

---

# Big Data Analytics Tools

## Applied to ATLAS Event Data

---

Ilija Vukotic  
University of Chicago

CHEP 2016, San Francisco



Big Data technologies have proven to be very useful for storage, visualization and analysis of the ATLAS distributed computing (meta)data. If we could apply the same technologies on real/MC data, we could:

- Do quick investigations without writing any code, directly in a web browser
- Deliver data to users:
  - Only events passing cuts and only variables user needs.
  - No need for experiment specific libraries - make it easier to use machine learning tools like Spark, Jupyter, R, SciPy, Caffe, TensorFlow. This would simplify Machine Learning Challenges: Higgs Boson ML Challenge, the Tracking challenge, ...
- Serve as a backend to:
  - Event Viewers (VP1, ATLANTIS, ATLASrift,...), make them platform independent
  - All kinds of education and outreach tools

The question is would the same technology map well on event data? Would event data structure map correctly, would data size be acceptable, what indexing and retrieval rate could we achieve?

# Related work



**EventIndex** - a complete catalogue of all ATLAS events.

- Collection:ActiveMQ, storage:Hadoop, indexing:HBase, custom WebUI, CLI.
- ~1300 B/event. Event identifier, trigger pattern, hit counts, pointers to the event in all processing stages.
- Used for: event picking, event skimming, Panda event service
- Contains  $\sim 40 \cdot 10^9$  data and  $25 \cdot 10^9$  MC event records. Searches at 10ms/event.

**Open data** - <http://atlas-opendata.web.cern.ch/>

- Provides access to real and simulated data, and educational tools, to the public
- Three levels of sophistication (simple visualizations, Jupyter based analysis, full set of analysis tools to download - documentation, datasets, software, VMs)

Lukas Heinrich's **project Aretha** <https://www.youtube.com/watch?v=S5SkTXhdCng>

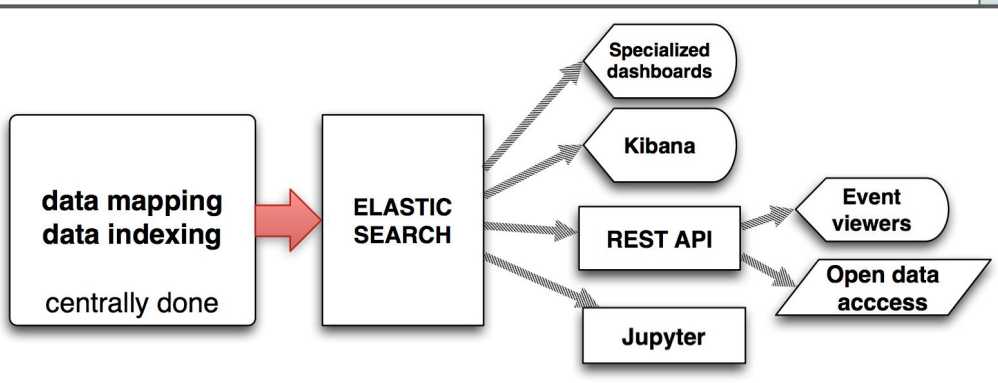
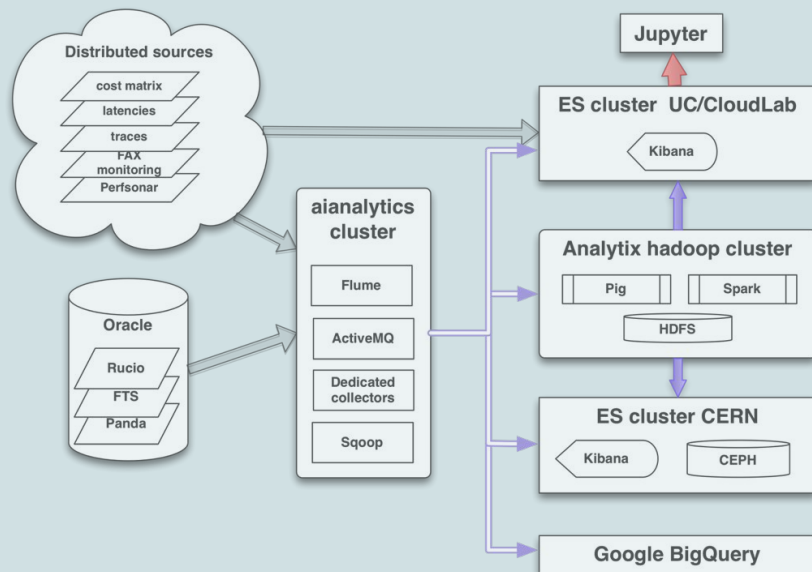
- Streams generated HepMC data from SHERPA, Herwig, MadGraph (locally or in cloud) to ZeroMQ
- Runs ZeroMQ data collector, Elasticsearch, Kibana all in docker
- Makes plots while jobs running, validate data using eg. Rivet

# Analytics platform

There is a number of possible workflows that would work for storing and processing event data while using already existing ADC analytics platform:

- Store data in Hadoop (avro) process using pig or Spark. Not interactive, not a data access solution.
- Store in HDF5, use Spark. Performant analysis, does not solve data access.
- ...

## ATLAS distributed computing analytics platform



Our selected workflow:

- Data indexed in Elasticsearch
- Fast visualization in Kibana
- Open data access through Elasticsearch REST API
- Data analysis on a co-located Jupyter cluster

# Elasticsearch cluster



Elasticsearch is a distributed, real-time data and analytics, high availability open source search engine built on top of Apache Lucene. Works with structured JSON documents, schema free, by default all fields are indexed, data are compressed.



Indices:	Memory: <b>246GB</b> /	Total Shards:	Unassigned Shards:	Documents:	Data:	Uptime: 4
<b>2088</b>	<b>433GB</b>	<b>18644</b>	<b>0</b>	<b>8,481,799,362</b>	<b>9TB</b>	<b>months</b>



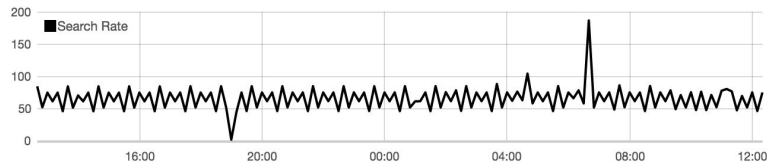
5 data + 3 head nodes

E5-2620v4, 64GB RAM, 4x800 GB SSD, 10Gbps NIC

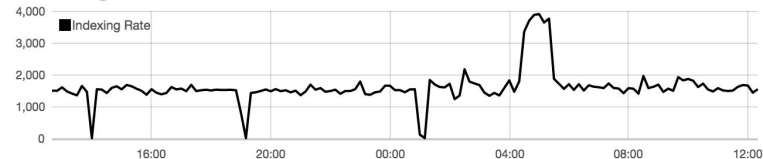
Used for ATLAS Distributed Computing analytics

Quite high base load

Search Rate: 74.46 /s



Indexing Rate: 1,536.95 /s



## At University of Chicago

- Intel(R) Xeon(R) CPU E5-2650 v2 @ 2.60GHz x2
- Tesla K20c x2
- 128GB RAM , 4.5TB RAID-5 for /scratch, 1TB RAID-1 for /
- If needed one more machine with 2 x Xeon Phi (64 core)
- Simple universal user, pass authentication, notebooks automatically checked in a git repo.

## At CERN

- SWAN (Service for Web based Analysis)
- Currently in Beta
- Code in CERNBox

## Advantages:

- Most of ML libraries installed, if there is anything missing we can easily add it.
- Local to Elasticsearch
- Learn from codes of others. See what was already done.

- Managed by CERN
- Your notebooks are yours only
- One can have the same environment anywhere you have CVMFS, just do:

```
setupATLAS
```

```
lsetup "lcgen v -p LCG_85swan3 x86_64-slc6-gcc49-opt ROOT rootpy ..."
```

# ATLAS Data Formats

- Consider only user oriented DxAOD format.
- xAOD is a common centrally produced data format ~PB size.
- DxAODs are ~TB scale derivations, centrally produced, skimmed/slided and customized for different physics groups use.
- Mini xAODs or flat n-tuples are subgroup specific and not considered here (~GB scale) .

Type	Event size [kb/event]	Branches
DAOD data15_13TeV	30	647
AOD data15_13TeV	184	1742
DAOD mc15_13TeV	25	999
AOD mc15_13TeV	216	2512

## ROOT structure

```
*Branch :EventInfoAux.  
*Entries : 6752 : BranchElement (see below)  
* .....*  
*Br 2 :EventInfoAux.xAOD::AuxInfoBase : BASE  
*Entries : 6752 : Total Size= 529349 bytes File Size = 21915  
*Baskets : 17 : Basket Size= 104448 bytes Compression= 24.12  
* .....
```

float32 AntiKt4EMTopoJetsAuxDyn Sp **numpy objects**

- File Size = 2124  
Compression= 13.46

```

"EventInfoAuxDyn": {
  "type": "object",
  "properties": {
    "HLT_xe70_mht": {"type": "boolean"},
    "HLT_xe70_tc_lcw": {"type": "boolean"},
    "HLT_xe80": {"type": "boolean"},
    "HLT_xe80_mht": {"type": "boolean"},
    "HLT_xe80_tc_lcw": {"type": "boolean"},
    "mpx": {"type": "double"},
    "name": {"type": "string", "index": "not_analyzed"},
  }
}

```

## Indexed event



# Indexing events

We indexed di-jet events xAODs, and DxAODs from DAOD\_EXOT2 derivation.

Used indexing method proved rather inefficient (parsing all the event structure for all events), also process was not tuned (batch size, etc).

type	Event size		Indexing rate [Hz]
	[kb/event]	increase	
DAOD data15_13TeV	146	5.4	21
AOD data15_13TeV	377	2.0	8
DAOD mc15_13TeV	116	4.6	19
AOD mc15_13TeV	385	1.8	9

# Data visualizations

## Bad

Visualizations not made with scientists in mind.  
Only simple operations supported.

## Good

Visualizations can be organized and saved in nice reusable dashboards.  
As all computation are distributed over all the ES nodes, visualizations are very fast even with billions of events.



# Filtering and aggregations

## Streaming full events:

```
"query" : { "match_all": {} }
```

## To return only some variables add:

```
"fields" : ["AntiKt4EMTopoJetsAuxDyn.ActiveArea", ...],
```

## To select on a variable (cut):

```
"query" : { "term" : { "EventInfoAuxDyn.HLT_xe80" : True } }
```

## Complex conditions (must, must\_not, should, filter):

```
'bool':{  
  'must':[  
    {"term" : { "EventInfoAuxDyn.HLT_xe80" : False }},  
    {"range" : { "AntiKt4EMTopoJetsAuxDyn.GhostAntiKt2TrackJetCount" : { "gt": 2 } }},  
    {"range" : { "AntiKt4LCTopoJetsAuxDyn.JetEMScaleMomentum_m": { "gt": 10000 } }},  
    {"exists" : { "field" : "HLT_xAOD__JetContainer_a4tcemsubjesFSAuxDyn.HECQuality"  
  ]  
}}
```

**Aggregations (cut and count analysis), very powerful, fast, can be nested. In our tests all are sub-second duration operations.**

## Simple count:

```
"aggs": { "docs": { "value_count": { "field": "_type" } } }
```

## Returning stats:

```
"aggs": { "ptStats": { "stats": { "field": "TrackPt" } } }
```

## Stats in bins:

```
"aggs": { "ptStats": { "histogram": {  
  "field": "TrackPt",  
  "Interval" : 50  
} } }
```

# Skimming & slimming performance



Read test	DAOD data15_13TeV		AOD data15_13TeV		DAOD mc15_13TeV		AOD mc15_13TeV	
	Rate [Hz]	CPU* [%]	Rate [Hz]	CPU [%]	Rate [Hz]	CPU [%]	Rate [Hz]	CPU [%]
1 variable from all events	486	1	264	1	520	1	225	1
10 variables from all events	453	3	263	1	491	3	227	1
10 variables from events passing cut ( ~2 % events)	606 / 32 kHz	22	251 / 13 kHz	10	971 / 50 kHz	15	226 / 11 kHz	20
Full events passing cut ( ~2% events)	161 / 6.3 kHz	77	35 / 1.8 kHz	84	121 / 6.2 kHz	80	32 / 1.5 kHz	82
Streaming all full events	116	86	39	87	120	88	31	86

Selections are basically free, cost is in data transfer back to client and client parsing JSON.

# Data modifications

**Simple fixes** can be done by:

- Adding corrected variable
- Defining “scripted” variables

**Re-calibration** from existing data would be performed by a central batch re-indexing.

Any kind of **reweighting** can be done by defining custom scoring function.

## Weighting:

```
"script_score": {  
  "script": {  
    "lang": "lang",  
    "params": {  
      "param1": value1,  
      "param2": value2  
    },  
    "Inline": "_score * doc['orig_pt'].value /  
      pow(param1, param2)"  
  }  
}
```

# Conclusions

---

- Feasibility of actual event indexing is tested
- Input/output rates are reasonable, increase in data volume is at the expected level.
- All kinds of optimizations are possible
  - Indexing directly from ATLAS analysis framework
  - Optimized mapping
  - Compressed responses
- Visualizations, skimming and sliming are very fast, very convenient for cut flow type analyses
- Full data retrieval while moderately slow, consumes too much CPU in parsing JSON returned.
- Elasticsearch is easy to scale up to a data center level (~ few hundreds of servers), bigger than that one need to add “tribe” nodes.