



Metadata Curation and Integration in High Energy and Nuclear Physics

[Maria Grigorieva](#), [Marina Golosova](#) , [Alexander Alexeev](#),
[Maxim Gubin](#), [Vasiliy Aulov](#), [Alexei Klimentov](#) and
[Torre Wenaus](#)

National Research Center "Kurchatov Institute"
Tomsk Polytechnic University
Brookhaven National Laboratory

Data Knowledge Base Highlights

DKB Basic Consideration

- ❑ Organizing metadata in ATLAS, so as to provide a holistic view on physics topics, including integrated representation of all ATLAS documents (papers, drafts, supporting documents, conference notes, Indico meetings, Twiki pages, etc) and corresponding data samples.

DKB Evolution

- ❑ The most important reports and summaries are made in **Twiki pages** (collaborative documentation) in semi-manual mode
- ❑ The metainformation in Twiki doesn't provide mechanisms for synchronization with database back-ends.
 - ❑ Provide **fully automatic** metadata search and aggregation by arbitrary set of parameters, synchronized with the existing database backends.
 - ❑ Index metadata and provide a **quick and flexible** google-like metadata **search**, categorization and aggregation.

ATLAS Metadata Sources

Public Results

- CERN Document Server
- CERN Twiki
- Indico
- GLANCE (Papers and ConfNotes)

Data Management and Analysis

- Production System Database
- ATLAS Metadata Interface (AMI)
- Rucio DDM
- Google Docs
- JIRA
- ATLAS SVN Repository
- BigPanDA Monitoring System

Physics and Experiment Environment

- ATLAS Geometry database
- ATLAS Conditions Database
- ATLAS Twiki Pages

ATLAS Metadata Issue

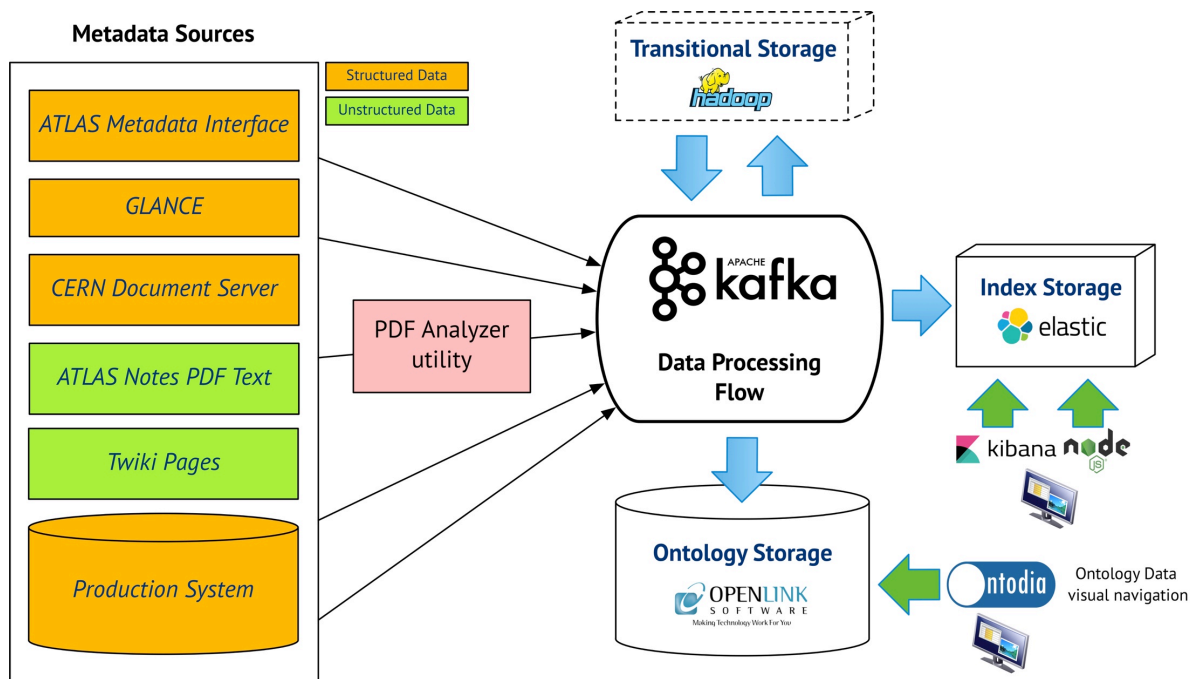
Most of these metadata sources exist *autonomously*. To gain the comprehensive information of research study, including experimental environment, data samples, and available results, scientists need to obtain *intersections among metadata* from different sources by themselves, as there are *no automatic tools* providing data integration.

Due to the complexity of modern HENP experiments, this task becomes more and more challenging.

DKB is aiming to fill this gap in meta-data integration.

DKB program code on GitHub:
<https://github.com/PanDAWMS/dkb>

DKB Architecture Prototype



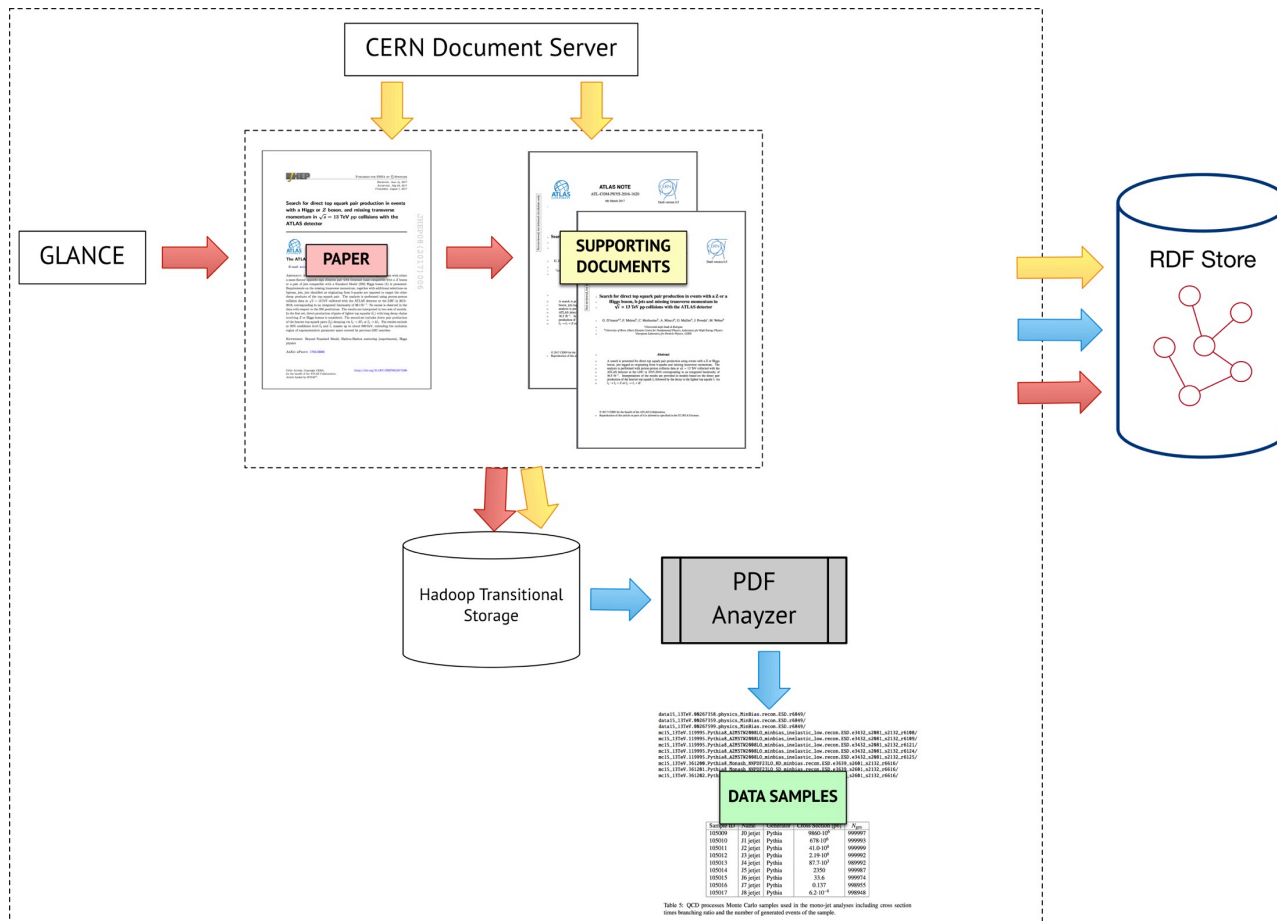
Technology evaluation and system prototype architecture:

- ✓ RDF storage: **Virtuoso**
- ✓ Transitional storage: **Hadoop**
- ✓ Metadata streaming&processing:
Apache Kafka
- ✓ Knowledge Base navigation
(Web GUI): **Ontodia**
- ✓ Index storage: **ElasticSearch**
- ✓ Index search web-interfaces:
Kibana, NodeJS
- ✓ Documents metadata mining:
**PdfMiner utility and
implemented PDF Analyzer
tool**

Basic programming language -



Initial Data Processing Flow



PDFAnalyzer tool extracts metadata from TXT and XML representation of PDF text, by regular expressions and context analysis.

28.09.17

Data Extraction from PDF Documents

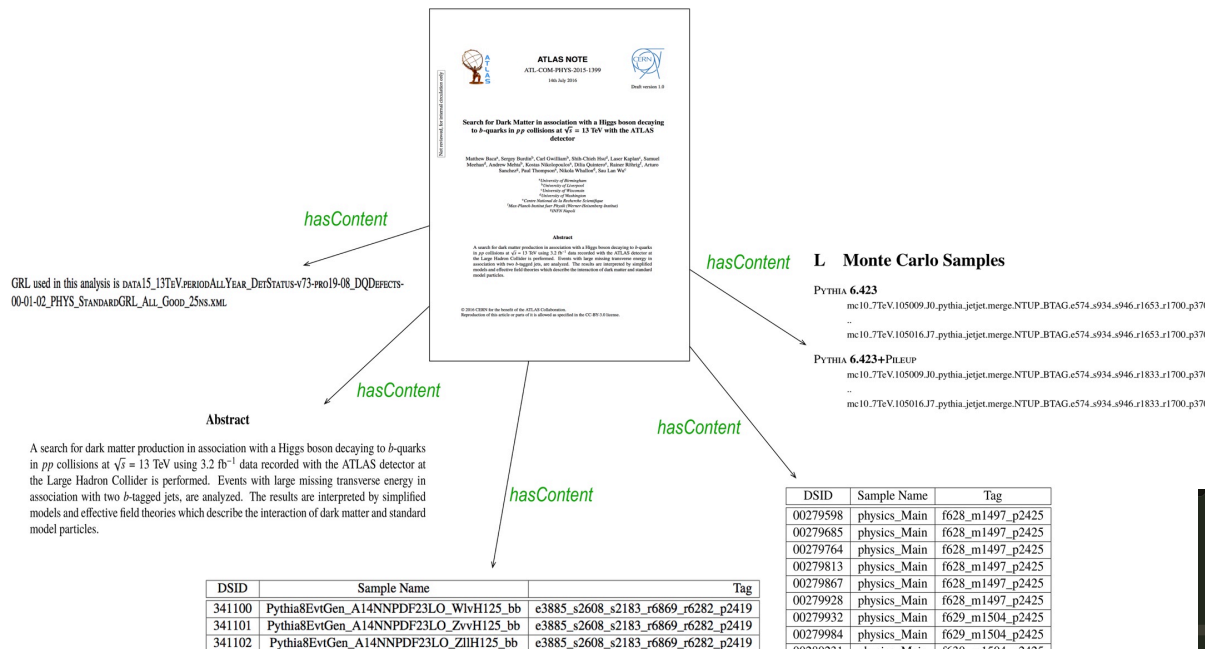


Table 19: Monte Carlo samples used as baseline for Standard Model VH(\rightarrow bb).

PDFAnalyzer extracts:

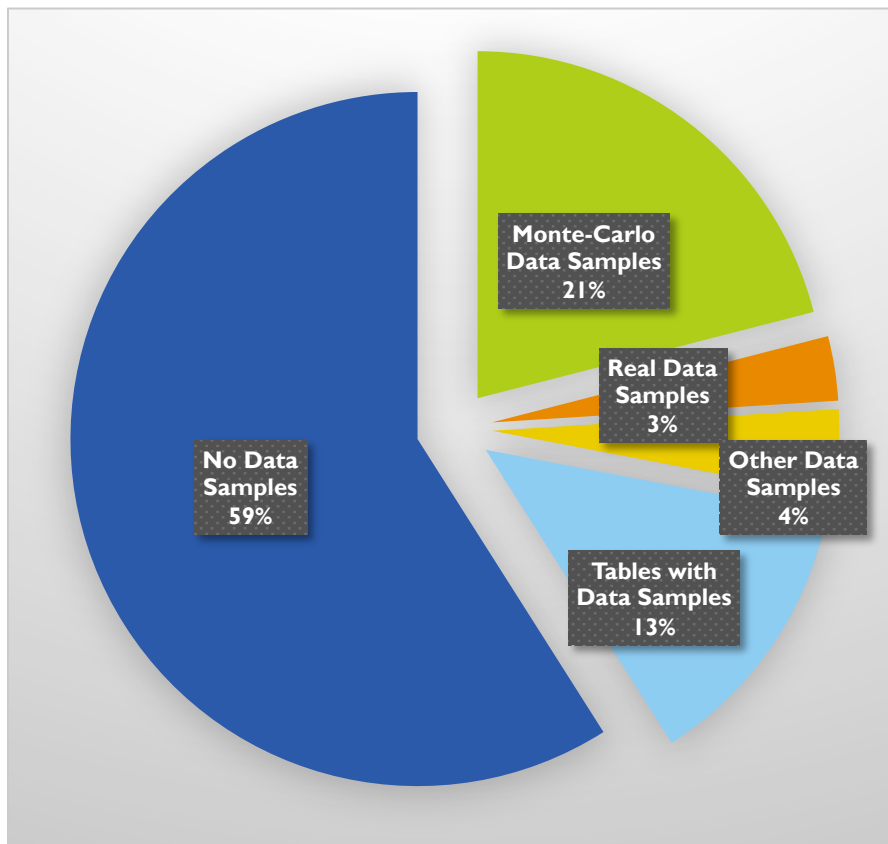
- dataset names by regular expression
- datasets metadata from tables
- experiment-specific metadata from text

Returns structured metadata in **JSON**.

Has GUI interface, providing manual correction of analysis results

```
{
  "content": {
    "real_datasets": [
      "data0_h1_00169175.physics_bulk.ROOT.NTUP_HI_r2114_r2768_p753",
      "data1_2p76TeV_00219257.physics_MinBias.merge.NTUP_HI_r519_m1313",
      "data1_2p76TeV_00219263.physics_MinBias.merge.NTUP_HI_r519_m1313",
      "data1_2p76TeV_00219295.physics_MinBias.merge.NTUP_HI_r519_m1313",
      "data1_2p76TeV_00219305.physics_MinBias.merge.NTUP_HI_r519_m1313",
      "data1_2p76TeV_00219322.physics_MinBias.merge.NTUP_HI_r519_m1313",
      "data1_2p76TeV_00219364.physics_MinBias.merge.NTUP_HI_r519_m1313"
    ],
    "plain_text": {
      "atlas_name": "ATL-COM-PHYS-2016-XXXX",
      "links": {},
      "data_taking_year": false,
      "luminosity": false,
      "energy": "2.76 TeV",
      "campaigns": [
        "mc12",
        "mc11"
      ],
      "collisions": "proton-proton",
      "mc_datasets": [
        "mc10_2TeV.119114.Hijing_PbPb_2p75TeV_MinBias_Flow_3JFV6.recon.NTUP_HI_e1846_s1160_s1161_d490_r2134",
        "mc12_2TeV.209836.Pythia_AUET2BCTEQ6L1_minbias_inelastic.recon.NTUP_HI_e2812_s1647_s1586_r6309_tid06162221_00"
      ]
    },
    "fname": "CDS_CERN-ATL-COM-PHYS-2016-257"
  }
}
```

PDF Analyzer Statistics

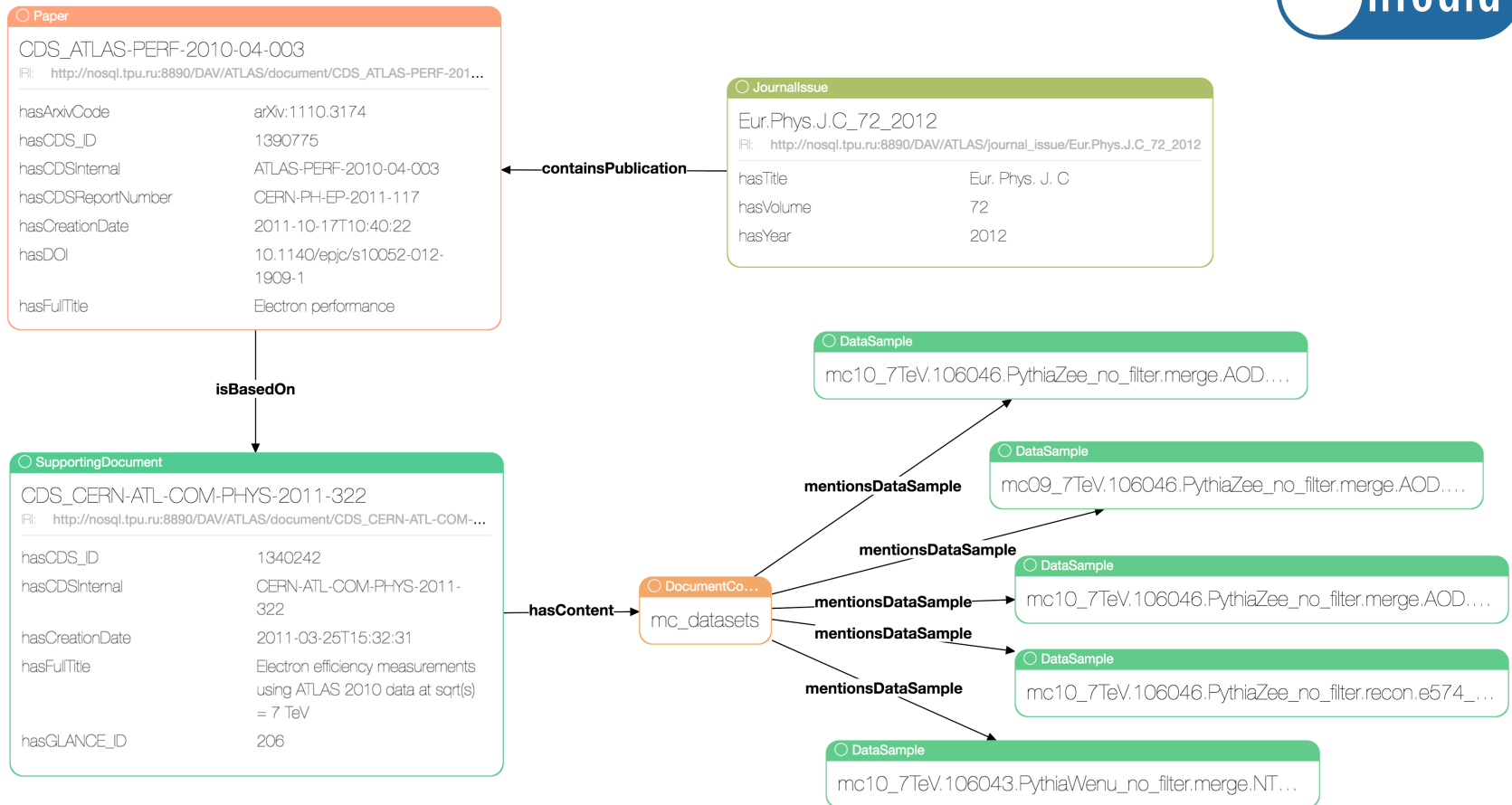


The statistics of ATLAS Internal Notes (500 documents) analysis showed that ~40% of documents allow to extract data samples automatically.

Why we can't extract data samples for for 60% of documents?

- ✓ No datasets in the text.
- ✓ Information about datasets is presented as instructions for human reader.
- ✓ Not all Papers have properly defined Supporting Documents.

Ontodia - JavaScript library for Virtuoso navigation



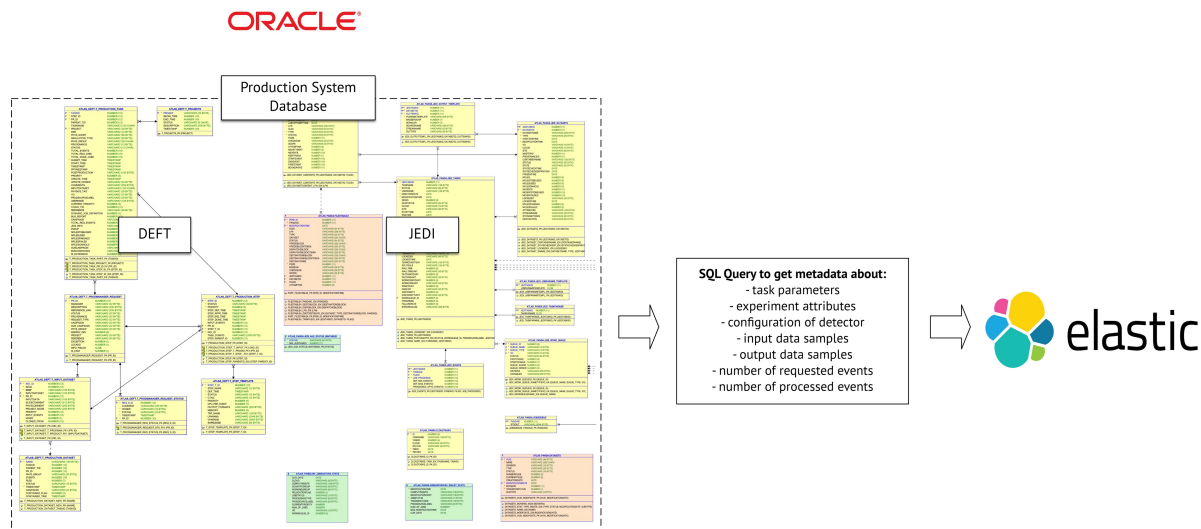
Metadata Indexing and Search Facilities

■ To provide fast and flexible data categorization, search and aggregation

1. Reproduce Event Summary report, but with physics category breakdown (like in Twiki Event Summary).
2. Implement google-like search of tasks and data samples by the arbitrary set of attributes, like campaign, project, ATLAS geometry, Condition Tags, hashtags, physics category, and others.

&MC16c_CP

	Evgen	Evgen Merge	Simul	Merge	Reco	Rec Merge
Step Name %						
Processed/Input events	1,913,235,610 /	1,236,857,070 /	1,934,968,180 /	313,326,050 /	1,925,870,230 /	1,914,226,230 /
running/pending/finished	1,956,298,240	1,284,022,270	1,989,539,330	310,709,550	1,930,593,230	1,915,887,230
	0.0%/0.0%/97.8%	0.0%/0.0%/96.33%	0.49%/0.0%/97.26%	0.08%/0.0%/100.84%	0.06%/0.0%/99.76%	0.03%/0.0%/99.91%



Metadata in ElasticSearch Index

Metadata categories:

- Task Parameters

- Taskid
- Taskname
- Status
- Timestamp
- Start time
- End time
- Request ID
- Ticket ID
- User Name

- Experiment parameters

- Energy
- Campaign/Subcampaign
- Project
- Physics group
- Physics category
- Hashtags
- Run number

- Configuration

- ATLAS geometry
- Conditions tags
- SW Release
- Trigger Config

- Events

- Requested
- Processed

- Data Samples

- Input
- Output

t _id	Q Q [] * 12115125
t _index	Q Q [] * prodsys
# _score	Q Q [] * -
t _type	Q Q [] * MC16
t architecture	Q Q [] * x86_64-slc6-gcc62-op
t campaign	Q Q [] * MC16
t conditions_tags	Q Q [] * OFLCOND-MC16-SDR-20
t core_count	Q Q [] * 8
t description	Q Q [] * MC16c PA sheet 1
o end_time	Q Q [] * September 12th 2017, 20:32:57.000
# energy_gev	Q Q [] * 13,000
t geometry_version	Q Q [] * ATLAS-R2-2016-01-00-01
t hashtag_list	Q Q [] * MC16c_PA
t input_datasets	Q Q [] * mc16_13TeV:mc16_13TeV.341080.PowhegPythia8EvtGen_CT10_AZNLOCTEQ6L1_VBFH125_WWlvlv_EF_15_5.simul.HITS.e3871_e5984_s3126_tid12023009_00, mc16_13TeV:mc16_13TeV.361238.Pythia8EvtGen_A3NNPDF23L0_minbias_inelastic_low.simul.HITS.e4981_s3087_s3111/, mc16_13TeV:mc16_13TeV.361239.Pythia8EvtGen_A3NNPDF23L0_minbias_inelastic_high.simul.HITS.e4981_s3087_s3111/
t output_datasets	Q Q [] * mc16_13TeV.341080.PowhegPythia8EvtGen_CT10_AZNLOCTEQ6L1_VBFH125_WWlvlv_EF_15_5.recon.A00.e3871_e5984_s3126_r9781_tid12115125_00
t phys_category	Q Q [] * Higgs
t phys_group	Q Q [] * MCGN
t pr_id	Q Q [] * 13326
# processed_events	Q Q [] * 16,750,000
t project	Q Q [] * mc16_13TeV
# requested_events	Q Q [] * 40,242,000
t run_number	Q Q [] * 341080
o start_time	Q Q [] * September 8th 2017, 21:56:56.000
t status	Q Q [] * done
t step_name	Q Q [] * Reco
t subcampaign	Q Q [] * MC16c
o task_timestamp	Q Q [] * September 14th 2017, 15:35:07.000
t taskid	Q Q [] * 12115125
t taskname	Q Q [] * mc16_13TeV.341080.PowhegPythia8EvtGen_CT10_AZNLOCTEQ6L1_VBFH125_WWlvlv_EF_15_5.recon.e3871_e5984_s3126_r9781
t ticket_id	Q Q [] * ATLPSTASKS-1144711
t trans_home	Q Q [] * Athena-21.0.32
t trans_path	Q Q [] * Reco_tf.py
t trans_uses	Q Q [] * Atlas-21.0.3
t trigger_config	Q Q [] * RD0toRD0Trigger-MCRECO:DBF:TRIGGERDBMC:2170,46,199
t user_name	Q Q [] * dsouth
t vo	Q Q [] * atlas

Google-like keyword search

new save open share 22/11/15 2015 this year

Summaries in Kibana

TTbarX: Category

Step ↕	Requested Events ↕	Processed Events ↕
Simul	3,200,000	3,200,000
Reco	442,512,000	170,637,000
Rec Merge	2,600,000	2,600,000
Merge	200,000	200,000
Evgen Merge	1,800,000	1,800,000

DrellYan: Category

Step ↕	Requested Events ↕	Processed Events ↕
Simul	194,913,900	194,913,900
Reco	12,444,403,000	11,579,977,000
Rec Merge	161,408,000	161,408,000
Merge	449,800	449,800
Evgen	2,202,000	2,202,000
Deriv Merge	176,500,000	176,500,000
Deriv	161,810,000	161,810,000

Higgs: Category

Step ↕	Requested Events ↕	Processed Events ↕
Simul	106,328,100	106,328,100
Reco	9,628,753,000	5,850,766,000
Rec Merge	86,838,000	86,838,000
Evgen Merge	138,399,700	138,399,700
Evgen	38,570,400	38,570,400
Deriv Merge	82,000,000	82,000,000
Deriv	77,510,000	77,510,000

```
"query": {
  "bool": {
    "must": [
      { "term": { "subcampaign.keyword": "MCI 6a" } },
      { "term": { "status": "done" } }
    ],
    "should": [
      { "term": { "hashtag_list": "MCI 6a" } },
      { "term": { "hashtag_list": "MCI 6a_CP" } }
    ]
  }
},
"aggs": {
  "category": {
    "terms": { "field": "phys_category" },
    "aggs": {
      "step": {
        "terms": { "field": "step_name.keyword" },
        "aggs": {
          "requested": {
            "sum": { "field": "requested_events" }
          },
          "processed": {
            "sum": { "field": "processed_events" }
          }
        }
      }
    }
  }
}
```


Web-interface prototype

DKB/DCC Whiteboard

Data Knowledge Base Prototype for ATLAS Collaboration. Processing papers & internal documents full text, search data samples, used in the data analysis.

Search form

Keywords

MC16a, Higgs, ATLAS-R2-2016-01-00-01, Reco, OFLCOND-MC16-SDR-16

Start Date

September 2017						
Su	Mo	Tu	We	Th	Fr	Sa
					1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28	29	30

eters

GET prodsys/MC16/_search

```
{
  "query": {
    "bool": {
      "must": {
        "query_string": {
          "query": "\"MC16a\" AND \"Higgs\" AND \"ATLAS-
R2-2016-01-00-01\" AND \"Reco\" AND \"OFLCOND-
MC16-SDR-16\"",
          "analyze_wildcard": true
        }
      }
    },
    "filter": {
      "range": {
        "task_timestamp": {
          "gte": "01-05-2017 00:00:00",
          "lt": "10-08-2017 00:00:00",
        }
      }
    }
  }
}
```

DKB/DCC Whiteboard

Data Knowledge Base Prototype for ATLAS Collaboration. Processing papers & internal documents full text, search data samples, used in the data analysis.

Dataset Search

"MC16a" AND "Higgs" AND "ATLAS-R2-2016-01-00-01" AND "Reco" AND "OFLCOND-MC16-SDR-16"

mc16_13TeV.345318.PowhegPythia8EvtGen_NNPDPF30_AZNLO_WpH125J_Hyy_Wincl_MINLO.recon.e5734_e5984_s3126_r9364					
TASK		EXPERIMENT		CONFIGURATION	
taskID	11828859	Campaign	MC16	Step Name	Reco
taskName	mc16_13TeV.345318.PowhegPythia8EvtGen_NNPDPF30_AZNLO_WpH125J_Hyy_Wincl_MINLO.recon.e5734_e5984_s3126_r9364	Subcampaign	MC16a	ticket_id	ATLPSTASKS-1108307
status	running	Project	mc16_13TeV	Architecture	x86_64-slc6-gcc62-op
Description	MC16a PA sheet 0	Energy		Core Number	
timestamp	31-08-2017 07:17:56	Physics Group	MCGN	ATLAS Geometry	ATLAS-R2-2016-01-00-01
start time	06-08-2017 21:24:08	Physics Category	Higgs	Conditions Tags	OFLCOND-MC16-SDR-16
end time		Hashtags	MC16a_PA	trigger_config	RDOtoRDOTrigger=MCRECO:DBF:TRIGGERDBMC:2136,35,16
Duration				trans_path	Reco_tf.py
EVENTS Requested / Processed / Done(%)				trans_path	AtlasOffline-21.0.20
40242000 / 10530000 / 26%				run_number	345318
DATASETS					
Input Datasets	mc16_13TeV:mc16_13TeV.345318.PowhegPythia8EvtGen_NNPDPF30_AZNLO_WpH125J_Hyy_Wincl_MINLO.simul.HITS.e5734_e5984_s3126_tid11828856_00 mc16_13TeV:mc16_13TeV.361238.Pythia8EvtGen_A3NNPDF23LO_minbias_inelastic_low.simul.HITS.e4981_s3087_s3111/ mc16_13TeV:mc16_13TeV.361239.Pythia8EvtGen_A3NNPDF23LO_minbias_inelastic_high.simul.HITS.e4981_s3087_s3111/				
Output Datasets	mc16_13TeV.345318.PowhegPythia8EvtGen_NNPDPF30_AZNLO_WpH125J_Hyy_Wincl_MINLO.recon.AOD.e5734_e5984_s3126_r9364_tid11828859_00				

mc16_13TeV.345433.PowhegPythia8EvtGen_NNPDPF3_AZNLO_WmH125J_MINLO_qqWWlvlv.recon.e5811_e5984_s3126_r9364

mc16_13TeV.345326.PowhegPythia8EvtGen_NNPDPF3_AZNLO_WmH125J_MINLO_lWWWlvlv.recon.e5823_e5984_s3126_r9364

mc16_13TeV.345337.PowhegPythia8EvtGen_NNPDPF3_AZNLO_ZH125J_MINLO_lIWWWlvlv.recon.e5810_e5984_s3126_r9364

mc16_13TeV.341429.PowhegPythia8EvtGen_CT10_AZNLO_WpH125J_MINLO_eveWWWlvlv.recon.e4210_e5984_s3126_r9364

mc16_13TeV.345319.PowhegPythia8EvtGen_NNPDPF30_AZNLO_ZH125J_Hyy_Zincl_MINLO.recon.e5743_e5984_s3126_r9364

mc16_13TeV.341433.PowhegPythia8EvtGen_CT10_AZNLO_WpH125J_MINLO_qqWWlvlv.recon.e3938_e5984_s3126_r9364

What's been done:

- Developed the method and implemented modules, providing automatic extraction of the data samples from ATLAS documents
 - Ontological data model for ATLAS documents, data samples and Experiment Attributes
 - Tools for metadata extraction/processing/conversion/importing
 - PDF Analyzer, providing metadata extraction from PDF documents
 - Kafka-based automation of dataflows execution
 - Virtuoso database was filled with metadata
 - Ontodia as GUI for Virtuoso navigation
- Provided fast and flexible data categorization and search
 - Production System database partly (MC16 campaign) indexed in ElasticSearch
 - ElasticSearch and Kibana infrastructure is installed at CERN
 - Kibana dashboard with Event Summary report, and a set of diagrams
 - Implemented NodeJS-based web-interface for ElasticSearch Storage

Future Plans

- ▣ DKB software transferring to CERN machines
- ▣ Synchronization between Production System database and ElasticSearch storage
- ▣ Performance tests of Production System and ElasticSearch storage:
 - ▣ Aggregations
 - ▣ Search by arbitrary set of parameters
- ▣ Web Interface for ElasticSearch storage
- ▣ Improve the mechanism of the Supporting Notes search, using CERN Document Server Publications history
- ▣ Improved Papers&Supporting Documents search results will be tested on neo4j or OrientDB graph database

Acknowledgment

- The work was supported by Russian Foundation for Basic Research (RFBR) under contract No. 16-37-00246.

Thanks

- Siarhei Padolski,
- Michail Borodin,
- Dimitry Krasnopevtsev,
- Dmitry Golubkov,
- Anastasia Kaida