

1 EI3 – THE ATLAS EVENTINDEX FOR LHC RUN 3

2 **F.V. Prokoshin^{1 a}, I.N. Aleksandrov¹, E.I. Alexandrov¹, Z. Baranowski², D.**
3 **Barberis³, G. Dimitrov⁴, A. Fernandez Casani⁵, E. Gallas⁶, C. Garcia**
4 **Montoro⁵, S. Gonzalez de la Hoz⁵, J. Hrivnac⁷, A.I. Kazymov¹, M.A.**
5 **Mineev¹, G. Rybkin⁷, J. Sanchez⁵, J. Salt⁵, M. Villaplana Perez⁸**

6 ¹ *Joint Institute for Nuclear Research, 6 Joliot-Curie St., Dubna, Moscow Region, 141980, Russia*

7 ² *CERN-IT, Geneva, Switzerland*

8 ³ *INFN Genova and Universita' di Genova, Dipartimento di Fisica, Genova, Italy*

9 ⁴ *CERN-ATLAS, Geneva, Switzerland*

10 ⁵ *Instituto de Fisica Corpuscular (IFIC), Centro Mixto Universidad de Valencia - CSIC, Valencia,*
11 *Spain*

12 ⁶ *Department of Physics, Oxford University, Oxford, United Kingdom*

13 ⁷ *LAL, Universite Paris-Sud, CNRS/IN2P3, Universite Paris-Saclay, Orsay, France*

14 ⁸ *Department of Physics, University of Alberta, Edmonton AB, Canada*

15 E-mail: ^afedor.prokoshin@cern.ch

16 The ATLAS Event Index provides since 2015 a good and reliable service for the initial use cases
17 (mainly event picking) and several additional ones, such as production consistency checks, duplicate
18 event detection and measurements of the overlaps of trigger chains and derivation datasets. LHC Run
19 3 will see increased data-taking and simulation production rates, with which the current infrastructure
20 would still cope but may be stretched to its limits by the end of Run 3. This talk describes a new
21 implementation of the front and back-end services that will be able to provide at least the same
22 functionality as the current one for increased data ingestion and search rates and with increasing
23 volumes of stored data. It is based on a set of HBase tables, with schemas derived from the current
24 Oracle implementation, coupled to Apache Phoenix for data access; in this way we will add to the
25 advantages of a BigData based storage system the possibility of SQL as well as NoSQL data access,
26 allowing us to re-use most of the existing code for metadata integration.

27 **Keywords:** Event Index, ATLAS computing, Database, BigData

28
29 Fedor Prokoshin

30 Copyright © 2019 for this paper by its authors.
31 Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).
32
33

34 1. Introduction

35 The ATLAS Experiment [1] produces large amounts of data: several billion events per year. In
36 addition, similar amount of simulated events is being produced. This data is kept on the hundreds of
37 sites all around the globe. A database containing references to the events is necessary in order to
38 efficiently access them from the distributed data storage. The ATLAS EventIndex [2, 3] provides a
39 way to collect and store event information using modern technologies. It also provides various tools to
40 access this information through command line, GUI and RESTful API interfaces

41 An infrastructure was created that allows fast and efficient selection of events of interest from
42 the billions of events recorded, based on various criteria. EventIndex provides an indexing system that
43 points to these events in millions of files scattered through a worldwide distributed computing system.
44 It automatically detects and indexes new data that was produced in this system from events collected
45 from the detector or simulated using Monte-Carl technique.

46 2. ATLAS EventIndex

47 EventIndex is a system designed to be a complete catalog of ATLAS events, real and
48 simulated data.

49 2.1 ATLAS Events

50 Event is the basic unit of ATLAS data. Each event contains result of a single triggered
51 interactions, plus eventually piled-up interaction.

- 52 • Signals from the detector
- 53 • Reconstructed particles with their parameters
- 54 • Trigger decisions

55 Event is uniquely identified by the run number and the event number. Event information is
56 stored in many instances, that are spread among the hundreds of GRID sites. Event information may
57 have different formats and contents to fit analyses needs.

58 Each event record contains the following fields:

- 59 • Event identifiers: run and event number, trigger stream, luminosity block, bunch crossing ID
60 (BCID)
- 61 • Trigger decisions: Trigger masks for each trigger level, and decoded trigger chains, specifying
62 trigger condition passed
- 63 • References to the events at each processing stage in all permanent files generated by central
64 productions. They can be used to retrieve specific events of interest (event picking).

66 2.2 ATLAS Datasets

67 ATLAS event data is written in files that are organized in datasets. Datasets can have different
68 format depending of the processing stage: Detector data is first written in the RAW format, then AOD
69 datasets are produced after reconstruction. Derived datasets (DAOD) are produced for use in the
70 specific analyses. In addition to datasets produced from the events from the detector, simulated (MC)
71 datasets are produced on the GRID, to be used in various analyses and system calibration. EVNT
72 datasets contain particles information and RDO contains simulated detector signals. There are various

73 versions of the datasets originating from the same detector or simulated events, with different
74 reconstruction settings and software version. Datasets are being reprocessed roughly yearly.

75 2.3 Use cases

76 Originally EventIndex was intended to be used mostly for Event Picking: user may ask for
77 event in specific format and processing version. Several user cases were added based on operation
78 experience and user requests. A trigger information stored in the event record allow counting and
79 selecting events based on trigger decisions. Trigger chains overlaps are also being counted for trigger
80 tables optimizations. EventIndex appears to be useful for production completeness and consistency
81 checks, looking for data corruption, missing and/or duplicated events, derivation overlaps counting .
82 EventIndex is used for dataset browsing: finding datasets of interest, dataset report and inspection.

83 Summary of use cases can be found on the project page:

84 <https://twiki.cern.ch/twiki/bin/view/AtlasComputing/EventIndexUseCases>

86 2.4 EventIndex architecture

87 Figure 2 shows the partitioned architecture of EventIndex.

88 The Data Production system extracts EventIndex information from ATLAS files produced at
89 CERN or on the Grid [4]. The process starts automatically as soon as a new dataset is completed. The
90 information on new datasets on the GRID is obtained from the ATLAS Metadata Interface database
91 (AMI, [5, 6]). These indexing jobs also provide a data integrity test, as they are the first to run on new
92 data files. All new datasets containing events in AOD format (Analysis Object Data, i.e. the
93 reconstruction output) are indexed by default, while other datasets derived from AOD (DAODs, i.e.
94 selected events with reduced information for specific analyses) are only indexed on demand. Since the
95 extracted metadata contains also the references to the corresponding raw data, it is always possible to
96 extract events in RAW data format too. For simulated data, all event generator files are indexed as
97 well.

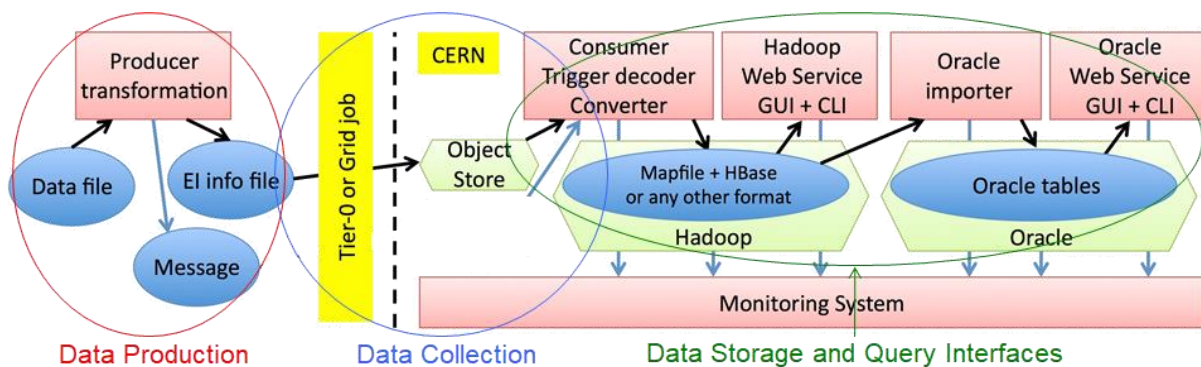


Figure 1. EventIndex architecture

102 The extracted event metadata is packed into small files and transferred to the CERN Object
103 Store by the Data Collection system [7, 8], which runs on servers at CERN. A Data Collection
104 Supervisor controls the transfers and validates indexed data. When a new dataset has been indexed, the
105 Supervisor informs the Consumer processes, which then fetch data from the Object Store and store it

106 in the Hadoop [9,10] cluster at CERN. The system is now in routine operation, with low number of
107 failures mostly originated from site problems or corrupted files.

108 Data Storage and Query system provide permanent storage for EventIndex data, fast access for
109 the most common queries, and reasonable-time response for complex queries. It uses Hadoop as a
110 baseline storage technology. The full information is stored in compressed MapFile format [11]. An
111 internal catalogue in HBase [12] keeps track of the status of each dataset and holds dataset metadata.
112 All event records are also inserted into a HBase table, which is used for fast event lookup for the event
113 picking use case. CLI, RESTful API and GUI interfaces available for data inspection, search and
114 retrieval. As of October 2019, the system holds 24 TB of data for real events and 11 TB for simulated
115 ones.

116 Reduced information from real data records (only event identification information and
117 pointers to event locations, no trigger information), are also copied to an Oracle database. The
118 connection of this database with other ATLAS databases stored in the same Oracle cluster is
119 especially useful to check dataset processing completeness and detect duplicates, providing added
120 value to this information. Oracle is also much faster in event lookup operations, if the data schema has
121 been well designed [13]. It is also being used for easy calculation of dataset overlaps. Information on
122 77 k datasets with 185 billion event records is stored there.

123 A monitoring system [14, 15] keeps tracks for health of the system components and the data
124 flow, providing a visualization of the system status and of the stored data volume. Current
125 implementation of the monitoring system uses InfluxDB [16] to store the data and Grafana [17] to
126 display information.

127 **3. EventIndex evolution towards LHC Run 3**

128 The current design of EventIndex, based on Hadoop and implemented using MapFiles and
129 HBase, has been in continuous operation since 2015 with satisfactory results. Nevertheless, the current
130 implementation of the EventIndex started showing scalability issues due to new use cases and the
131 increasing amount of stored data (the event rate increased steadily throughout Run 2 beyond the initial
132 expectation). The fast data querying for real events based on Oracle is no longer sufficient to cover all
133 requests. Therefore, an R&D programme was started to explore new technologies that would allow the
134 EventIndex to include new functionalities and to keep up with the future demanding event rates.

135

136 **3.1 Next generation EventIndex**

137 EventIndex technologies must evolve to future demanding rates of data. Currently ALL
138 ATLAS processes produce ~30 billion events records/year (up to 350 Hz on average). This is update
139 rate through the whole system, all years, real and simulated data. EventIndex read 8M files/day and
140 produce 3M files. In the future due to expected trigger rates, need to scale for next ATLAS runs. For
141 Run3 (2021-2023) an expected increase is at least half an order of magnitude: 35 B new real
142 events/year and 100 B new MC event/year. And we expect an order of magnitude increase for Run4
143 (2026-2029): 100 B new real events and 300 B new MC events per year. In addition, sum up replicas
144 and reprocessing.

145 An evolution of the EventIndex concepts is being planned. Currently the same event across
146 each processing step (RAW, ESD, AOD, DAOD, NTUP) is physically stored at different HADOOP
147 HDFS files. In the future EventIndex will be one and only one logical record per event, that will
148 contain event identification, immutable information (trigger, lumiblock, ...), and for each processing
149 step there will be additional information: link to algorithm (processing task configuration), pointer(s)
150 to output(s) and possibly flags for offline selections (derivations).

151 The EventIndex evolution includes support of virtual datasets. A virtual dataset is a logical
152 collection of events, created either explicitly (giving a collection of Event Ids) or implicitly (selection
153 based on some other collection or event attributes). Labeling individual events by a process or a use
154 with attributes (key:value)

155

156 3.2 New Use Cases

157 A number of new use cases were proposed

- 158 • Massive event picking: Selection of many events, touching a large fraction (or all) of the files
159 in a dataset. May need a dedicated service, especially if input on tape (RAW)
- 160 • Adding “offline trigger” information: Store the results of selections that can be used to form
161 derived datasets. This requires the ability to add info to part of event record.
- 162 • Using offline trigger information: Select events using online and offline trigger information to
163 build a “Virtual dataset”. Use massive event picking to physically retrieve events belonging to
164 a virtual dataset (probably in AOD format, but also RAW if very few) and continue the
165 analysis with more info on reduced size datasets. This is useful if selecting <1% of the events
- 166 • Partial processing for production tests. May skip some input checks and then assign a finite
167 lifetime to the information (delete once the test is done).

169 3.3 EI Evolution: SQL on HBase

170 Apache HBase is the Hadoop database, a distributed, scalable, big data store. It is open-
171 Source, distributed, versioned, non-relational database modeled after the Google BigTable paper.
172 HBase was built on top of HDFS and provides fast record lookups (and updates) for large tables.

173 HBase organizes data into tables. Tables have rows and columns, which store values (like a
174 spreadsheet). Rows are identified uniquely by their row key. Each row can have a different schema.
175 Data within a row is grouped by column family. Must be defined up front and not easily modified

176 HBase belongs to noSQL database family. NoSQL databases allow to deal with scalability
177 problems that relational databases were traditionally suffered. When data model is trivial - simple key-
178 value store could satisfy it. On the other hand, SQL/structured schemas provide their own advantages:
179 structured data are easy to understand and maintain, and standard declarative query logic are
180 ‘optimized’ for complex queries.

181 Various possibilities for SQL on HBase were considered: Apache Impala, Hive (handling of a
182 row key mapping must be on the application side) and Spark (mainly for batch jobs). A series of tests
183 were performed with prototype based on Apache Kudu [18, 19]. Finally, Apache Phoenix [20] was
184 chosen as the most promising platform for use in the new EventIndex. As SQL layer on top of HBase
185 it provides structured schema of the tables instead of schemaless freeride, mapping of columns to
186 HBase cells and serialization of data types to bytes. It also has SQL planner and optimizer with built-
187 in HBase related optimizations, server-side (optimized) executions and provides convenient access via
188 JDBC (Java DB Connector).

189 Phoenix provides OLTP and operational analytics for HBase through SQL. It takes SQL
190 query, compiles it into a series of HBase scans. It allows direct use of the HBase API, along with
191 coprocessors and custom filters and produces regular JDBC result sets. HBase RowKey design must
192 be adapted to Phoenix’s types and sizes, losing “some” performance. Phoenix allows use of RowKey
193 fields in queries, but they are stored as one entity in HBase.

Not reviewed, for internal circulation only

194 Several tests has been performed: loading Atlas EventIndex data to HBase via Phoenix and
195 Phoenix queries on loaded data. Results are encouraging: Single event picking in 30 ms, 1 full dataset
196 queries in 6-10 seconds. Some basic functions are ready, and further work on performance and user
197 interfaces is ongoing.

198 4. Conclusion

199 The EventIndex project started in 2012 at the end of LHC Run 1 driven by the need of having
200 a functional event picking system for ATLAS data. The data storage and search technology selected in
201 the first phase of the project (Hadoop MapFiles and HBase, in 2013-2014) was the most advanced
202 available at that time in the fast-growing field of BigData and indeed after a couple of initial hiccups it
203 proved reliable and performed satisfactorily. Part of the data are replicated also to Oracle for faster
204 access but mainly to have a uniform environment between event and dataset metadata.

205 Nevertheless, the current implementation of the EventIndex started showing scalability issues
206 as the amount of stored data increases: slower queries, high storage usage (now eased by
207 compression). Also, significant increase in the data rates expected in future LHC runs demands
208 transition to a new technology

209 Phoenix queries and HBase new event table prototypes have been tested and show
210 encouraging results. There is table schema candidate, basic functionality is ready, working towards
211 improved performance and better interfaces. Need to keep testing with more data and get performance
212 metrics. The plan is to have the new system operational by the middle of 2020 in parallel with the old
213 one and phase out the old system at the end of 2020 (well in advance of the start of LHC Run 3)
214 article.

215 References

- 216 [1] ATLAS Collaboration. The ATLAS Experiment at the CERN Large Hadron Collider // Journal of
217 Instrumentation, Volume 3, Issue 08, pp. S08003 (2008) - DOI: 10.1088/1748-0221/3/08/S08003
- 218 [2] Barberis D. *et al.* The ATLAS EventIndex: an event catalogue for experiments collecting large
219 amounts of data // Journal of Physics: Conference Series, Volume 513, Issue 4, article id. 042002
220 (2014) - DOI: 10.1088/1742-6596/513/4/042002
- 221 [3] Barberis D. *et al.* The ATLAS EventIndex: architecture, design choices, deployment and first
222 operation experience // Journal of Physics: Conference Series, Volume 664, Issue 4, article id. 042003
223 (2015) - DOI: 10.1088/1742-6596/664/4/042003
- 224 [4] Bird I. *et al.* LHC Computing Grid. Technical Design Report // CERN-LHCC-2005-024 LCG-
225 TDR-001. - Geneva : CERN, 2005. - 149 p.
226 Available at: <https://cds.cern.ch/record/840543/files/lhcc-2005-024.pdf> (accessed 19.11.2019)
- 227 [5] Albrand S., Fulachier J. and Lambert F. The ATLAS metadata interface // Journal of Physics:
228 Conference Series, Volume 219, Issue 4, article id. 042030 (2010) - DOI:10.1088/1742-
229 6596/219/4/042030
- 230 [6] Fulachier J., Odier J. and Lambert F. ATLAS Metadata Interface (AMI), a generic metadata
231 framework // Journal of Physics: Conference Series, Volume 898, Issue 6, article id. 062001
232 (2017) - DOI: 10.1088/1742-6596/898/6/062001
- 233 [7] Sanchez J., Fernandez Casani J., González de la Hoz S. on behalf of the ATLAS Collaboration.
234 Distributed Data Collection for the ATLAS EventIndex, // Journal of Physics: Conference Series,
235 Volume 664, Issue 4, article id. 042046 (2015) - DOI: 10.1088/1742-6596/664/4/042046

- 236 [8] Fernandez Casani A. *et al.* Distributed Data Collection for the Next Generation ATLAS
237 EventIndex Project // The European Physical Journal Conferences Volume 214, article id. 04010
238 (2019) - DOI: 10.1051/epjconf/201921404010
- 239 [9] K Shvachko *et al.* The Hadoop Distributed File System, // Proceedings of the 2010 IEEE 26th
240 Symposium on Mass Storage Systems and Technologies (MSST), MSST'10, pages 1-10 - DOI:
241 10.1109/MSST.2010.5496972,
- 242 [10] Apache Hadoop. Available at <http://hadoop.apache.org> (accessed 19.11.2019)
- 243 [11] Favareto A. *et al.* Use of the Hadoop structured storage tools for the ATLAS EventIndex event
244 catalogue // Physics of Particles and Nuclei Letters, September 2016, Volume 13, Issue 5, pp
245 621–624 - DOI:10.1134/S1547477116050198
- 246 [12] Apache HBase. Available at: <https://hbase.apache.org> (accessed 19.11.2019)
- 247 [13] Gallas E. *et al.* An Oracle-based EventIndex for ATLAS // Journal of Physics: Conference Series,
248 Volume 898, Issue 4, article id. 042033 (2017) - DOI: 10.1088/1742-6596/898/4/042033
- 249 [14] Barberis D. *et al.* ATLAS EventIndex monitoring system using the Kibana analytics and
250 visualization platform // Journal of Physics: Conference Series, Volume 762, Issue 1, article id.
251 012004 (2016) - DOI: 10.1088/1742-6596/762/1/012004
- 252 [15] Alexandrov E. *et al.* BigData tools for the monitoring of the ATLAS EventIndex // Proceedings
253 of the VIII International Conference "Distributed Computing and Grid-technologies in Science
254 and Education" (GRID 2018), Dubna, Moscow region, Russia, September 10 -14, 2018,
255 Available at: <http://ceur-ws.org/Vol-2267/91-94-paper-15.pdf> (accessed 19.11.2019)
- 256 [16] InfluxDB: Purpose-Built Open Source Time Series Database. Available at:
257 <https://www.influxdata.com> (accessed 19.11.2019)
- 258 [17] Grafana: The open observability platform. Available at: <https://grafana.com> (accessed
259 19.11.2019)
- 260 [18] Baranowski Z. *et al.* A study of data representation in Hadoop to optimize data storage and search
261 performance for the ATLAS EventIndex // Journal of Physics: Conference Series, Volume 898,
262 Issue 6, article id. 062020 (2016) - DOI: /10.1088/1742-6596/898/6/062020
- 263 [19] Fernandez Casani A. *et al.* A prototype for the evolution of ATLAS EventIndex based on Apache
264 Kudu storage // The European Physical Journal Conferences Volume 214, article id. 04057
265 (2019) - DOI: 10.1051/epjconf/201921404057
- 266 [20] Apache Phoenix: <https://phoenix.apache.org> (accessed 19.11.2019)
- 267