# Event Index Core (Hadoop & HBase)

- **Import**
- Physical Containers
- Trigger Info
- Datasets Overlap
- Journal
- Important Missing Pieces

Andrea Favareto (testing)

<u>Julius Hrivnac</u>

Fedor Prokoshin (trigger)

Rainer Toebbicke (infrastructure, HBase EL)

Ruijun Yuan (graphics)

El WS. 14 June 2016



### **Import**

MC 2012 3

MC 2013 0

MC 2016 0

MC 2015 36692

IMPORTED

- Some MC datasets imports is very slow due to their size

  x\*100000 MapFiles per dataset
- Upload to HBase runs twice a day (and lasts several minutes)
  - Can be run more often

Runs smoothly

Temporary files (\*.2) cleaned automatically

500-700 imports per day

El datasets imported quickly

- Nightly checks for various inconsistencies
  - Then found problems handled by hand
    - Currently 4 datasets with different nevents in EI and AMI
- Hourly checks for duplicated events
  - Almost 100 already found, all EI15.1\*DAOD or MC15.1
  - Tagged in Catalog as multi:<multiplicity>x<number>
    - catalog -query 'id:MC15.1 multi:2' -filter 'id multi'
  - Should be fixed ? (easy)
- aiatlas016 had problems (do anyone know why ?)
  - All services transparently moved to aiatlas054 during that period



### Physical Containers (Event Lookup)

- Physical containers introduced (in Catalog)
  - As \*.0 TagSets (TagSet = a set of TagFiles)
- Existing containers created
  - Most of them are incomplete (missing:<n> tag in Catalog records that)
- Unclear control flow (who/when/how records new Containers)
- Available to EL soon

- EL -indexer is no more equivalent to EL -simple/-rich
  - Indexer implements agreed strategies (for AMI tags, etc...)
  - Simple/Rich gives always everything (from HDFS TagFiles)
    - Useful for more detailed studies in case of problems



### Trigger Info

- Graphical output prototype
  - o In 1/10000 scale
- TI is now available as Python client and direct URL (the same syntax as for EL)

### **Event Index** Trigger Info Global Help -legend Year Projets Stream Name Prod Step Data Type Version Run Number · Catalog El15.1 ▼ data15 13TeV ▼ physics Main merge ▼ AOD ▼ f641 m151′ 00284213 Event Index (Expert Mode) -query · Event Lookup □ L1 □ L2 ● EF or HLT □ all -tlevel Trigger Info -tdecision 🖲 L1 AfterVeto or L2/EF/HLT Physics 🔾 L1 AfterPrescale or L2/EF/HLT PassThrough 🔾 L1 BeforePrescale or L2/EF/HLT Resurrected Bookmarks -toverlap | HLT e120 | loose, HLT 2e17 | vloose, HLT 2g50 | loose, HLT 2j55 | bloose L14J20.0ETA4 (disables -tlevel and -tdecision) System Journal (for admins) -email (implies asynchronous execution) -tpercent O Complete Result Percent of overlaps Search Reset

Overlaps	HLT_e120	HLT_2e17	HLT_2g50	HLT_2j55	HLT_g120
HLT_e120	10000	1.0	2.0	1.0	86.0
HLT_2e:.HLT_e120_loose	1.0	10000	3.0	0.0	1.0
HLT_2g50	2.0	3.0	10000	0.0	14.0
HLT_2j55	1.0	0.0	0.0	10000	5.0
HLT_g120	86.0	1.0	14.0	5.0	10000

Progress map: 100%



### Datasets Overlap

- Functionality available since a long time:
- Number of overlapping events between two datasets
   number of duplicated events in merged dataset
  - There should be no duplicated events in either dataset
- Uses aux M/R job
- Use -email option as calculation can be slow for big datasets
  - => runs offline and sends results in email
- The API will be generalised to accept more datasets in one task.
- More user-friendly wrapper & GUI will be provided
- Should start systematic creation of overlap matrix for new / old datasets (could start immediately)
  - Recorded (in Catalog) ?
  - Available from Journal anyway (once executed).

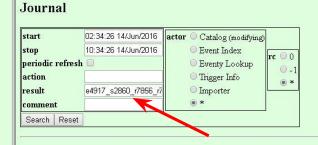




- Free text search possible for action, result and comment
- A lot of information already in Journal (& Catalog) - all results are there:
  - Can be exploited

### **Event Index**

- Global Help
- · Catalog
- Event Index (Expert Mode)
- Event Lookup
- Trigger Info
- · Bookmarks
- System Journal (for admins)



Entries of , e4917 s2860 r7856 r7676 in result between 02:34:26 14/Jun/2016 and 10:34:26 14/Jun/2016:

actor ▲▼	rc AV	<u>time</u> ▲▼	timesstamp AV	action AV	result ▲▼	<u>comment</u> <b>▲▼</b>	<u>id</u> ▲▼
EICLI	0	51	03:33:32.669 14/Jun/2016	{-outname=MC15.1/mc15_5TeV.410013.PowhegPythiaEvtGmore	Info: 50000 events found, result MC15.1.mc15_5TeV more		2016.06.14.03.33.32.569/atlevind
CatalogCLI	0	2	03:33:37.496 14/Jun/2016	(-modify=project.mc15_5TeV runNumber:410013 stream more	Info: 2s spend Result: 1 : TagFile(id: MC15.1.mc1 more		2016.06.14.03.33.37.398/atlevind



### Important Missing Pieces

- Commands with Shibboleth protection
  - Currently only Graphical GUI and direct URL are protected and available outside CERN
  - Prevents updates of CVMFS distribution!
- Dataflow logistics
  - Updating/removing datasets
  - Synchronising with AMI
    - AMI status changes too often
    - Maybe should be checked for each command
      - That would make it slower
- Import of BIG MC datasets
- Update TWiki (some obsolete info)
  - Up-to-date info is always available via CLI help
    - and TagConvertor Web Home
    - and Web Service GUI



## Info

Web Service: <a href="https://atlas-event-index.cern.ch/EIHadoop">https://atlas-event-index.cern.ch/EIHadoop</a>

**Documentation & Distribution:** <a href="https://atlas-event-index.cern.ch/doc">https://atlas-event-index.cern.ch/doc</a>

Twiki: <a href="https://twiki.cern.ch/twiki/bin/viewauth/AtlasComputing/EventIndex">https://twiki.cern.ch/twiki/bin/viewauth/AtlasComputing/EventIndex</a>

Twiki Tutorial: <a href="https://twiki.cern.ch/twiki/bin/viewauth/AtlasComputing/EventIndexTutorial">https://twiki.cern.ch/twiki/bin/viewauth/AtlasComputing/EventIndexTutorial</a>

**Sources:** svn+ssh://svn.cern.ch/reps/atlasoff/Database/TAGHadoop/TagConvertor

AFS: /afs/cern.ch/sw/lcg/external/Java/TagConvertor

**CVMFS:** \$ Isetup eiclient