# The LHCb Data Management System

**Philippe Charpentier**

**CERN**

**On behalf of the LHCb Collaboration**

# LHCb Computing Model in a nutshell
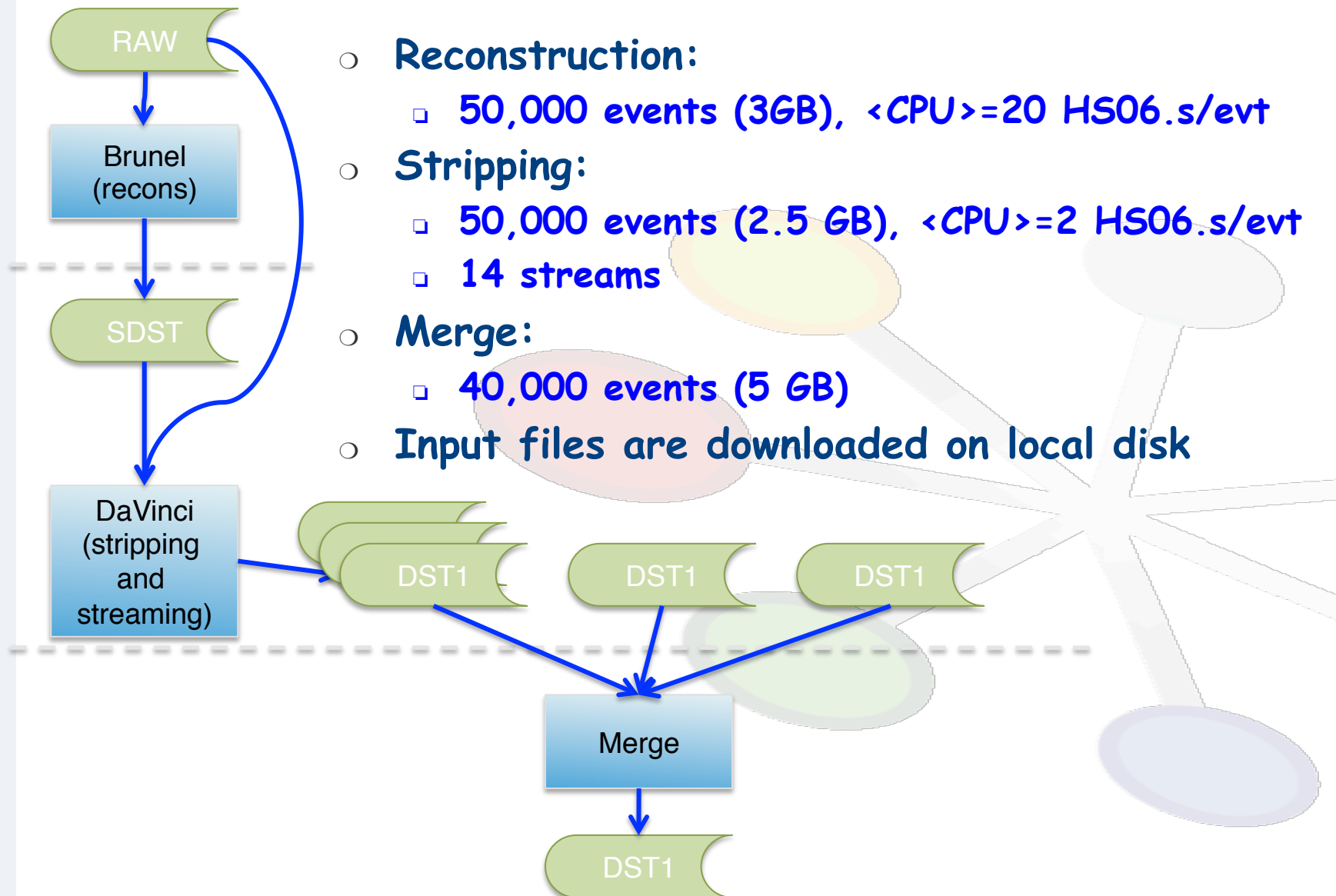
- ○ **RAW files (3GB) are transferred to Tier0**
  - ❑ **Verify migration and checksum**
- ○ **Transfer to Tier1s**
  - ❑ **Each file replicated once (whole run at a single site)**
- ○ **Reconstruction**
  - ❑ **At CERN and Tier1s (up to 300 HS06.hours)**
    - ✩ **If needed Tier2s can also be used as "Tier1 co-processor"**
- ○ **Stripping**
  - ❑ **On Tier1 where SDST is available**
- ○ **MC simulation**
  - ❑ **Complex workflow: simulation, digitization, reconstruction, trigger, filtering**
  - ❑ **Running everywhere with low priority**
    - ✩ **Tier2, Tier1, CERN and unpledged resources (some non-Grid)**
- ○ **User analysis**
  - ❑ **Running at Tier1s for data access, anywhere for MC studies**
- ○ **Grid activities under control of LHCbDirac**
  - ❑ **LHCb extensions of the DIRAC framework (cf Poster**
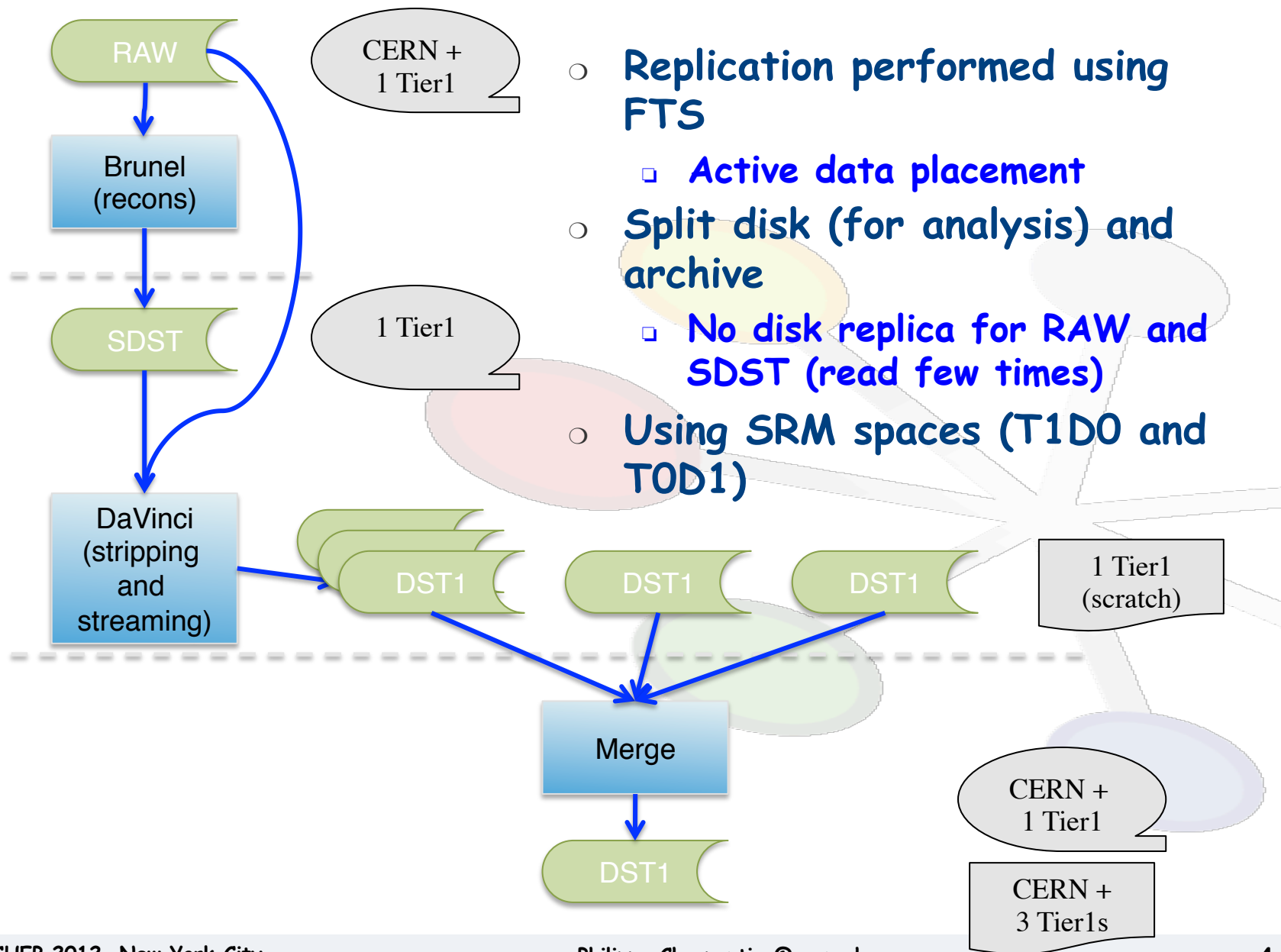  - ❑ **Many presentations at CHEP (again this time, talks and posters)**

Poster #272, S.Roiser

Poster #145, F.Stagni

RAW

Brunel
(recons)

SDST

DaVinci
(stripping
and
streaming)

DST1    DST1    DST1

Merge

DST1

○ **Reconstruction:**
- ❑ **50,000 events (3GB), <CPU>=20 HS06.s/evt**

○ **Stripping:**
- ❑ **50,000 events (2.5 GB), <CPU>=2 HS06.s/evt**
- ❑ **14 streams**

○ **Merge:**
- ❑ **40,000 events (5 GB)**

○ **Input files are downloaded on local disk**

- Replication performed using FTS
  - **Active data placement**
- **Split disk (for analysis) and archive**
  - **No disk replica for RAW and SDST (read few times)**
- **Using SRM spaces (T1D0 and T0D1)**

○ **Granularity at the file level**

   ❑ **Data Management operations (replicate, remove replica, delete file)**

   ❑ **Workload Management: input/output files of jobs**

○ **LHCbDirac perspective**

   ❑ **DMS and WMS use LFNs to reference files**

   ❑ **LFN namespace refers to the origin of the file**

      ✩ **Constructed by the jobs (uses production and job number)**

      ✩ **Hierarchical namespace for convenience**

      ✩ **Used to define file class (tape-sets) for RAW, SDST, DST**

      ✩ **GUID used for internal navigation between files (Gaudi)**

○ **User perspective**

   ❑ **File is part of a dataset (consistent for physics analysis)**

   ❑ **Dataset: specific conditions of data, processing version and processing level**

      ✩ **Files in a dataset should be exclusive and consistent in quality and content**

o **Logical namespace**

- **Reflects somewhat the origin of the file (run number for RAW, production number for output files of jobs)**
- **File type also explicit in the directory tree**

o **Storage Elements**

- **Essential component in the DIRAC DMS**
- **Logical SEs: several DIRAC SEs can physically use the same hardware SE (same instance, same SRM space)**
- **Described in the DIRAC configuration**
  - ☆ **Protocol, endpoint, port, SAPath, Web Service URL**
  - ☆ **Allows autonomous construction of the SURL**
  - ☆ `SURL = srm:<endPoint>:<port><WSUrl><SAPath><LFN>`
- **SRM spaces at Tier1s**
  - ☆ **Used to have as many SRM spaces as DIRAC SEs, now only 3**
  - ☆ **LHCb-Tape (T1D0) custodial storage**
  - ☆ **LHCb-Disk (T0D1) fast disk access**
  - ☆ **LHCb-User (T0D1) fast disk access for user data**

- ○ **Currently using the LFC**
  - ❑ **Master write service at CERN**
  - ❑ **Replication using Oracle streams to Tier1s**
  - ❑ **Read-only instances at CERN and Tier1s**
    - ☆ **Mostly for redundancy, no need for scaling**
- ○ **LFC information:**
  - ❑ **Metadata of the file**
  - ❑ **Replicas**
    - ☆ **Use "host name" field for the DIRAC SE name**
    - ☆ **Store SURL of creation for convenience (not used)**
      - ❆ **Allows lcg-util commands to work**
  - ❑ **Quality flag**
    - ☆ **One character comment used to set temporarily a replica as unavailable**
- ○ **Testing scalability of the DIRAC file catalog**
  - ❑ **Built-in storage usage capabilities (per directory)**

*Poster: A.Tsaregorodtsev*

LHCb DATA MANAGEMENT

- o **User selection criteria**
  - □ **Origin of the data (real or MC, year of reference)**
    - ☆ **LHCb/Collision12**
  - □ **Conditions for data taking of simulation (energy, magnetic field, detector configuration…**
    - ☆ **Beam4000GeV-VeloClosed-MagDown**
  - □ **Processing Pass is the level of processing (reconstruction, stripping…) including compatibility version**
    - ☆ **Reco13/Stripping19**
  - □ **Event Type is mostly useful for simulation, single value for real data**
    - ☆ **8 digit numeric code (12345678, 90000000)**
  - □ **File Type defines which type of output files the user wants to get for a given processing pass (e.g. which stream)**
    - ☆ **RAW, SDST, BHADRON.DST (for a streamed file)**
- o **Bookkeeping search**
  - □ **Using a path**
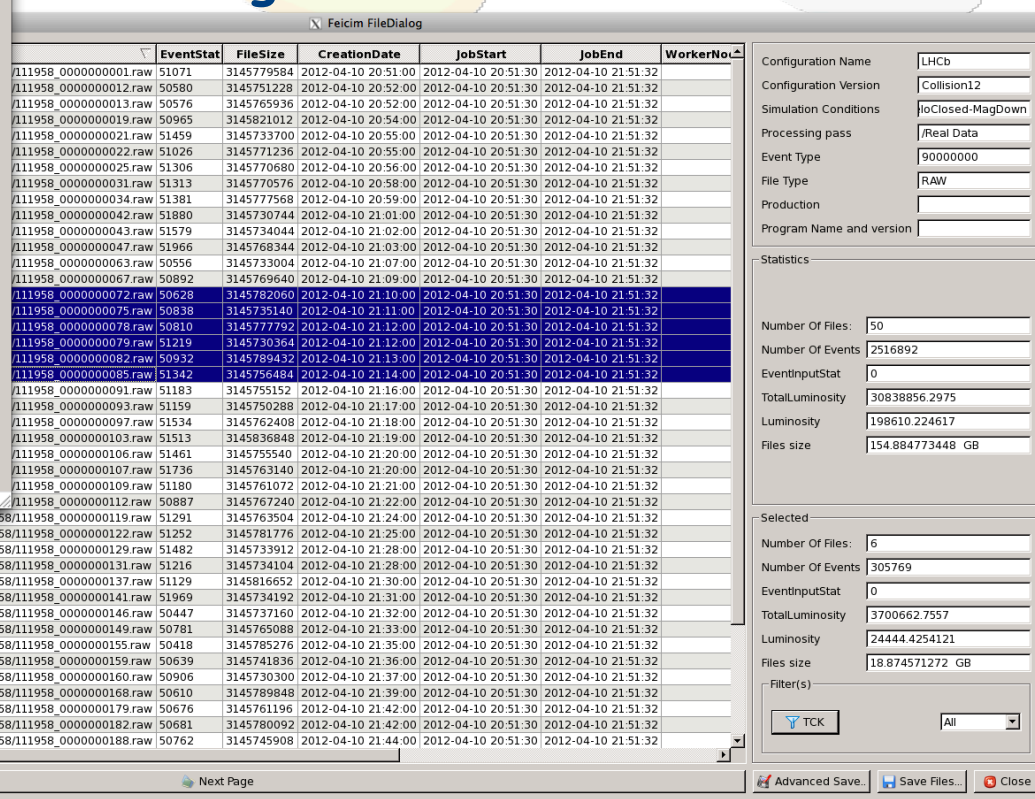    - ☆ **/<origin>/<conditions>/<processing pass>/<event type>/<file type>**

- Much more than a dataset catalog!
- Full provenance of files and jobs
  - Files are input of processing steps ("jobs") that produce files
  - All files ever created are recorded, each processing step as well
    - Full information on the "job" (location, CPU, wall clock time…)
- BK relational database
  - Two main tables: "files" and "jobs"
  - Jobs belong to a "production"
  - "Productions" belong to a "processing pass", with a given "origin" and "condition"
  - Highly optimized search for files, as well as summaries
- Quality flags
  - Files are immutable, but can have a mutable quality flag
  - Files have a flag indicating whether they have a replica or not

# Bookkeeping browsing
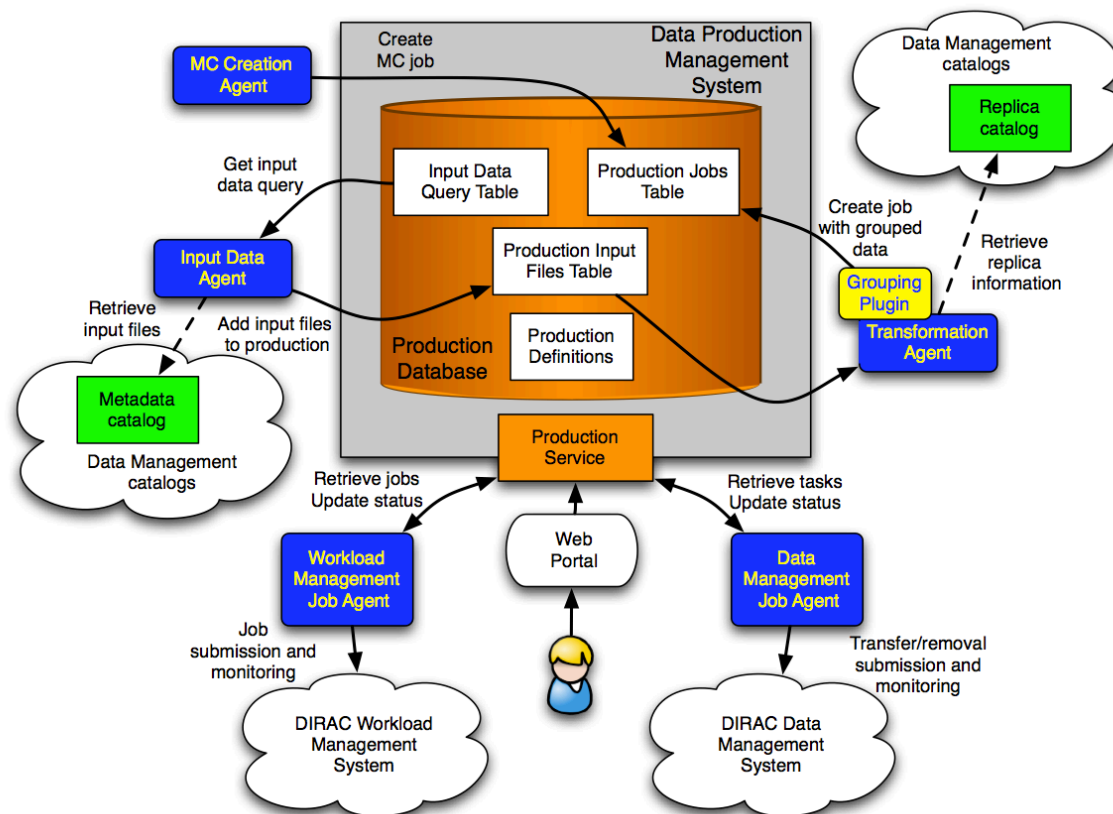
○ **Allows to save datasets**
  ❑ **Filter, selection**
  ❑ **Plain list of files**
  ❑ **Gaudi configuration file**
○ **Can return files with only replica at a given location**

# Dataset based transformation



- Same mechanism used for jobs (workload management tasks) and data management
- Input datasets based on a bookkeeping query
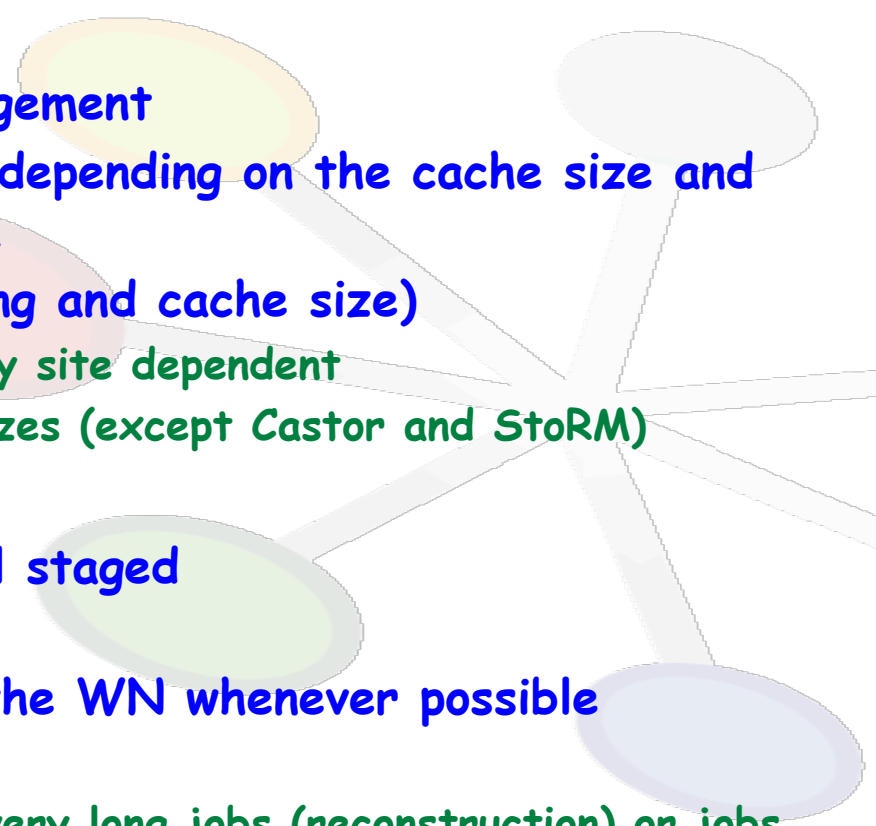  - Tasks are data driven, used by both WMS and DMS

# Data management transformations

- ○ **Replication transformations**
  - ❑ **Uses a policy implementing the Computing Model**
  - ❑ **Creates transfer tasks depending on the original location and space availability**
  - ❑ **Replication using a tree (not all from the initial source)**
    - ✰ **Optimized w.r.t. channel recent bandwidth and SE locality**
  - ❑ **Transfers are whenever possible performed using FTS**
    - ✰ **If not possible, 3rd party gridftp transfer**
      - ❄ **If no FTS channel or for user files (special credentials)**
  - ❑ **Replicas are automatically registered in the LFC**
- ○ **Removal transformations**
  - ❑ **Used for retiring datasets or reducing the number of replicas**
  - ❑ **Used exceptionally to completely remove files (tests, bugs…)**
  - ❑ **Replica removal protected against last replica removal!**

  - ❑ **Transfers and removal use the Data Manager credentials**

# Staging: using files from tape

- If jobs use files that are not online (on disk)
  - Before submitting the job
  - Stage the file from tape, and pin it on cache
- Stager agent
  - Performs also cache management
  - Throttle staging requests depending on the cache size and the amount of pinned data
  - Requires fine tuning (pinning and cache size)
    - Caching architecture highly site dependent
    - No publication of cache sizes (except Castor and StoRM)
- Jobs using staged files
  - Check first the file is still staged
    - If not reschedule the job
  - Copies the file locally on the WN whenever possible
    - Space is released faster
    - More reliable access for very long jobs (reconstruction) or jobs using many files (merging)

- Transfer accounting
  - **Per site, channel, user…**
- Storage accounting
  - **Per dataset, SE, user…**
  - **User quotas**
    - Not strictly enforced…

**LFN size per Processing Pass**
30 Days from 2012-04-07 to 2012-05-07

TB

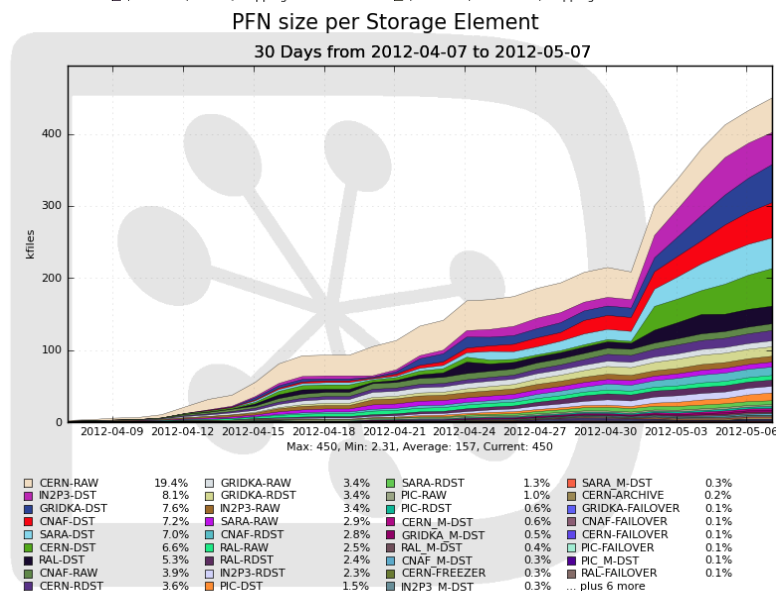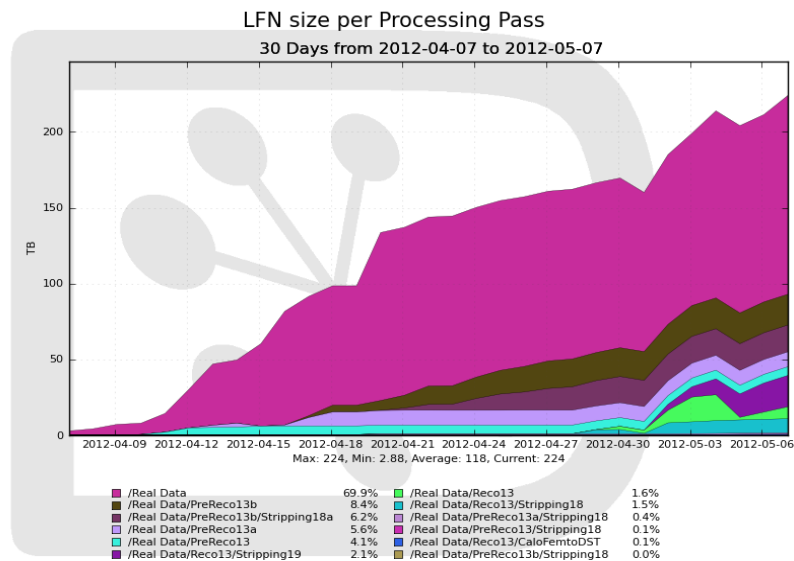Max: 224, Min: 2.88, Average: 118, Current: 224

| | | | |
|---|---|---|---|
| /Real Data | 69.9% | /Real Data/Reco13 | 1.6% |
| /Real Data/PreReco13b | 8.4% | /Real Data/Reco13/Stripping18 | 1.5% |
| /Real Data/PreReco13b/Stripping18a | 6.2% | /Real Data/PreReco13a/Stripping18 | 0.4% |
| /Real Data/PreReco13a | 5.6% | /Real Data/PreReco13/Stripping18 | 0.1% |
| /Real Data/PreReco13 | 4.1% | /Real Data/Reco13/CaloFemtoDST | 0.1% |
| /Real Data/Reco13/Stripping19 | 2.1% | /Real Data/PreReco13b/Stripping18 | 0.0% |

**Throughput by Channel**
24 Hours from 2012-04-14 08:30 to 2012-04-15 08:30 UTC

MB / s

Max: 2,123, Average: 231

Generated on 2012-04-15 08:34:43 U

| | | | | | |
|---|---|---|---|---|---|
| CERN-RAW -> IN2P3-RAW | 22.7% | CERN-RAW -> SARA-RAW | 17.3% | CERN-RAW -> CERN-FREEZER | 0.7% |
| CERN-RAW -> CNAF-RAW | 22.3% | CERN-RAW -> RAL-RAW | 15.4% | | |
| CERN-RAW -> GRIDKA-RAW | 18.7% | CERN-RAW -> PIC-RAW | 3.0% | | |

**PFN size per Storage Element**
30 Days from 2012-04-07 to 2012-05-07

kfiles

Max: 450, Min: 2.31, Average: 157, Current: 450

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| CERN-RAW | 19.4% | GRIDKA-RAW | 3.4% | SARA-RDST | 1.3% | SARA_M-DST | 0.3% |
| IN2P3-DST | 8.1% | GRIDKA-RDST | 3.4% | PIC-RAW | 1.0% | CERN-ARCHIVE | 0.2% |
| GRIDKA-DST | 7.6% | IN2P3-RAW | 3.4% | PIC-RDST | 0.6% | GRIDKA-FAILOVER | 0.1% |
| CNAF-DST | 7.2% | SARA-RAW | 2.9% | CERN_M-DST | 0.6% | CNAF-FAILOVER | 0.1% |
| SARA-DST | 7.0% | CNAF-RDST | 2.8% | GRIDKA_M-DST | 0.5% | CERN-FAILOVER | 0.1% |
| CERN-DST | 6.6% | RAL-RAW | 2.5% | RAL_M-DST | 0.4% | PIC-FAILOVER | 0.1% |
| RAL-DST | 5.3% | RAL-RDST | 2.4% | CNAF_M-DST | 0.3% | PIC_M-DST | 0.1% |
| CNAF-RAW | 3.9% | IN2P3-RDST | 2.3% | CERN-FREEZER | 0.3% | RAL-FAILOVER | 0.1% |
| CERN-RDST | 3.6% | PIC-DST | 1.5% | IN2P3_M-DST | 0.3% | … plus 6 more | |

Generated on 2012-05-07 15:20:11 UTC

TC

- Improvements on staging
  - Improve the tuning of cache settings
    - Depends on how caches are used by sites
  - Pinning/unpinning
    - Difficult if files are used by more than one job
- Popularity
  - Record dataset usage
    - Reported by jobs: number of files used in a given dataset
    - Account number of files used per dataset per day/week
  - Assess dataset popularity
    - Relate usage to dataset size
  - Take decisions on the number of online replicas
    - Taking into account available space
    - Taking into account expected need in the coming weeks
  - First rely on Data Manager receiving an advice
    - Possibly move to automated dynamic management

# Conclusions

○ **LHCb uses two views for data management:**

- ❑ **Dataset view as seen by users and productions**
- ❑ **Files view as seen by Data Management tools (replication, removal)**

○ **Datasets are handled by the LHCb Bookkeeping system (part of LHCbDirac)**

○ **File view is generic and handled by the DIRAC DMS**

○ **The LHCb DMS gives full flexibility for managing data on the Grid**

○ **In the future LHCb expect to use popularity criteria for deciding on the number of replicas for each dataset**

- ❑ **Should give more flexibility to the Computing Model**