# Distributed Data Management in LHCb

**Ph. Charpentier**

**CERN-LHCb**

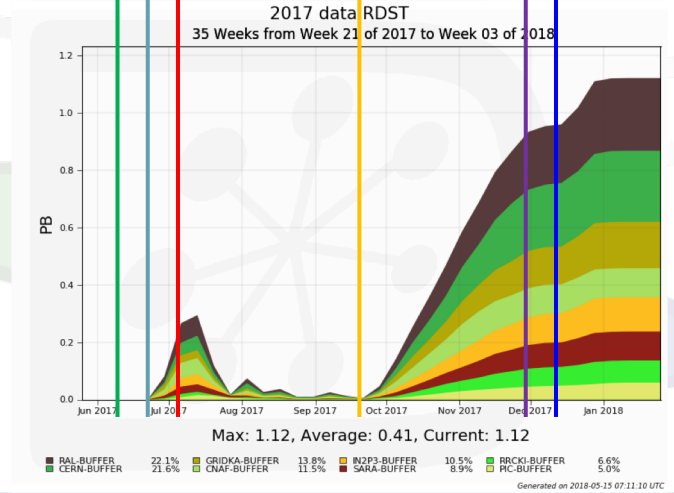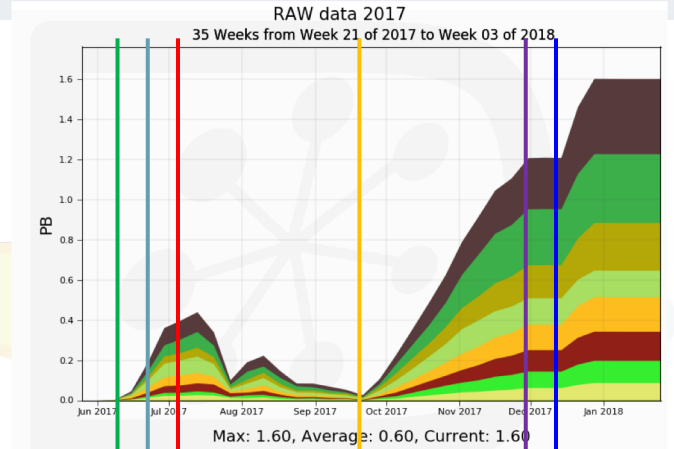# Data Management model: physics analysis

- Physics data: real data ($\mu$)DST
  - Merge into ~5GB files: keep run granularity
  - Each run has its destination site: upload to disk at destination
  - Replicate to a second disk SE + one "archive" tape SE
- MC physics data ($\mu$)DST
  - Small files uploaded with geographical mapping to disk @ Tier1
  - Files merged (5GB) at Tier1
  - Replicate to a second disk SE + one "archive" tape SE
- Selecting second disk SE
  - Randomly selected using free disk space at site as a weight
  - Advantage: no risk to fill up a site with replication
- For real data: keep only last processing pass
  - Sometimes last 2 processing passes
  - Additionally some incremental processing passes (fixing and adding selections)
- Use popularity for retiring datasets (real data and MC)

- **All production data are centrally managed**
  - **No user / site / group is allowed to replicate production datasets**
  - **Users are handling their own private data**
  - **Working groups also have dedicated storage, but WG productions are regular productions, i.e. centrally managed**
- **Duties of the production manager**
  - **RAW data replication and removal from disk after processing**
  - **RDST and FULL.DST (reco output) replication to tape / disk**
  - **Stripping output replication**
  - **MC output replication**
  - **Obsolete or bad data removal (complete removal or only disk removal)**
  - **RAW and RDST pre-staging prior to re-stripping campaigns**
    - ❖ **In conjunction with removal after processing**
- **Mostly handled by one person**
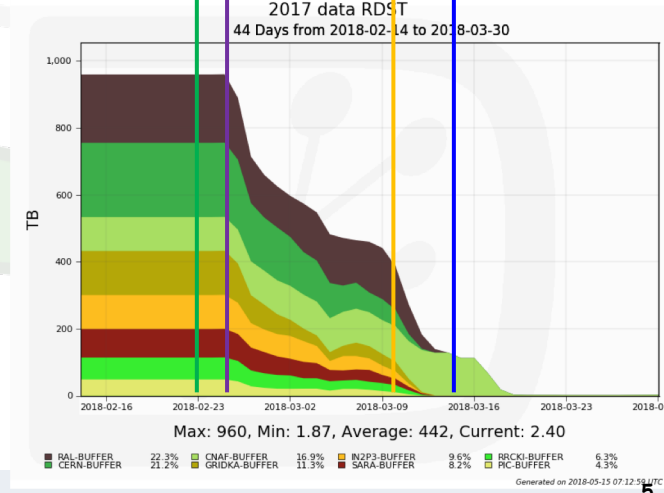  - **Was me for long, now Chris**

- **Start of data taking**
  - **Replicating RAW data to buffer**
- **Start reconstruction**
  - **RDST to Buffer**
- **Start stripping**
  - **Removing RAW and RDST from Buffer**

- **Stop removing RAW+RDST from buffer after stripping**

- **End of data taking**

- **Stage first part of the year (but CNAF)**



RAW data 2017
35 Weeks from Week 21 of 2017 to Week 03 of 2018

Max: 1.60, Average: 0.60, Current: 1.60

2017 data RDST
35 Weeks from Week 21 of 2017 to Week 03 of 2018

Max: 1.12, Average: 0.41, Current: 1.12

| | | | | | |
|---|---|---|---|---|---|
| ■ RAL-BUFFER | 22.1% | ■ GRIDKA-BUFFER | 13.8% | ■ IN2P3-BUFFER | 10.5% | ■ RRCKI-BUFFER | 6.6% |
| ■ CERN-BUFFER | 21.6% | ■ CNAF-BUFFER | 11.5% | ■ SARA-BUFFER | 8.9% | ■ PIC-BUFFER | 5.0% |

*Generated on 2018-05-15 07:11:10 UTC*

# Real data re-stripping

- **Start re-stripping**
- **Start removal of RAW+RDST**
- **Start staging at CNAF (after flood)**
- **Start CNAF processing**



2017 physics data (re-stripping)
39 Days from 2018-02-19 to 2018-03-30

Max: 1.84, Average: 1.08, Current: 1.84

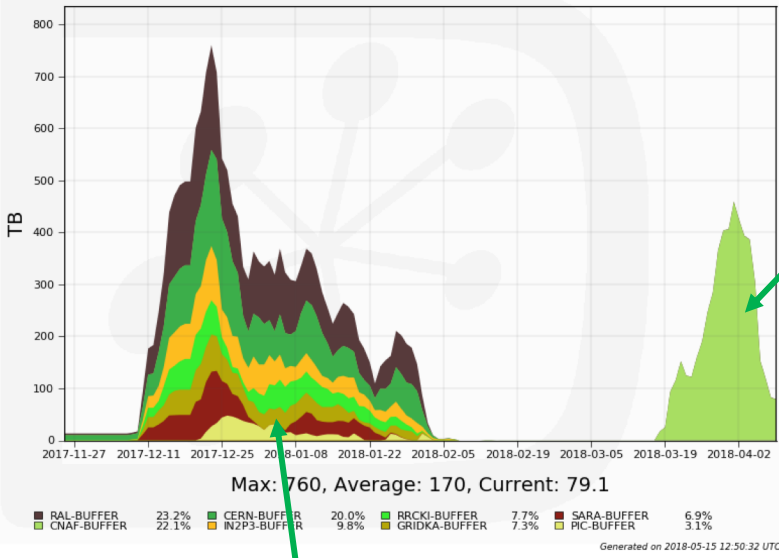| | | | | | | |
|---|---|---|---|---|---|---|
| CERN-DST-EOS | 14.1% | CNAF-DST | 3.2% | Manchester-DST | 0.7% |
| IN2P3-ARCHIVE | 10.6% | RRCKI-DST | 3.2% | RAL-HEP-DST | 0.7% |
| RAL-ARCHIVE | 10.5% | UKI-LT2-IC-HEP-DST | 2.3% | LAL-DST | 0.5% |
| CERN-ARCHIVE | 10.4% | NCBJ-DST | 1.8% | LPNHE-DST | 0.3% |
| RAL-DST | 10.3% | CNAF-ARCHIVE | 1.8% | CBPF-DST | 0.3% |
| GRIDKA-DST | 9.2% | CPPM-DST | 1.5% | Glasgow-DST | 0.3% |
| IN2P3-DST | 6.3% | CSCS-DST | 1.3% | IHEP-DST | 0.1% |
| SARA-DST | 5.0% | UKI-LT2-QMUL-DST | 1.3% | RRCKI-FAILOVER | 0.0% |
| PIC-DST | 3.8% | Liverpool-DST | 0.9% | ... plus 8 more | |

Generated on 2018-05-15 13:47:19 UTC

- **Physics data replication**



RAW data 2017
44 Days from 2018-02-14 to 2018-03-30

Max: 1,208, Min: 7.70, Average: 553, Current: 7.70

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| RAL-BUFFER | 21.9% | CNAF-BUFFER | 17.7% | IN2P3-BUFFER | 9.5% | RRCKI-BUFFER | 6.2% |
| CERN-BUFFER | 21.1% | GRIDKA-BUFFER | 11.3% | SARA-BUFFER | 8.0% | PIC-BUFFER | 4.3% |



2017 data RDST
44 Days from 2018-02-14 to 2018-03-30

Max: 960, Min: 1.87, Average: 442, Current: 2.40

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| RAL-BUFFER | 22.3% | CNAF-BUFFER | 16.9% | IN2P3-BUFFER | 9.6% | RRCKI-BUFFER | 6.3% |
| CERN-BUFFER | 21.2% | GRIDKA-BUFFER | 11.3% | SARA-BUFFER | 8.2% | PIC-BUFFER | 4.3% |

Generated on 2018-05-15 07:12:50 UTC

2016 data RAW+RDST
19 Weeks from Week 47 of 2017 to Week 14 of 2018

**Files at CNAF
with mesh processing**

**No files from CNAF**

2016 re-stripping jobs
19 Weeks from Week 47 of 2017 to Week 14 of 2018

# Some technical details: Transformation plugins

○ **All DM (and WMS) operations are handled through the TransformationSystem**
  □ **Specific set of LHCb plugins**

○ **Plugins used to create replication tasks**
  □ **RAWReplication**
    ❖ Replicates RAW files to a Tier1 tape SE and to a disk buffer either at that Tier1 or at CERN
    ❖ A run (1 hour of data taking) is assigned to a site
  □ **ReplicateWithAncestors**
    ❖ Replicates files from tape to disk buffer as well as their ancestors (useful for staging RAW+RDST data prior to re-stripping)
  □ **LHCbDSTBroadcast**
    ❖ Replicates files (grouped by runs) to a given number of disk SEs (current default is one more) and a number of archive SE (current default is 1)
  □ **LHCbMCDSTBroadcastRandom**
    ❖ Replicates files randomly to a number of additional disk SEs (default: 1) and a number of archive SEs (default: 1)
  □ **ReplicateDataset[ToRunDestination] / ArchiveDataset**
    ❖ Creates additional replicas on specific SEs (disk or archive)

- **Plugins used to create removal tasks**
  - **RemoveReplicasWhenProcessed**
    - Creates tasks to remove files from buffer after they have been processed
  - **RemoveReplicasWithAncestors**
    - Creates tasks to remove files from buffer after they have been processed and remove their ancestors as well (useful after stripping)
  - **RemoveDatasetFromDisk**
    - Creates tasks to remove all disk replicas, provided there is a tape replica (useful to remove unpopular / useless datasets)
  - **DestroyDataset**
    - Removes all replicas of all files

- **DM tasks are transformed into requests in the RMS**
  - **Limit the number of files per task / request**
  - **DM operations executed by the RequestExecutingAgent (+FTSAgent for replication)**
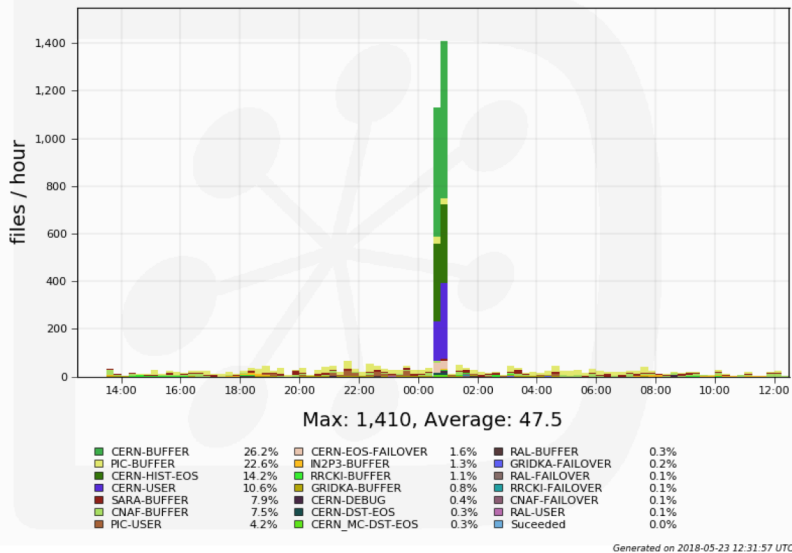
○ **Many plugin parameters are configurable**

- **Decreasing order of precedence:**
  - ❖ **Default in code (usually None)**
  - ❖ **Global defaults in CS (`/Operations/TransformationPlugins` section)**
  - ❖ **Plugin-level defaults (`/Operations/TransformationPlugins/<plugin>` section)**
  - ❖ **Transformation level values (additional transformation parameters)**
- **Examples:**
  - ❖ **`DestinationSEs` / `NumberOfReplicas` for replication plugins / transformations**
  - ❖ **`FromSEs` for removal plugins / transformations**
- **Pre-staging throttling**
  - ❖ **Limit the total number of files pre-staged for a given directory**
  - ❖ **Uses the current occupancy (from StorageUsage DB) and number of files in the pipeline**
  - ❖ **Limit per site computed using site shares**
  - ❖ **Watermark limit on each SE (typically 30-40 TB free space)**
  - ❖ **Allows to launch pre-staging with many files and let the system throttle**
    - ➢ **New files are replicated when old files are removed (timescale of days)**

# Feeding the TransformationSystem

- LHCb uses the Bookkeeping system (BK) to add files to the TS
  - Usually expressed as a BK "path"
    - Can also be an explicit production number
  - Example: `/LHCb/Collision16//RealData/Reco16/Stripping28r1//BHADRON.MDST`
  - Possibly add data quality criteria
    - E.g. avoid data flagged as BAD

- BK query associated to a given transformation
  - Incremental periodic queries, add files to transformation
  - Daily full query (for safety / recovery of failures)

- Files can also be added "manually" to transformations
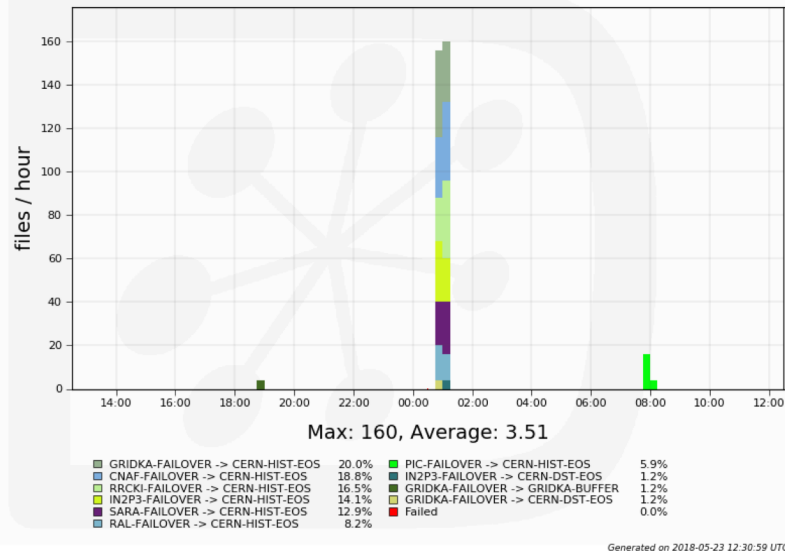  - Good practice to not mix with BK query!

- **If destination SE is banned, or if transfer fails**
  - Upload the file to any of our "failover SEs" (one per Tier1)
  - Create an RMS request to transfer back the file to its final destination
- **(Almost) no loss of output data**
  - Unless there is a connectivity problem on the WN site
  - If this happens the job is Failed
- **Other failover**
  - If registration fails, a request is set (possibly for specific catalogs)
- **If there is a request, the jobs remains in an intermediate status**
  - Completed / Pending requests
  - This is a very bad name and should be changed to ??? (Completing?)
  - The RMS callback will eventually set the job to Done when the request is successful
- **FTS is used for failover transfers**

Failed DM uploads by destination
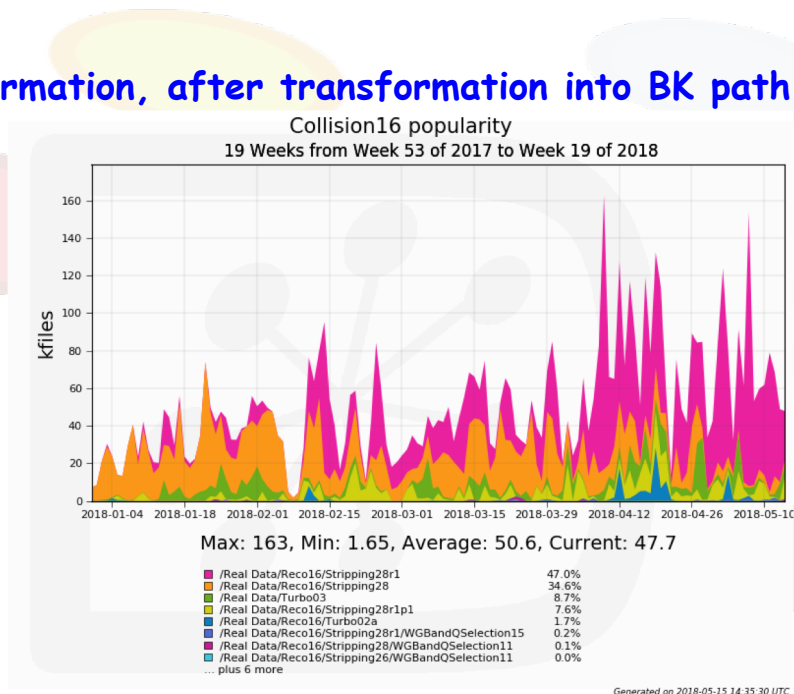24 Hours from 2018-05-22 12:30 to 2018-05-23 12:30 UTC

Max: 1,410, Average: 47.5

| | | | | | |
|---|---|---|---|---|---|
| CERN-BUFFER | 26.2% | CERN-EOS-FAILOVER | 1.6% | RAL-BUFFER | 0.3% |
| PIC-BUFFER | 22.6% | IN2P3-BUFFER | 1.3% | GRIDKA-FAILOVER | 0.2% |
| CERN-HIST-EOS | 14.2% | RRCKI-BUFFER | 1.1% | RAL-FAILOVER | 0.1% |
| CERN-USER | 10.6% | GRIDKA-BUFFER | 0.8% | RRCKI-FAILOVER | 0.1% |
| SARA-BUFFER | 7.9% | CERN-DEBUG | 0.4% | CNAF-FAILOVER | 0.1% |
| CNAF-BUFFER | 7.5% | CERN-DST-EOS | 0.3% | RAL-USER | 0.1% |
| PIC-USER | 4.2% | CERN_MC-DST-EOS | 0.3% | Suceeded | 0.0% |

Generated on 2018-05-23 12:31:57 UTC



Successful Failover transfers
24 Hours from 2018-05-22 12:30 to 2018-05-23 12:30 UTC

Max: 160, Average: 3.51

| | | | |
|---|---|---|---|
| GRIDKA-FAILOVER -> CERN-HIST-EOS | 20.0% | PIC-FAILOVER -> CERN-HIST-EOS | 5.9% |
| CNAF-FAILOVER -> CERN-HIST-EOS | 18.8% | IN2P3-FAILOVER -> CERN-DST-EOS | 1.2% |
| RRCKI-FAILOVER -> CERN-HIST-EOS | 16.5% | GRIDKA-FAILOVER -> GRIDKA-BUFFER | 1.2% |
| IN2P3-FAILOVER -> CERN-HIST-EOS | 14.1% | GRIDKA-FAILOVER -> CERN-DST-EOS | 1.2% |
| SARA-FAILOVER -> CERN-HIST-EOS | 12.9% | Failed | 0.0% |
| RAL-FAILOVER -> CERN-HIST-EOS | 8.2% | | |

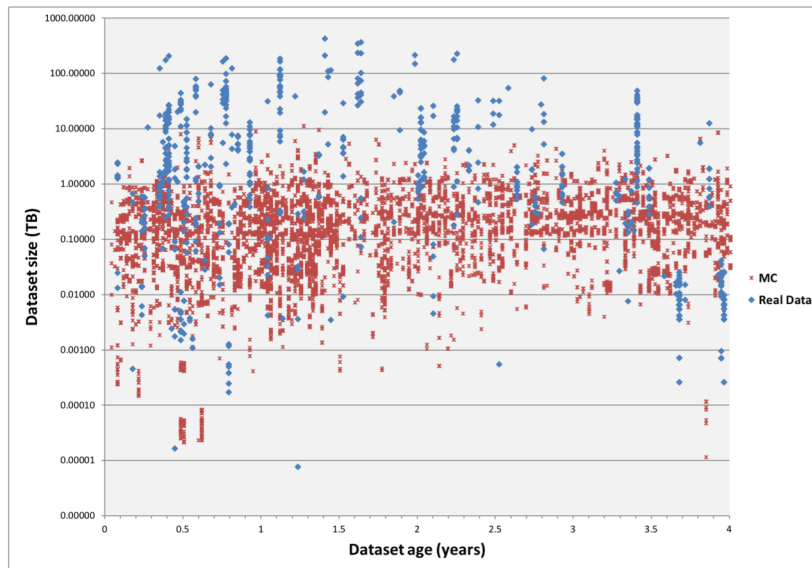Generated on 2018-05-23 12:30:59 UTC

- **LHCb specific extension**
- **StorageUsageAgent: twice a day DFC scan to record per directory and per SE**
  - **Number of files**
  - **Size**
- **StorageHistoryAgent: uses output of the above but translates DFC directory into BK path**
  - **Records each entry into the accounting DB**
  - **Allows to make plots shown previously in this talk**

- Each user job reports, per directory:
  - Number of files processed
  - SE used
- PopularityAgent:
  - Records in accounting DB this information, after transformation into BK path (similar to StorageHistoryAgent)
- PopularityAnalysisAgent
  - Creates huge CSV table
    - Used for manual DM checks and operations
  - May give hints for data removal
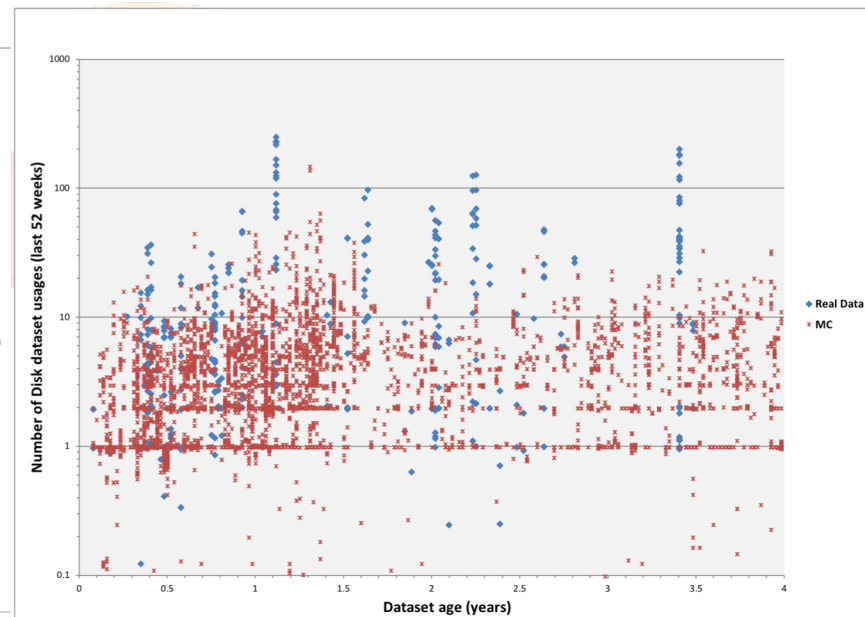    - Based on big data algorithms
    - Developed by our Yandex colleagues



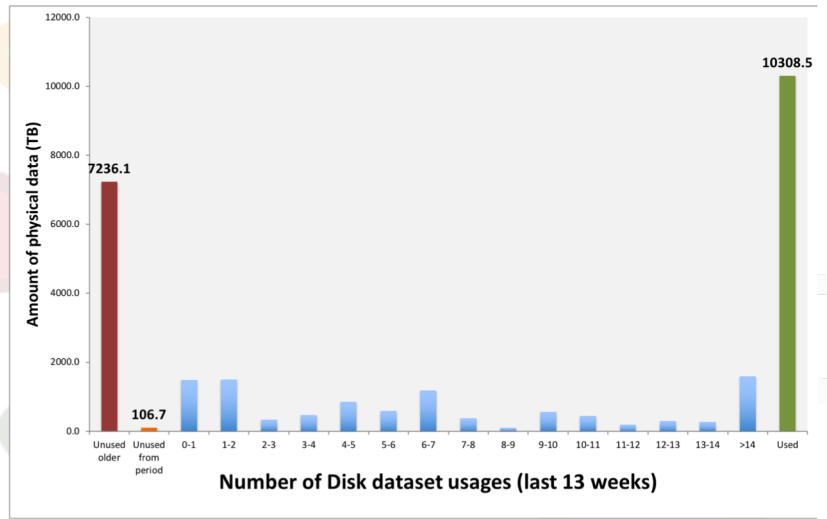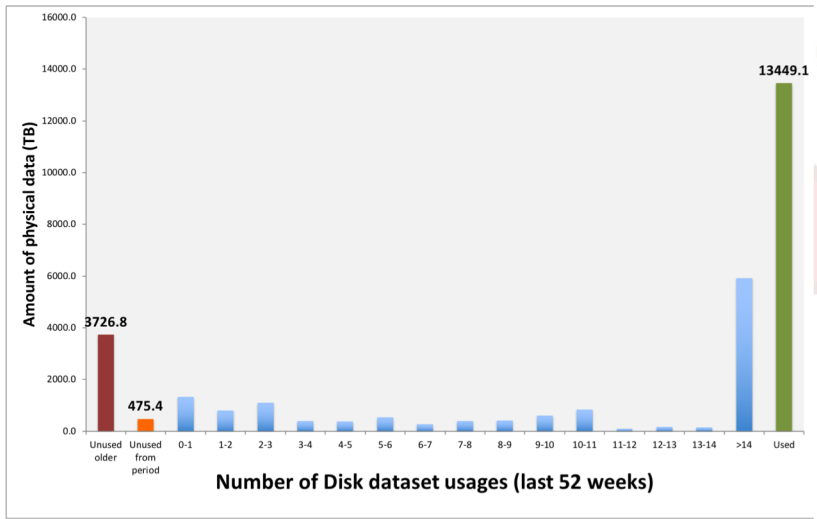Collision16 popularity
19 Weeks from Week 53 of 2017 to Week 19 of 2018

Max: 163, Min: 1.65, Average: 50.6, Current: 47.7

| | | |
|---|---|---|
| /Real Data/Reco16/Stripping28r1 | | 47.0% |
| /Real Data/Reco16/Stripping28 | | 34.6% |
| /Real Data/Turbo03 | | 8.7% |
| /Real Data/Reco16/Stripping28r1p1 | | 7.6% |
| /Real Data/Reco16/Turbo02a | | 1.7% |
| /Real Data/Reco16/Stripping28r1/WGBandQSelection15 | | 0.2% |
| /Real Data/Reco16/Stripping28/WGBandQSelection11 | | 0.1% |
| /Real Data/Reco16/Stripping26/WGBandQSelection11 | | 0.0% |
| ... plus 6 more | | |

*Generated on 2018-05-15 14:35:30 UTC*

## Dataset size vs age

## Dataset usages vs age

- **Distributed Data management in LHCb is fully handled by LHCbDirac**

- **Heavy usage of the TransformationSystem + RMS + FTS**
  - **Input queries from the LHCb bookkeeping**

- **LHCb-specific plugins developed for each purpose**
  - **However quite flexible (many configurable parameters)**

- **Re-processing campaigns quite successful using data from tape**
  - **Pre-staging launched a few weeks before the production**
  - **Automatic removal of pre-staged data after processing**