

Event Index Core *(Hadoop & HBase)*

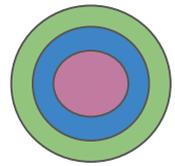


- Architecture
- Import
- Trigger
- Event Lookup
- Export
- Data Compression
- Documentation
- Plans

Andrea Favareto
[Julius Hrivnac](#)
Fedor Prokoshin
Grigorij Rybkin
Rainer Toebicke
Ruijun Yuan

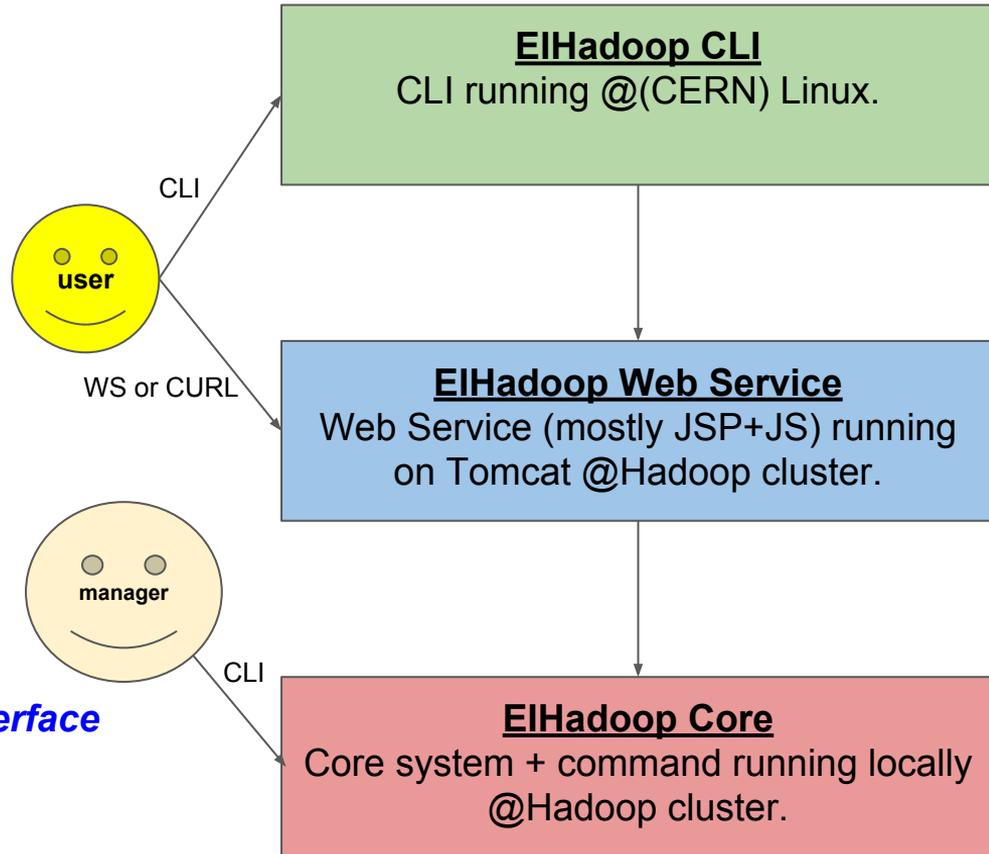


Architecture



- Interface-implementation separation
- The same logic in all clients (WS, CLI)
- **Linux-like** principle CLI client
 - Interface style familiar to users
 - Compatible with other Atlas components
 - Very powerful
 - Usage patterns easy to document
 - Universal Portable CLI under development (to be run on any platform)
- Automatic deployment of all clients
- Integrated help
- All code is public and designed to be called from other components
- Some bad (legacy) names and conventions :-)

Users seem to be happy with the existing CLI interface when they know it does exist





Architecture - Deployment

- All code contained in one JAR files, packed in different ways according to intended use
 - EIHadoop.jar for the main distribution
 - EIHadoopEL.jar,... as Java executables (differs from main distribution by *Manifest*, which specifies the entry point & run conditions)
 - EIHadoop.war to be deployed on Tomcat
- All commands (catalog, ei, el, ti, inspect) implement the same style CLI interface
- All Help is included (help = man page + script used for nightly test)
- Full deployment (Core, WS, CLI, ...) via ant install; ant deploy-prod
 - + mail to CVMFS maintainers
- WebService & (remote) CLI are decorated clients of the Core CLI
 - WebService specifies additional constraints for users + graphics
 - Three incarnations:
 - Graphical “sexy” view
 - CLI server
 - CURL server
 - CLI connects to WebService, automatically finds out which one
- Remote commands are somehow simplified versions of the core commands
 - Dangerous operations are not available



Architecture - Deployment - Example

- TICLI.java implements ti command @Hadoop cluster (in EIHadoop.jar)
- TI.jsp implements WebService server for remote CLI
- TIMenu.jsp, TIView.jsp implements graphical WebService
- TIdirect.jsp implements server for CURL clients
- TIRMT.java implements remote ti command
 - Trivial => CVMFS update rarely needed

in EIHadoopTI.jar

```
public class TIRMT extends RMTCLI {  
    public static void main(String[] args) throws Exception {  
        new TIRMT().doit(args);  
    }  
}
```

in EIHadoop.war

```
RMTWSUtils utils = new RMTWSUtils(request);  
utils.setArguments(new String[]{"query",  
                                "tlevel",  
                                "tdecision",  
                                "toverlap",  
                                "tpercent",  
                                "export",  
                                "email",  
                                "help"});  
utils.setDefault("tlevel", "3");  
utils.setDefault("tdecision", "0");  
utils.setHelpCond(new String[]{"query"});  
out.println(utils.execute());
```

New functionality easy to introduce and deploy on all clients.

The same pragmatics on all clients.



Architecture - Deployment - Example

```
# to find all 2016 datasets with duplicated events
```

```
catalog -query 'id:EI16.1 multi:~null' -filter 'id multi nevents amiNevents'
```

```
# to check a statistical distribution of LumiBlockN of a dataset for some fired trigger
```

```
ei -query 'dataset:data15_13TeV.00279515.physics_Main.merge.AOD.r7562_p2521' from  
-mr 'trigFired("HLT_tau35_medium1_tracktwo_tau25_medium1_tracktwo_L1TAU20IM_2TAU12IM")' select  
-aux 'net.hep.atlas.Database.EIHadoop.Accessor.Aux.VarStat;-v BunchId' execute/accumulate
```

***Too complicated for users, who have managed Athena & Root ?
(all is documented on man-like help & FAQ)***



Import

- Re-configured to merge before import
 - Much faster
 - Used already before to import big MC datasets
 - Further simplification will follow
 - Soon, may not be needed ?
- Currently importing mostly MC
 - 10x datasets per day
 - Not used

Imported EVENTS

```
=====
```

| | | | |
|---------|-----|-----|-----|
| EI 2009 | 34 | 444 | 939 |
| EI 2010 | 1 | 196 | 985 |
| EI 2011 | 1 | 839 | 177 |
| EI 2012 | 3 | 313 | 189 |
| EI 2013 | 331 | 475 | 303 |
| EI 2014 | 436 | 763 | 032 |
| EI 2015 | 17 | 798 | 354 |
| EI 2016 | 20 | 888 | 617 |
| MC 2009 | | | 0 |
| MC 2010 | | | 0 |
| MC 2011 | | | 0 |
| MC 2012 | | 300 | 000 |
| MC 2013 | | | 0 |
| MC 2014 | | | 0 |
| MC 2015 | 27 | 248 | 486 |
| MC 2016 | 10 | 020 | 000 |



Trigger

- Trigger statistics tables generated for all 2016 runs
 - Now processing 2015,... runs
- New Catalog variable `trigger:bad` for incorrectly decoded triggers
 - 17 in 2016 (*express_express*, *physics_BphysicsDelayed*, *physics_CosmicCalo*)
- Thinking about coupling Trigger with Event Lookup
 - Either on the search or on the report part



Event Lookup

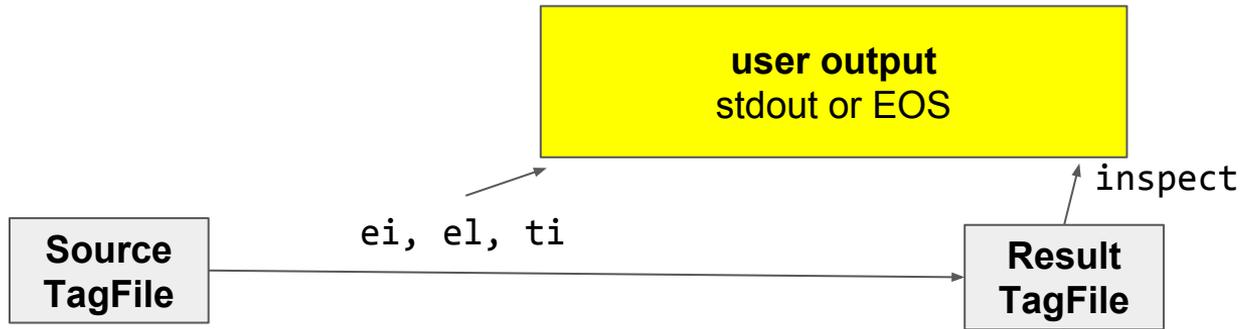
- Grigori is now managing EL code
 - Consolidation
 - Aligning -simple/-rich with -indexer
 - Indexer uses HBase - faster, but derived
 - Simple/Rich uses HDFS - slower, but more primary
 - Using trigger information
 - What to do with MC ?



Export

Results can be exported to EOS using `-export eos` argument to all commands

- `ei`, `catalog`, `el`, `ti`, `inspect`
 - Small problem with `el`, because it uses different client (following PAnda request)
- Exporting results (ie. answer from CLI command), not the resulting TagFiles
 - Resulting TagFiles can be exported using `inspect` command
- Useful together with `-email` argument
 - User gets information about results in EOS by email
- Some problems with authentication



Example:

*Overlap table is created by `ti` command
and viewed by `inspect` command*

Ordinary dataset-TagFiles, Overlap tables, Statistics tables,...



Data Compression

- Compression:
 - Transformation uncompressed -> compressed TagFiles tested
 - 5-10x space gain
- Another possibility: Use different data format
 - Kudu ?
 - As evaluated by Zbysek
 - Could just convert each TagFiles separately so that existing framework keeps working
 - Some modifications needed when code relies on Map file format
 - Open for possible future migration to more modern system
 - Maybe could merger EI.1 -> EI.0
 - Was planned, but difficult to execute with Map files
 - Could have fun with Spark & comp.



Documentation

➤ Home page

- Architecture + practical instructions (where to find, how to compile,...)

➤ Integrated help

- Linux-like man pages + set of actual nightly tests
- Available from CLI & Web Service

➤ Use Case page

- High level description of functionality
 - Not much details
- To be linked from Use Case TWiki:
 - <https://atlas-event-index.cern.ch/doc/usecases/#EL>
 - <https://atlas-event-index.cern.ch/doc/usecases/#TStats>
 -

➤ Frequently Asked Questions

- Detailed recipe-like description of typical tasks
- Should contain all information for a user to use the system
 - Cut-and-paste style
- Quickly growing
 - Tricky to keep up-to-date & consistent
- Can be linked from outside:
 - <https://atlas-event-index.cern.ch/doc/faq/Right.html#TOverlaps>
 - <https://atlas-event-index.cern.ch/doc/faq/Right.html#RecentDatasets>
 - ...

*All documentation is inside
EIHadoop distribution - to create
consistent view across all
platforms & clients.*

What can be done

How to do it



Frequently Asked Questions

Event Index Frequently Asked Questions

Access

- [How to setup command line environment for Event Index ?](#)
- [Where is Event Index Web Service ?](#)
- [Where is Event Index Documentation ?](#)

Datasets

- [How to find whether a datasets is available and what are their properties ?](#)
- [How to find which datasets have been recently imported ?](#)
- [How to search for specific datasets ?](#)
- [How to count number of datasets with some characteristics ?](#)
- [Is there a Web Service for datasets exploration ?](#)
- [How to get event overlaps between datasets ?](#)
- [How to find if a dataset has duplicated events ?](#)
- [How to get a statistic of a variable \(and which variables are available\) ?](#)
- [Which variables can be accessed ?](#)

Events

- [How to perform detailed event-level search ?](#)
- [How to get GUIDs for specific events ?](#)
- [How to get GUIDs for specific events using just CURL in CERN ?](#)
- [How to get GUIDs for specific events using just CURL on any SLC machine ?](#)

Event Index Frequently Asked Questions

How to setup command line environment for Event Index ?

Using CVMFS:

```
export ATLAS_LOCAL_ROOT_BASE=/cvmfs/atlas.cern.ch/repo/ATLASLocalRootBase
alias setupATLAS='source ${ATLAS_LOCAL_ROOT_BASE}/user/atlasLocalSetup.sh'
setupATLAS
asetup 19.0.0
localSetupEIClient
```

Within CERN AFS:

```
$ source /afs/cern.ch/project/jps/repos/atlas-eies/TagConventor/bin/setup.sh
```

This setup will give the Event Index commands catalog, ei, el, ti and inspect. Those commands (except el) give a simple help with examples when called without arguments. *It is possible to install Event Index environment on non-CVMFS/non-CERN-AFS machines. Contact [J.Hrivnac](#) for details.*

Where is Event Index Web Service ?

It is [here](#). It is protected by the CERN SSO.

Where is Event Index Documentation ?

It is [here](#). You can also consult [Event Index TWiki Pages](#) and [Event Index Tutorial](#).

How to find which datasets are available and what are their properties ?

```
# to get all available information about a dataset
catalog -query dataset:data16_hip8TeV.00314170.physics_ZeroBias.recon.AOD.f781_m1741
# to get some information about a dataset (status, number of events and when it was imported)
catalog -query dataset:data16_hip8TeV.00314170.physics_ZeroBias.recon.AOD.f781_m1741 -filter 'status nevents imported'
```

Other interesting characteristics can be asked via -filter argument:

```
imported: the date and time of the import
```



FAQ - Current Status

Access

How to setup command line environment for Event Index ?

Where is Event Index Web Service ?

Where is Event Index Documentation ?

Datasets

How to find whether a datasets is available and what are their properties ?

How to find which datasets have been recently imported ?

How to search for specific datasets ?

How to count number of datasets with some characteristics ?

Is there a Web Service for datasets exploration ?

How to get event overlaps between datasets ?

How to find if a dataset has duplicated events ?

How to get a statistic of a variable (and which variables are available) ?

Which variables can be accessed ?

Events

How to perform detailed event-level search ?

How to get GUIDs for specific events ?

How to get GUIDs for specific events using just CURL in CERN ?

How to get GUIDs for specific events using just CURL on any SLC machine ?

Triggers

How to get Trigger Overlaps for a dataset ?

How to get Trigger Statistics for a dataset ?

How to check triggers ?

How to create Trigger different Overlap or Statistics tables ?

Misc

How to execute a command which takes a long time ?

How to get big results ?



Plans

- Universal Portable CLI
 - Will run everywhere, even on iPhone
- Trigger statistics for all data
- Connection between Trigger and Event Lookup
 - Would it be useful ?
- Event Lookup Consolidation (Grigorij)
- Data compression / New format (Zbysek ?)
- Finish implementation of *TagSets* (= sets of *TagFiles*) to handle Physics Containers
 - In the system since the beginning, needs re-animation
- Commands with Shibboleth protection on any Linux
 - Not just SLC
- More data-oriented Web Service (Justin)
 - Current WS is service-oriented
- Code & Doc & Test-suit cleaning
 - A lot of dead wood



Info

Web Service: <https://atlas-event-index.cern.ch/ElHadoop>

Documentation & Distribution: <https://atlas-event-index.cern.ch/doc>

Use Cases: <https://atlas-event-index.cern.ch/doc/usecases>

Frequently Asked Questions: <https://atlas-event-index.cern.ch/doc/faq>

Sources: `svn+ssh://svn.cern.ch/repos/atlasoff/Database/TAGHadoop/TagConvertor`

AFS: `/afs/cern.ch/project/jps/repos/atlas-eies/TagConvertor`

CVMFS: `$ lsetup eicient`

EOS: `/eos/atlas/atlascerngroupdisk/proj-evind/Results`