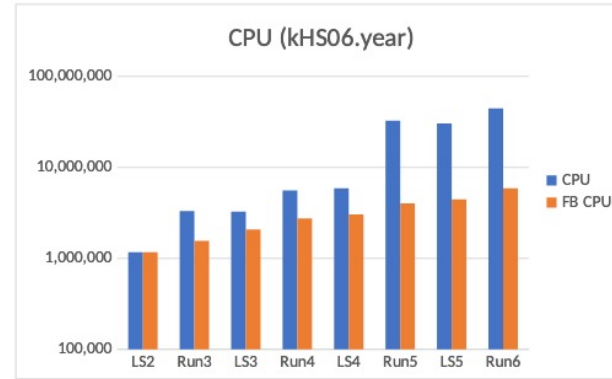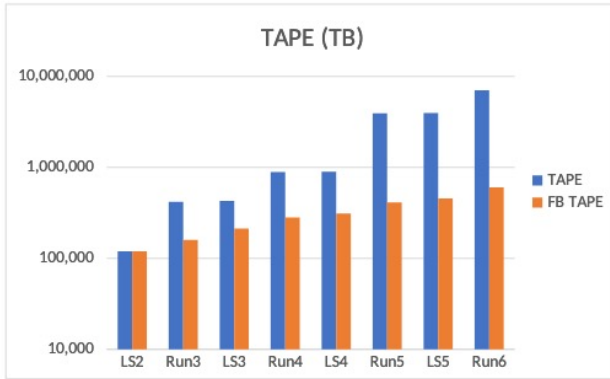Reality check needed vs. e.g. LHC planning and LHCb readiness

# 2023 preliminary requests shown at the C-RRB

| LHCb | | LHCb-PUB-2021-002 | | THIS DOCUMENT | |
| --- | --- | --- | --- | --- | --- |
| | | **2022** | | **2023 (prelim.)** | |
| | | Request | 2022 req. / 2021 CRSG | Request | 2023 req. / 2022 CRSG |
| **WLCG CPU** | **Tier-0** | 189 | 108% | 361 | 190% |
| | **Tier-1** | 622 | 108% | 1185 | 191% |
| | **Tier-2** | 345 | 107% | 657 | 190% |
| | **HLT** | 50 | 100% | 50 | 100% |
| | **Sum** | 1206 | 108% | 2252 | 187% |
| **Others** | | 50 | 100% | 50 | 100% |
| **Total** | | **1,256** | **107%** | **2,302** | **183%** |
| **Disk** | **Tier-0** | 26.5 | 141% | 42.8 | 162% |
| | **Tier-1** | 52.9 | 141% | 85.6 | 162% |
| | **Tier-2** | 10.2 | 141% | 16.5 | 162% |
| | **Total** | **89.6** | **141%** | **144.9** | **162%** |
| **Tape** | **Tier-0** | 81 | 184% | 132 | 164% |
| | **Tier-1** | 139 | 184% | 228 | 164% |
| | **Total** | **219.9** | **184%** | **360.5** | **164%** |

# Upgrade I and II computing model assumptions

| Model assumptions | | |
|---|---|---|
| | Upgrade I | Upgrade II |
| Peak L $(\mathrm{cm}^{-2}\mathrm{s}^{-1})$ | $2 \times 10^{33}$ | $1.5 \times 10^{34}$ |
| Yearly integrated luminosity ($\mathrm{fb}^{-1}$) | 10 | 50 |
| Logical bandwidth to tape (GB/s) | 10 | 50 |
| Logical bandwidth to disk (GB/s) | 3.5 | 17.5 |
| Running time (s) | $5 \times 10^6$ | |
| Trigger rate fraction (%) | 26 / 68 / 6 Full / Turbo / TurCal | |
| Ratio Turbo/Full event size | 16.7% | |
| Ratio full/fast/param. MC | 40:40:20 | |
| CPU work per event full/fast/param. MC (HS06.s) | 1200 / 400 / 20 | |
| Number of simulated events | $4.8 \times 10^9 / \mathrm{fb}^{-1} /$year | |
| Data replicas on tape | 2 (1 for derived data) | |
| Data replicas on disk | 2 (Turbo); 3 (Full, TurCal) | |
| MC replicas on tape | 1 (MDST) | |
| MC replicas on disk | 0.3 (MDST, 30% of the total dataset) | |

28

# Resources required for Run 4,5,6



TAPE (TB)



CPU (kHS06.year)



DISK (TB)

New resources can not be acquired in a scheme where funding is flat and performance increase by 10% each year.

# Mitigation strategies

- Similar to ATLAS and CMS, huge R&D effort of the HEP community
  - Simulation
    - GEANT4 running on GPU
    - Calorimeter cluster simulation using ML techniques and/or shower libraries
    - …
  - Reduce storage requirements
    - nanoAOD format
    - Lossless data compression
    - Improves data placement
    - …
  - Skilled manpower is the key for the success !

# Conclusions

- Run3 is a huge challenge for the new LHCb detector, trigger, DAQ and offline processing

- Smooth ongoing offline computing operations dominated by the MC production

- Pledges for the coming years are below the LHCb requests

- Ongoing effort to optimise the MC software, data production procedures, onvolve new opportunistic resources including HPCs