

1 Towards a new conditions data infrastructure in ATLAS

2 *Evgeny Alexandrov*¹, *Luca Canali*², *Davide Costanzo*^{3,*}, *Andrea Formica*^{4,**}, *Elizabeth*
3 *J.Gallas*^{5,***}, *Mikhail Mineev*¹, *Nurcan Ozturk*^{6,****}, *Shaun Roe*², *Vakho Tsulaia*⁷, and
4 *Marcelo Vogel*⁶

5 ¹Joint Institute for Nuclear Research, Joliot-Curie 6, 141980 Dubna (Russia)

6 ²CERN, CH-1211 Geneva 23 (Switzerland)

7 ³Department of Physics and Astronomy, University of Sheffield, Sheffield (United Kingdom)

8 ⁴IRFU, CEA, Université Paris-Saclay, F-91191 Gif-sur-Yvette (France)

9 ⁵University of Oxford, Denys Wilkinson Bldg, Keble Rd, Oxford OX1 3RH (United Kingdom)

10 ⁶University of Texas at Arlington, 701 South Nedderman Drive, Arlington, TX 76019 (USA)

11 ⁷Lawrence Berkeley National Laboratory, Berkeley, CA 94720 (USA)

12 **Abstract.** The ATLAS experiment is preparing a major change in the condi-
13 tions data infrastructure in view of LHC Run 4. In this paper we describe the
14 ongoing changes in the database architecture which have been implemented for
15 Run 3, and describe the motivations and the on-going developments for the
16 deployment of a new system (called CREST for Conditions Representational
17 State Transfer, as a reference to REST architectures). The main goal is to set up
18 a parallel infrastructure for full scale testing before the end of Run 3.

19 1 Introduction

20 The processing of ATLAS [1] event data necessitates the retrieval of a collection of auxiliary
21 non-event data stored within database systems. This data, referred to as "conditions data,"
22 generally exhibits variations over time and encompasses elements such as detector alignment,
23 calibration, and configuration information. The complexity escalates due to the requirement
24 of disseminating this information across the global ATLAS computing grid, along with the
25 sheer multitude of concurrently operating processes on the grid. Each process demands a
26 distinct set of conditions to advance.

27 Our focus is directed towards the foundational database infrastructure, which underwent
28 a redesign for ATLAS in Run 3. This redesign involved the consolidation of resources within
29 the online Oracle cluster, coupled with the necessary developments to ensure secure access.

30 This reorganization resulted in an architecture resembling the one the experiment is
31 preparing for Run 4, known as the CREST project [2]. Here we expound upon the archi-
32 tecture and the current development status of this project. The first significant milestone
33 involves deploying a functional demonstrator by fall 2023, with the intention of testing seg-
34 ments of data processing workflows using real conditions data migrated from the existing
35 system.

*e-mail: davide.costanzo@cern.ch

**e-mail: andrea.formica@cern.ch

***e-mail: elizabeth.gallas@cern.ch

****e-mail: nurcan.ozturk@cern.ch



2 Conditions Database infrastructure

The infrastructure for managing and accessing condition data within ATLAS consists of the following components:

- Database clusters: Oracle databases store the conditions data according to the LCG Conditions database infrastructure [3] which includes C++ and python methods within its COOL API for managing database content and for client read-only access. Over the years using this infrastructure, further methods have been developed on top of the COOL API to suit ATLAS-specific requirements.
- Database read-only copies: Database replicas kept in sync with the source utilizing Oracle technologies which include Active Data Guard (ADG, the Oracle-provided technology for physical replication) and Golden Gate (an Oracle solution for logical replication of selected schemas).
- Generic database access via a middle-tier server: This is achieved through the Frontier system [4] which mediates client data selection requests with the underlying database storage system. Benefits of this layer include the ability to monitor requests as well as to moderate intermittent spikes in load.
- Web Proxy: To ensure efficient access and optimal resource utilization, a series of Squid proxies are deployed. These proxies screen the Frontier server and the database, filtering requests from individual clients (jobs) that are accessing condition data.

2.1 Architecture before Run 3

The conditions data are organized in two different Oracle clusters, depending on their usage (online data taking and High Level Trigger processing, or any other "offline" workflow):

1. Online Oracle cluster (named ATONR) on the ATLAS technical network, for conditions data to be consumed in real time workflows.
2. Offline Oracle cluster (named ATLR) on the CERN GPN network, for conditions data to be consumed in offline workflows, from bulk processing to reprocessing and Monte Carlo simulation.

Conditions that are primarily stored in the online cluster are additionally replicated using Golden Gate streaming technology to the offline cluster. Once in the offline cluster, these conditions can be accessed in read-only mode. The diagram 1 provides a simplified overview of the ATLAS Conditions Data infrastructure prior to Run 3.

2.2 Architecture during Run 3

For Run 3 a major operation of databases consolidation has been prepared. Two main aspects were considered in this plan:

- Oracle license model: the license model until 2023 was covering all Tier-1s Oracle nodes used by ATLAS to keep a copy (via Golden Gate) of the conditions data. After April 2023 a new license model was adopted by CERN IT, based on a "per-core" license for Oracle nodes.
- ATLAS was the only client inside CERN of the Oracle Golden Gate replication technology: its usage was discouraged by CERN IT, considering the licensing costs, the additional load on the support team, and the redundancy with the other available replication technology: Oracle Data Guard.

78 Subsequently, the ATLAS Database and Metadata (ADAM) team made the decision to shift
 79 real data conditions workflows from the offline to the online cluster in order to phase out the
 80 need for the Golden Gate replication. The architecture that was implemented is depicted in
 81 figure 2. A few major alterations in the architecture were required to eliminate the Oracle
 82 Golden Gate replication (from the online cluster to the offline cluster and Tier-1s): the centralization
 83 of all real data conditions into the Oracle ATONR cluster and the establishment of
 84 an intermediary service (called COOL-Proxy) designed to manage user requests for the storage
 85 of new conditions data. In this freshly devised infrastructure, all conditions data situated
 86 outside the online ATLAS network (ATCN) can be accessed solely through read-only (ADG)
 87 replicas of the ATONR nodes.

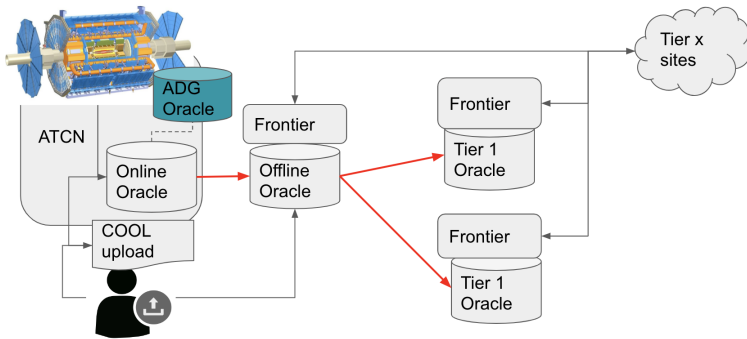


Figure 1. Conditions Data system architecture before Run 3; the red arrows depict data copies using the Oracle Golden Gate technology.

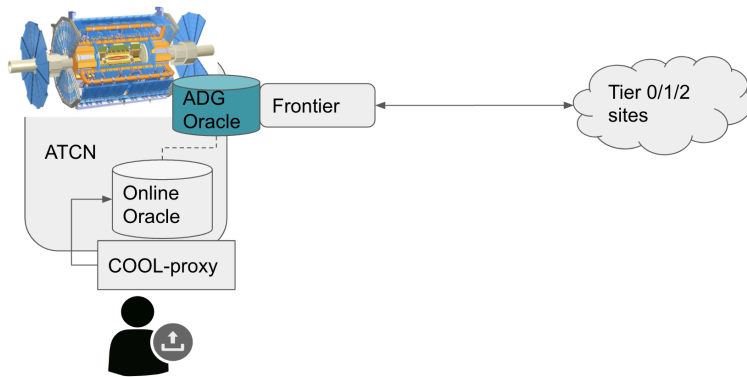


Figure 2. Conditions Data system architecture for Run 3; the COOL-proxy is accessible from CERN General Public Network (GPN).

88 This novel architecture brings about the ability to optimize Oracle licensing costs.
 89 Through the consolidation of all conditions data usage within the online cluster (ATONR),
 90 a distinct and dedicated environment for all data processing workflows is established. This
 91 separation prevents the mingling of conditions usage with other applications, as the ATLR
 92 cluster supports a wide array of applications ranging from detector construction to authorship

93 and metadata. Additionally, this approach has led to a reduction in Oracle administration
94 burden and associated expenses for Tier-1s.

95 To ensure secure access to the ATONR cluster from the CERN GPN network, we devised
96 an intermediary server that acts as a custom proxy system. This system, known as COOL-
97 Proxy, is intended for use by experts responsible for uploading new conditions data. Notably,
98 the COOL-Proxy system employs the new CERN SSO authentication mechanism. Conditions
99 data experts are linked to specific e-groups, and a token mechanism is utilized to grant
100 them writing privileges solely within the Oracle schema corresponding to the e-group(s) to
101 which they belong.

102 On the hardware front, enhancements were made to the online cluster, involving the ad-
103 dition of an extra node. This improved configuration ensures extra capacity to critical online
104 system to process all new conditions upload workflows, while maintaining the same load on
105 the rest of the online environment as experienced during Run 2. Moreover, this extra node
106 improves the available redundancy in case of cluster node failures.

107 **3 CREST: a Conditions Database infrastructure after Run 3**

108 The system implemented during Run 3 closely resembles the architecture that is being tested
109 for the ATLAS runs commencing from Run 4. The new initiative for managing conditions
110 data is named CREST, originating as a progression from the CMS conditions database. It in-
111 herits fundamental concepts for the data model and the design of relational tables from its pre-
112 decessor. The development of the CREST system also benefited from discussions within the
113 HEP Software Foundation [5] working group on cross-experiment conditions data manage-
114 ment systems [6]. The intention behind CREST is to replace the existing COOL conditions
115 database in satisfying the conditions data requests of all offline data processing and Monte
116 Carlo simulations from Athena [7] jobs (Athena is the the ATLAS software framework for
117 event processing). This comprehensive system is composed of several integral components:

- 118 • **Relational Database:** Data in the CREST database is stored in relational tables utilizing a
119 straightforward schema. Conditions data payloads are stored within the database as Large
120 Objects (LOBs) and referenced through unique keys stored as related metadata in CREST.
121 In-depth information regarding the data model can be found in references [2] and [8].
- 122 • **A REST API for the conditions data management system,** accompanied by an implemented
123 web server and corresponding client libraries.
- 124 • **A web proxy system designed to offer a caching layer,** thereby diminishing the utilization
125 of the web server and database by clients (Athena jobs).

126 While the architecture closely mirrors that of Run 3, there are notable improvements (refer
127 to figure 3). Apart from a significant simplification in the data model, resulting in a substantial
128 reduction in the number of tables, the CREST system introduces a REST API for conditions
129 data management. This innovation entirely decouples client code from the underlying storage
130 implementation. Consequently, clients are no longer obligated to understand how the storage
131 system is internally structured.

132 The most significant distinction lies in the capability to maximize the utilization of the
133 caching layer by establishing a clear demarcation between metadata, such as the validity
134 intervals for each individual conditions data payload, and the conditions payload itself. This
135 segregation is achieved at the level of the REST API definition through access to identical
136 conditions data sets via a unique *key*.

137 To grasp the benefits of such a data model design, we can examine the current utilization
138 of database servers in both ATLAS and CMS. We assume that both experiments possess a
139 comparable amount of conditions data.

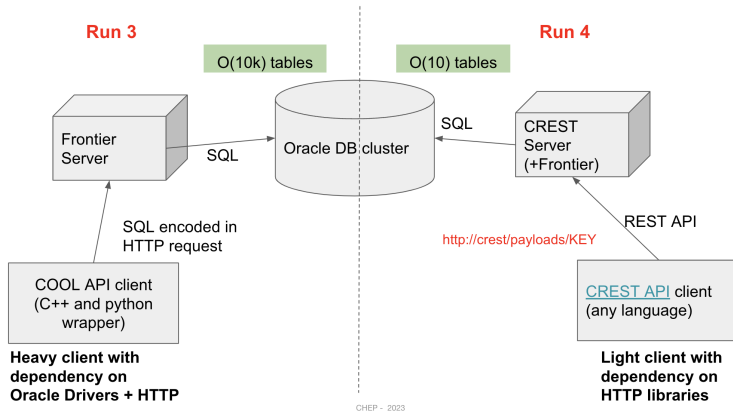


Figure 3. Conditions Data: comparison between present and proposed architectures.

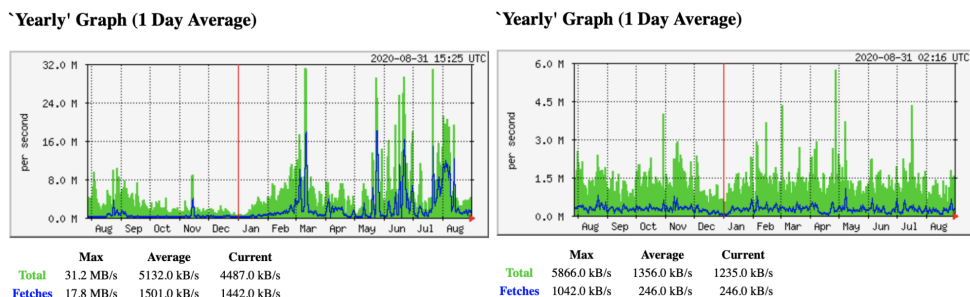


Figure 4. Frontier/Squid usage in ATLAS and CMS.

140 We have extracted the volumes of data retrieved from Oracle ("Fetches") and from the
 141 Squid system ("Total") for both the ATLAS and CMS experiments, utilizing official monitoring
 142 plots. These figures are presented in table 1. Additionally, the yearly (single day average)
 143 graph is displayed in figure 4. The considerable variability observed in the case of ATLAS
 144 suggests a less optimal utilization of the caching system, likely stemming from the manner
 145 in which clients request the necessary data. A more comprehensive investigation into the
 146 underlying causes of these inefficiencies has been conducted within ATLAS, leveraging the
 147 logging data from the Frontier servers [9].

Table 1. ATLAS and CMS Frontier/Squid monitoring.

Experiment	Type	Fetches (MB/s)	Total (MB/s)	Ratio
ATLAS	Year Avg	1.5	5.1	30%
ATLAS	Year Max	17	31	54%
CMS	Year Avg	0.25	1.35	20%
CMS	Year Max	1	5.8	17%

148 **4 CREST architecture**

149 The CREST system follows a multi-tier model architecture. In this arrangement, the back-
150 end remains a relational database that employs a concise collection of tables to oversee the
151 management of conditions data metadata and payloads. Simultaneously, a web server is
152 fashioned as the front-end. This web server actualizes a REST API, abstracting the direct
153 interaction with the database. A collection of client libraries has been prepared to facilitate
154 the utilization of the REST API from various programming languages. Notably, a C++ client
155 has been meticulously developed for usage from Athena clients. The existing state of the
156 system is elaborated upon in this section.

157 **4.1 CREST REST API**

158 The REST API is documented using OpenAPI specifications [10] in a YAML format. This
159 API essentially outlines the URL paths made accessible through the CREST server, as well
160 as the data objects exchanged between the server and the client (JSON is employed for data
161 transmitted via HTTP). Opting for a standardized set of specifications for API description
162 offers the advantage of being compatible with a diverse range of tools, enabling the gener-
163 ation of code for both server stubs and clients across various programming languages and
164 frameworks.

165 The API description encompasses a comprehensive array of metadata elements essential
166 for conditions data management, including tags, intervals of validity (IOVs), and global tags.
167 Remarkably, these metadata components exhibit high similarity between the current ATLAS
168 system (COOL) and the CMS data model.

169 **4.2 CREST server**

170 The CREST server is constructed using established Java technologies [11], specifically rely-
171 ing on specifications like JAX-RS and JPA, alongside the Spring Boot [12] framework.

172 The description of the REST API via OpenAPI facilitates the generation of server stubs
173 within the Jersey framework [13], employing standard generation tools [14].

174 The robust support and seamless interoperability within the Java ecosystem contribute to
175 the stability of the server code over time. This ecosystem's flexibility allows for effortless
176 transitions between various sets of implementation libraries. For instance, a switch between
177 web servers such as Tomcat [15] and Undertow [16] can be accomplished through a simple
178 adjustment in the CREST server's build file, without necessitating internal code alterations.

179 Utilizing JPA implementations, such as Hibernate [17], for standardized database access
180 provides the advantage of concise object-relational mapping syntax, while still retaining the
181 option for deeper optimization of specific queries.

182 The project's source code (for the server and the related libraries) is hosted on GitLab at
183 CERN¹.

184 **4.3 CREST client libraries and tools**

185 Interactions with the CREST server occur through utilization of the REST API. The official
186 client takes the form of a C++ implementation, and it is integrated into Athena conditions
187 data services.

188 To assess the functionality of the current software for a specific subsystem, we can carry
189 out trials by migrating the conditions data for that given subsystem into the CREST database.

¹<https://gitlab.cern.ch/crest-db>

190 To facilitate this migration from the existing COOL DB, a dedicated tool has been crafted.
 191 This converter tool can be configured to selectively copy "tags" from COOL to CREST. This
 192 operation involves employing the COOL API for reading and the CREST C++ client for
 193 data insertion via the CREST server into Oracle. The tool also provides a set of logging
 194 information, offering the added benefits of profiling and debugging the copying process.

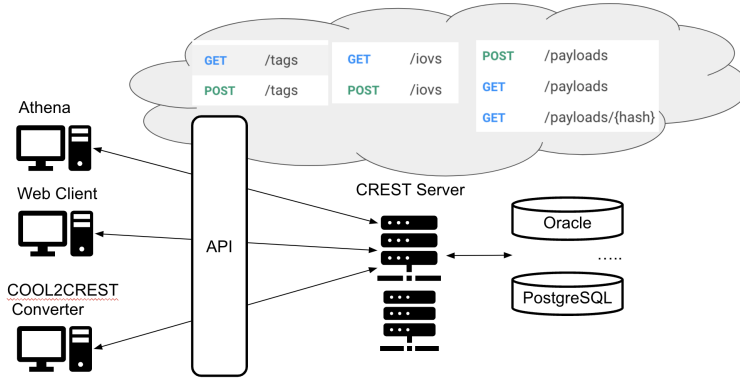


Figure 5. CREST deployment and test infrastructure

195 4.4 CREST deployment and test infrastructure

196 To assess the functionality of the CREST system, we have employed a cloud-based deployment
 197 approach. A range of distinct CREST servers, each employing different underlying
 198 schemes, are accessible to developers and Athena testers. This setup allows us to experiment
 199 with various database platforms, including Oracle and Postgres, ensuring that the server code
 200 remains well separated from the specifics of the underlying storage technology. A depiction
 201 of this deployment scheme is presented in figure 5.

202 As of now, the official deployment resides in a Kubernetes [18] cluster, utilizing machines
 203 within the CERN openstack infrastructure. This cluster also serves to deploy a caching system
 204 that relies on Varnish. This caching system plays a crucial role in validating the architecture
 205 during our initial large-scale tests.

206 ATLAS has laid out plans to introduce a demonstrator for CREST utilization by the conclusion
 207 of 2023. This demonstrator is set to undergo testing within the High-Level Trigger (HLT)
 208 workflow. This choice stems from the fact that online workflows, such as the HLT,
 209 impose more demanding requirements in terms of caching. Notably, conditions data like luminosity
 210 and beam-spot need to be refreshed regularly, sometimes even at the granularity of
 211 each luminosity block ².

212 5 Conclusions

213 We have detailed the modifications made to the ATLAS conditions data management infras-
 214 tructure in preparation for the Run 3 data acquisition phase. These adjustments are geared
 215 towards readying both the experiment and conditions data users for an enhanced architecture
 216 set to be employed in the upcoming data acquisition (during Run 4). Additionally, we have

²A luminosity block is defined as a period with stable luminosity (generally about one minute in duration).

217 elucidated the novel architecture known as the CREST project and highlighted its distinctions
218 from the current system. We have emphasized the core differences and enhancements that the
219 new architecture aims to tackle.

220 References

- 221 [1] ATLAS Collaboration, JINST **3**, S08003 (2008), [https://dx.doi.org/10.1088/](https://dx.doi.org/10.1088/1748-0221/3/08/S08003)
222 1748-0221/3/08/S08003
- 223 [2] P.J. Laycock, D. Dykstra, A. Formica, G. Govi, A. Pfeiffer, S. Roe, R. Sipos, Journal of
224 Physics: Conference Series **1085**, 032040 (2018), [https://dx.doi.org/10.1088/](https://dx.doi.org/10.1088/1742-6596/1085/3/032040)
225 1742-6596/1085/3/032040
- 226 [3] A. Valassi, R. Basset, M. Clemencic, G. Pucciani, S.A. Schmidt, M. Wache, *COOL,*
227 *LCG conditions database for the LHC experiments: Development and deployment status,*
228 *in IEEE Nuclear Science Symposium Conference Record, 2008. NSS '08* (2008),
229 pp. 3021–3028
- 230 [4] D. Dykstra, J. Phys.: Conf. Ser. **331**, 042008 (2011), [http://iopscience.iop.org/](http://iopscience.iop.org/1742-6596/331/4/042008)
231 1742-6596/331/4/042008
- 232 [5] *HEP Software Foundation*, <https://hepsoftwarefoundation.org/>
- 233 [6] M. Bracko, M. Clemencic, D. Dykstra, A. Formica, G. Govi, M. Jouvin, D. Lange,
234 P. Laycock, L. Wood, TBD (2019), <https://www.osti.gov/biblio/1527431>
- 235 [7] G.A. Stewart, et al., J.Phys.Conf.Ser. **762**, 012024 (2016), [https://iopscience.](https://iopscience.iop.org/article/10.1088/1742-6596/762/1/012024)
236 [iop.org/article/10.1088/1742-6596/762/1/012024](https://iopscience.iop.org/article/10.1088/1742-6596/762/1/012024)
- 237 [8] L. Rinaldi, A. Formica, E.J. Gallas, N. Ozturk, S. Roe, EPJ Web Conf. **214**, 04052
238 (2019), <https://doi.org/10.1051/epjconf/201921404052>
- 239 [9] A. Formica, N. Ozturk, M. Si Amer, J.L. Bahilo, E.J. Gallas, I. Vukotic, EPJ Web Conf.
240 **245**, 04032 (2020), <https://doi.org/10.1051/epjconf/202024504032>
- 241 [10] *OpenAPI v3.1.0*, <https://spec.openapis.org/oas/latest.html>
- 242 [11] *Java EE 8*, <https://javaee.github.io/javaee-spec/>
- 243 [12] *Spring Boot documentation*, [https://docs.spring.io/spring-boot/docs/](https://docs.spring.io/spring-boot/docs/current/reference/htmlsingle/)
244 [current/reference/htmlsingle/](https://docs.spring.io/spring-boot/docs/current/reference/htmlsingle/)
- 245 [13] *Jersey 2 JAX-RS API implementation*, <https://eclipse-ee4j.github.io/jersey>
- 246 [14] *OpenApi tools and code generation*, [https://github.com/OpenAPITools/](https://github.com/OpenAPITools/openapi-generator)
247 [openapi-generator](https://github.com/OpenAPITools/openapi-generator)
- 248 [15] *Tomcat Web server*, <https://tomcat.apache.org/>
- 249 [16] *Undertow Web server*, <https://undertow.io/>
- 250 [17] *Hibernate ORM library*, <https://hibernate.org>
- 251 [18] *Kubernetes*, <https://kubernetes.io>