

# 1 ATLAS Event Data Organization and I/O Framework

## 2 Capabilities in Support of Heterogeneous Data

### 3 Access and Processing Models

---

**David Malon,<sup>a</sup> Jack Cranshaw\*,<sup>a</sup> Peter van Gemmeren,<sup>a</sup> and Marcin Nowak<sup>b</sup> on behalf of the ATLAS Collaboration<sup>†</sup>**

<sup>a</sup>Argonne National Laboratory

<sup>b</sup>Brookhaven National Laboratory

E-mail: [malon@anl.gov](mailto:malon@anl.gov)

Choices in persistent data models and data organization have significant performance ramifications for data-intensive scientific computing. In experimental high energy physics, organizing file-based event data for efficient per-attribute retrieval may improve the I/O performance of some physics analyses but hamper the performance of processing that requires full-event access. In-file data organization tuned for serial access by a single process may be less suitable for opportunistic sub-file-based processing on distributed computing resources. Unique I/O characteristics of high-performance computing platforms pose additional challenges. The ATLAS experiment at the Large Hadron Collider employs a flexible I/O framework and a suite of tools and techniques for persistent data organization to support an increasingly heterogeneous array of data access and processing models.

*38th International Conference on High Energy Physics*  
*3-10 August 2016*  
*Chicago, USA*

---

\*Speaker.

<sup>†</sup>Argonne National Laboratory's work was supported by the U.S. Department of Energy, Office of Science, Office of Basic Energy Sciences, under contract DE-AC02-06CH11357.



## 4 **1. ATLAS Run 2 event data model**

5 The ATLAS experiment[1] substantially redesigned its event data model for Run 2 of the Large  
6 Hadron Collider, orienting it more directly toward user analysis. Details of the Run 2 data model  
7 have been described elsewhere[2]. Among its key features are:

- 8 • a shared pattern for state definition for event data model objects;
- 9 • optimization for attribute-level retrieval and histogramming;
- 10 • column orientation, aggregating attributes across objects in a collection, with a preference  
11 for structs of vectors over vectors of structs;
- 12 • design with direct mapping of the transient data model to its persistent representation in  
13 ROOT[3] in mind;
- 14 • support for dynamic runtime addition of new attributes and decorations.

15 Such a design is well matched to end-user analysis use cases. There are a few limitations:  
16 event-by-event variation in content types is not supported (e.g., if one event contains a jet collection,  
17 every event must contain a jet collection, even if it is empty), and schema evolution support is  
18 deliberately limited to that provided by ROOT for the sake of simplicity and efficiency in late-stage  
19 analysis.

## 20 **2. ATLAS persistent data model infrastructure**

21 ATLAS employs a highly capable and very general persistent data model infrastructure that  
22 has been described in detail elsewhere [4] [5]. Among its features are:

- 23 • support for direct navigation to and retrieval of arbitrary data objects across supported per-  
24 sistence technologies;
- 25 • uniform reference model for event, sub-event and non-event data, in-file and cross-file refer-  
26 ences with runtime cross-file navigation, and runtime back navigation to upstream data;
- 27 • support for persistent state representation of arbitrarily complex transient data objects at a  
28 level above the choice of persistence technology;
- 29 • support for objects that serve as event entry points, which also record provenance, support  
30 access to upstream data, and allow restoration of the state of the transient event store inde-  
31 pendent of persistence technology.

32 The capabilities of this infrastructure are more extensive and more general than ATLAS has  
33 tended to exploit in practice. In Run 1, for example, while the software provided the possibility  
34 of runtime access to optional additional non-local data while processing an input dataset, such  
35 access would not have been readily supported by then-current production system components and  
36 site configurations. These capabilities, though, are well suited to a world of distributed object  
37 stores, and to an environment in which any data at any level of granularity should be readable from  
38 anywhere via wide-area access protocols.

### 39 **3. Performance: tradeoffs and tuning**

40 While the ATLAS Run 2 event data model is well suited to direct user analysis, it is less tailored  
41 to full-event processing (reconstruction, input to the derivation framework that produces ATLAS  
42 analysis data products, and so on) than to selective content retrieval. Among other considerations, a  
43 side effect of the model is substantial memory consumption for writing and reading when individual  
44 attributes have first-class status in the persistent data model. This effect is decidedly non-trivial  
45 when full events are processed, and would without optimization cause ATLAS reconstruction and  
46 other processing to exceed the available memory on many of the computing resources upon which  
47 it relies. While adding attributes and decorations as "dynamic" is appealing to users, such additions  
48 can exacerbate the problem, as each such attribute is stored in an output ROOT file in such a way  
49 that it contributes its own nonnegligible buffer and memory footprint consequences. Because of  
50 the requirement that all events have the same content objects, the value of such dynamic attributes  
51 may be diminished in practice in any case, as equivalent (albeit null) attributes must be added by  
52 back-filling or by other means to the events that lack them.

53 Performance tuning must balance such trade-offs. Examples of tunable parameters include  
54 buffer and (ROOT) basket sizes, commit intervals, and buffer flush settings. A program of careful  
55 measurement for a variety of use cases is required, and has been undertaken under various local  
56 and remote data access scenarios for the most important production workflows. Reordering of  
57 data within files can also help, and can be optimized for specific use cases. Efficient aggregation  
58 and de-aggregation of attributes in transient-persistent conversions (an area of ongoing work) also  
59 promises the potential for substantial savings in storage footprint.

### 60 **4. Emerging workflows**

61 Opportunistically-available resources are playing an increasingly important role in ATLAS  
62 computing. Efficient use requires a scatter-gather architecture capable of delivering one or a  
63 few events rather than full files of events to ephemerally-available resources. The ATLAS event  
64 service[6] supports such a fine- and variable-grained model. Feeding the very large numbers of  
65 processors on high-performance computing (HPC) platforms is another increasingly important use  
66 case.

67 I/O components have been successfully adapted to support the ATLAS event service model,  
68 with concomitant support for multi-process worker jobs [7]. There remains fertile ground for opti-  
69 mization, though, both in persistent data layout and in I/O components that support such process-  
70 ing. As the current production mapping of the ATLAS Run 2 data model to ROOT persistence  
71 aggregates objects of the same type across events, care must be taken by I/O components to avoid  
72 substantial inefficiencies. These may arise, for example, when different processes handle different  
73 events that have been compressed together in the same sets of buffers. A simple alternative persis-  
74 tent data model strategy to support wide-area event-by-event data distribution might store all data  
75 for a single event in a single contiguous block of bytes, making the process of sending one event  
76 to one processor and another event to another relatively straightforward. This contiguous block of  
77 bytes strategy is in fact employed in the so-called "bytestream" files written by the trigger farm  
78 for offline reconstruction. A disadvantage to such an approach, though, is a larger storage footprint

79 because of reduced compression (no compression across events). At LHC data volumes and given  
80 collaboration storage resource constraints, such a disadvantage may or may not be decisive. For  
81 processing that requires access only to a handful of event data attributes, there could be further dis-  
82 advantages (e.g., reading unneeded data) that must be balanced with the use cases that read entire  
83 events.

84 These and other tradeoffs serve to emphasize the advantages of a flexible persistent data model  
85 architecture as well as a need for an ongoing program of optimization. Current ATLAS production  
86 does not adjust or tune the persistent data model differently to support event service versus more  
87 traditional workflows, or to adapt to an object store versus a more standard file system as an output  
88 destination, but it could.

## 89 5. Conclusions and next steps

90 Performance tuning in ATLAS is an ongoing activity, not only for I/O and persistence, but for  
91 almost all software used in large-scale event processing. The approach to I/O and persistence tun-  
92 ing must balance analysis and production use cases, wide area and local access, and standard and  
93 emerging workflows. Tuning choices may change as use case and workflow balances shift. Among  
94 strategies under investigation is a storage-chunk-aware event streaming service, which could dis-  
95 tribute parcels of events matched to how events are bundled in the input file (as determined by  
96 flush settings and commit intervals and other I/O options). Among the challenges is to implement  
97 such a capability in a way that keeps framework components as technology-independent as possi-  
98 ble so that ATLAS may benefit, when appropriate, from serialization and persistence technology  
99 developments both within and without the HEP software community.

## 100 References

- 101 [1] ATLAS Collaboration 2008, *The ATLAS Experiment at the CERN Large Hadron Collider, JINST 3*  
102 **S08003** [doi:10.1088/1748-0221/3/08/S08003].
- 103 [2] A. Buckley et al., *Implementation of the ATLAS Run 2 event data model, J. Phys.: Conf. Ser. 664*  
104 (2015) no.7, 072045 [doi:10.1088/1742-6596/664/7/072045].
- 105 [3] R. Brun et al., *ROOT: An object oriented data analysis framework, Nucl.Instrum.Meth. A389* (1997)  
106 81-86
- 107 [4] P. van Gemmeren et al., *Next-Generation Navigational Infrastructure and the ATLAS Event Store, J.*  
108 *Phys. Conf. Ser. 513*, 052036 (2014) [doi:10.1088/1742-6596/513/5/052036].
- 109 [5] P. van Gemmeren and D. Malon, *Persistent Data Layout and Infrastructure for Efficient Selective*  
110 *Retrieval of Event Data in ATLAS*, Proceedings of the DPF-2011 Conference, Providence, RI, August  
111 8-13, 2011 [physics.data-an/1109.3119].
- 112 [6] P. Calafiura et al., *The ATLAS Event Service: A new approach to event processing, J. Phys. Conf. Ser.*  
113 **664**, no. 6, 062065 (2015). [doi:10.1088/1742-6596/664/6/062065].
- 114 [7] P. van Gemmeren et al., *I/O strategies for multicore processing in ATLAS, J. Phys. Conf. Ser. 396,*  
115 *022054* (2012), [doi:10.1088/1742-6596/396/2/022054].