

1 The new inter process communication middle-ware for the 2 ATLAS Trigger and Data Acquisition system

3 **Serguei Kolos¹ and Reiner Hauser² on behalf of the ATLAS TDAQ Collaboration**

- 4 1. University of California Irvine, 4172 Frederick Reines Hall, Irvine, CA, USA.
5 2. Michigan State University, Est. 1855. East Lansing, Michigan, USA.
6 serguei.kolos@cern.ch

7 **Abstract.** The ATLAS Trigger & Data Acquisition (TDAQ) project was started almost twenty
8 years ago with the aim of providing scalable distributed data collection system for the
9 experiment. While the software dealing with physics data flow was implemented by directly
10 using the low-level communication protocols, like TCP and UDP, the control and monitoring
11 infrastructure services for the TDAQ system were implemented on top of the CORBA
12 communication middle-ware. CORBA provides a high-level object oriented abstraction for the
13 inter process communication, hiding communication complexity from the developers. This
14 approach speeds up and simplifies development of communication services but incurs some
15 extra cost in terms of performance and resources overhead. Our experience of using CORBA
16 for control and monitoring data exchange in the distributed TDAQ system was very successful,
17 mostly due to the outstanding quality of the CORBA brokers, which have been used in the
18 project: omniORB for C++ and JacORB for Java. However, due to a number of shortcomings
19 and technical issues the CORBA standard has been gradually losing its initial popularity in the
20 last decade and the long term support for the open source implementations of CORBA
21 becomes questionable. Taking into account the time scale of the ATLAS experiment, which
22 goes beyond the next two decades, the TDAQ infrastructure team reviewed the requirements
23 for the inter process communication middle-ware and performed the survey of the
24 communication software market in order to access the modern technologies which raised in the
25 past years. Based on the result of that survey several technologies were evaluated for
26 estimating the long-term benefits and drawbacks of using them as a possible replacement for
27 CORBA during the next long LHC shutdown, which is scheduled in 2 years from now. The
28 evaluation concluded recently with the recommendation of using communication library called
29 ZeroMQ in place of CORBA. The article presents the methodology and the results of the
30 evaluation as well as the plans of organizing the migration from CORBA to ZeroMQ.

31 1. Introduction

32 The TDAQ [1] online system of the ATLAS [2] experiment is composed of tens of thousands of
33 software processes distributed over several thousand computers. For the system to function properly
34 all of these processes must be operated in a coherent way, thus making Inter-Process Communication
35 (IPC) a crucial task. The current implementation of the TDAQ control system, which was born in
36 1998, is based on the CORBA [3] communication middleware. Two CORBA implementations have
37 been used: JacORB [4] for Java and omniORB [5] for C++. They both satisfied the performance and
38 scalability requirements and simplified development and maintenance of the TDAQ software.
39 However, after more than 10 years of successful experience with the CORBA software, we have
40 decided that the time is right to explore if there are new products on the IPC software market which
41 can improve our system performance and maintainability.



42 **2. CORBA in the light of modern software practices**

43 CORBA is an open standard for distributed object computing, which was proposed in 1991 by the
44 Object Management Group (OMG). This was the first attempt to provide a broad high-level standard
45 for information exchange in a distributed software environment. The standard was quite successful and
46 played an important role in the overall evolution of distributed software systems. However many key
47 features of the CORBA standard have a number of built-in drawbacks, which have become more and
48 more prominent in recent years, making CORBA less attractive for modern software development.

49 *2.1. The Interface Definition Language*

50 CORBA proposed a dedicated language called Interface Definition Language (IDL) for
51 communication protocol description. The code for a specific programming language can be
52 automatically generated from such a description. This approach provided a powerful yet simple
53 solution for establishing communication between different programming languages and operating
54 systems. While IDL was originally one of the strongest points of the CORBA standard, the passage of
55 time has seen it become one of the weakest. The issue was that the standard is very strict with defining
56 the mapping from the IDL declarations to the programming languages, thus practically excluding any
57 opportunity to benefit from language evolution. While the most popular programming languages, like
58 C++ and Java, evolved significantly in the last decade the CORBA API couldn't benefit from that so
59 now the CORBA API looks archaic and awkward. The C++ IDL mapping efficiency also suffers from
60 the absence of some recently introduced features like zero-copy or move-assignment semantics.

61 *2.2. CORBA brokers interoperability*

62 Another strong point of the CORBA standard was interoperability between different CORBA
63 implementations, which is required by the standard. This is a good feature, which unfortunately led to
64 some issues with communication efficiency as all CORBA implementations were forced to use the
65 same communication protocol, called IIOP, to support interoperability. The problem is that IIOP has
66 some drawbacks, like for example the data size and processing time overhead due to the data
67 alignment requirements.

68 *2.3. The source code compatibility*

69 The standard precisely defines the API for any possible operation including object creation,
70 registration, activation and so on for assuring source code compatibility between different CORBA
71 implementations. This requirement forced any CORBA broker to provide a high-level object-oriented
72 API, which completely hides all aspects of the underlying communication. This of course simplifies
73 software development and maintenance but at the same time adds a noticeable performance overhead
74 and reduces the flexibility of the communication implementation. In practice that would also make any
75 end user CORBA applications strongly dependent on the quality of the chosen CORBA
76 implementation, thus making problematic a transparent migration from one CORBA broker to another.

77 **3. The modular communication architecture**

78 The rapid evolution of programming languages and the increasing popularity of the open source
79 software development model drastically changed the landscape of the IPC software domain in the last
80 decade. The modern open source market offers a large variety of high quality software packages,
81 which can be used for implementing the basic components of a high-level communication system. A
82 combination of such packages is an attractive alternative to a monolithic heavyweight communication
83 system like CORBA. Such a solution implies that the communication system is organized into a
84 hierarchy of software layers with well-defined interfaces between them, which allow changes to the
85 implementation for these layers in a transparent way for end user applications. Another important
86 property of this design is the transparency of the layers, which means that the APIs of all layers are
87 exposed to the end users and can be used independently of each other. This gives the full advantage of
88 using the high-level API for a simple communication implementation while is still leaving open the
89 possibility of implementing performance-critical applications using the low-level communication API
90 to increase efficiency.

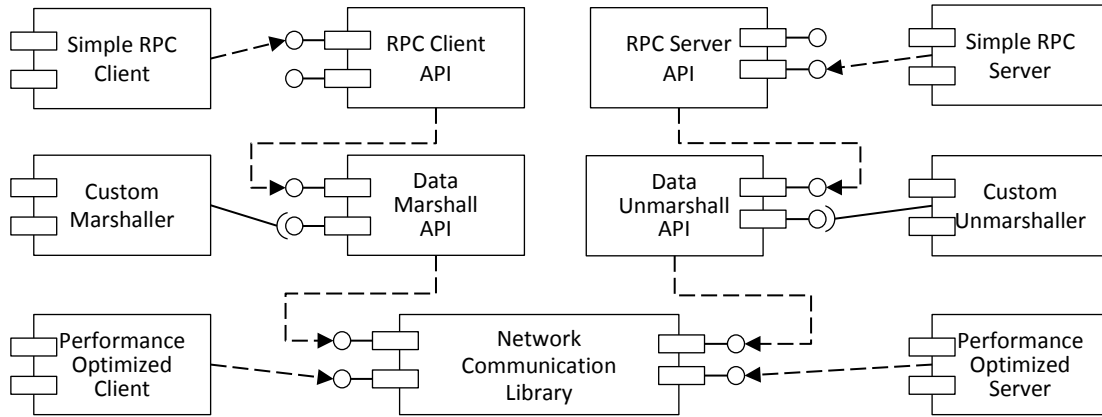


Figure 1. Modular communication architecture for the new ATLAS TDAQ IPC software.

91

92 Figure 1 shows the architecture of the new IPC software for the ATLAS TDAQ system, which has
 93 been designed using this approach. This software consists of three main layers:

- 94
- 95 1. A low-level communication library, which provides a simple and efficient API for exchanging
 96 unstructured data, i.e. messages, between software applications independently of their
 97 location. This component provides abstraction for the network communication layer.
 - 98 2. The data marshalling and unmarshalling components provide a way of passing structured
 99 information between applications. These components define the API for the conversion of an
 100 arbitrary data structure into a message, which can be passed over a network and vice versa.
 101 This API can have multiple implementations, which may be interchanged transparently for the
 102 end users. The only limitation is that any implementation of this API has to support all
 103 programming languages used in the ATLAS TDAQ system. A user application can use these
 104 components directly to convert arbitrary information to the network specific format before
 105 passing it to another application using the first layer API.
 - 106 3. The Remote Procedure Call (RPC) layer provides the top level API for inter-process
 107 communication, making remote calls look like the normal local ones in a given programming
 108 language. This is the top-most API layer, which is simple to use but incurs a relatively high
 109 overhead with respect to the lower layers due to the generic code, which is capable of
 transparently mapping an arbitrary user function to an RPC procedure.

110 **4. The new IPC software implementation**

111 To simplify implementation and maintenance of the new IPC software we decided to use the existing
 112 open software projects as much as possible. As a result the fully functional IPC implementation has
 113 only a few hundred lines of custom code, which was provided mostly for the RPC layer
 114 implementation, while the implementations of the first two layers are almost entirely based on external
 115 software.

116 *4.1. Using ZeroMQ as the network transport layer*

117 After comparing a number of communication libraries available on the open software market we have
 118 chosen ZeroMQ [6] as the implementation of the network communication layer. ZeroMQ is a low-
 119 level C-style library for I/O based on an object resembling a standard socket, but which hides all the
 120 state management and error handling complexity one would normally expect. ZeroMQ supports most
 121 of the widely used programming languages including C, C++, C#, Java, Python, Ruby, and many
 122 others. ZeroMQ is open source software with large and very active user and developer communities.
 123 ZeroMQ provides a simple yet mature API for remote communication, offering the full power of a
 124 classic socket, becoming a de facto standard for many modern communication software projects. In
 125 order to support some common services used in the TDAQ system, some extra classes have been

126 added on top of the ZeroMQ API, but in agreement with our design approach they don't imply any
127 limitations to the direct use of the ZeroMQ native classes.

128 4.2. Data serialization

129 Serialization is an important ingredient of a software communication system, which has a major
130 impact on its performance. For the sake of flexibility we have defined a simple interface for
131 *marshalling*, i.e. converting a programming language structure to the sequence of bytes and
132 *unmarshalling*, which does the opposite operation.

```
public interface ISerializer {
    public class Request {
        private static final AtomicLong gid = new AtomicLong(0);

        public final Method method;
        public final List<Object> params;
        public final long id;
    }

    String serializeRequest(Request req);

    Request deserializeRequest(String str, Map<String, Method> methods);

    String serializeResponse(Object resp, Exception e, long id);

    <T> T deserializeResponse(String str, Class<T> clazz) throws Exception;
}
```

Figure 2. Example IPC API serialisation interface (Java).

133 Figure 2 shows how such an interface is declared for the IPC API in Java. The interface has four
134 methods: two for marshalling and unmarshalling requests on the client side and another two doing the
135 same operations for a server response. The default implementation of this interface uses the Google
136 Gson [7] library to pass the data over a network in Json format. The default implementation works fine
137 for simple requests, which don't carry too much structured data. On the other hand for performance-
138 optimized applications one can provide another implementation of this interface, which can be easily
139 plugged into the IPC library and used transparently without affecting the code of the communication
140 applications.

141 4.3. The RPC API

142 The TDAQ online software system consists of a number of communication services providing specific
143 APIs for exchanging different types of information. To simplify development and maintenance of such
144 services we tried to implement the RPC communication layer for the IPC software in such a way that a
145 remote request would look as much as possible like a local one. It wasn't possible to achieve this goal
146 completely in the C++ API, but in Java this was successfully implemented using the Java Reflection
147 API [8].

148 Following the standard Java approach a user shall first declare a new interface with all methods, which
149 can be called remotely. Figure 3 shows a simple example of such an interface.

```
public interface IHello {
    void say(String greetings);
}
```

Figure 3. A simple custom communication interface (Java).

150

151 Figure 4 demonstrates how a simple IPC Java server can be implemented. The *HelloImpl* class
152 implements the *IHello* interface in the same way as a normal Java interface. Then the instance of the

153 *HelloImpl* class is given as a parameter to the IPC Server class constructor, which is bound to the
154 specific network endpoint.

```
public class HelloImpl implements IHello
{
    @Override
    public void say(String greetings) {
        System.out.println(greetings);
    }

    public static void main(String[] a){
        Server server = Server.build(
            "tcp://localhost:5555",
            new HelloImpl());

        server.wait();
    }
}
```

Figure 4. Simple RPC server (Java).

```
public static void main(String[] a)
{
    IHello hello = Client.build(
        IHello.class,
        "tcp://localhost:5555");

    hello.say("Hello, World");
    ...
}
```

Figure 5. Simple RPC Client (Java).

155 The client implementation is shown in Figure 5. It uses the static *Client.build* function giving it two
156 parameters: the server's endpoint and the interface class. The call returns an instance of the special
157 implementation of the *IHello* interface, which can be used to translate a local call to an interface's
158 method to the invocation of the corresponding function on the remote server.

159 5. ZeroMQ performance and scalability tests

160 Before choosing ZeroMQ as the communication layer implementation for the TDAQ IPC software we
161 performed several tests to verify its performance and understand how it scales with the number of
162 communicating applications. For comparison we also repeated the same tests with the omniORB
163 CORBA broker and the ICE [9] framework from ZeroC company. ICE is a modern CORBA-like
164 object-oriented RPC framework, which is free of many CORBA drawbacks, but still provides a
165 complex high-level monolithic IPC solution.

166 In order to test how well the communication software scales with the number of communicating
167 clients, the following configurations was tested for each of them:

- 168 • A single server application running on a dedicated computer and answering to all the clients'
169 requests. For all tests the servers were running the same number of threads: 1 I/O thread for
170 incoming request routing and 20 worker threads for executing request code.
- 171 • Client applications were equally distributed on a cluster of computers connected to the server
172 via the local network. Each client sent a request containing a 1-byte string to the server and
173 received the same string back as fast as possible.

174 5.1. Hardware configuration for the tests

175 All tests were performed at a facility comprising around 200 computers connected to a local network
176 via 10Gb Ethernet. Server applications ran on a commodity server with Intel Xeon E5645 4 core 2.4
177 GHz CPU [10] and 24 GB RAM. Client applications were equally distributed over 100 computers
178 with Intel Xeon E5420 4 core 2.5 GHz CPU [11] and 16 GB RAM.

179 5.2. Test results

180 Figure 6 shows the average request execution time for all three implementations with different
181 numbers of concurrent clients. The results clearly indicate that all the systems have excellent
182 scalability and offer very good performance. ZeroMQ uses a bit more time for a single request, which
183 can be explained by the differences in handling requests on the client side between ZeroMQ and the
184 other two systems.

185 Contrary to omniORB and ICE brokers ZeroMQ does not write data to the network socket from the
186 user thread. Instead the *write* function just places the data into a queue which is then handled by

187 another thread dedicated to IO operations. This thread reads this data from the queue and sends it to
188 the network.

189 Figure 7 shows the average CPU time which the servers spent in handling a single remote request. For
190 these tests the implementation of the remote method did nothing apart from returning its input
191 parameter. The results therefore show the pure overhead of the communication software itself. While
192 ICE and omniORB are very similar in that respect, ZeroMQ has much smaller overhead due its
193 simplicity.

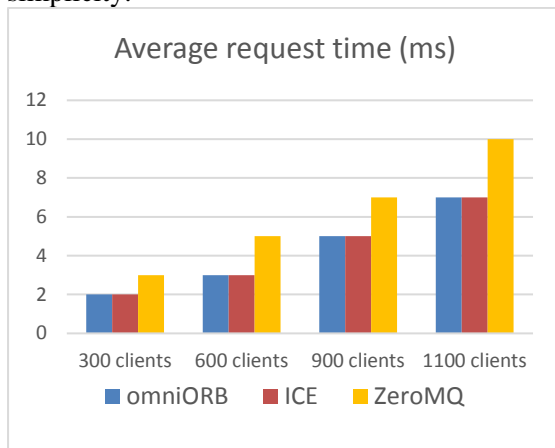


Figure 6. Average wall clock time for a single request execution (measured on the client).

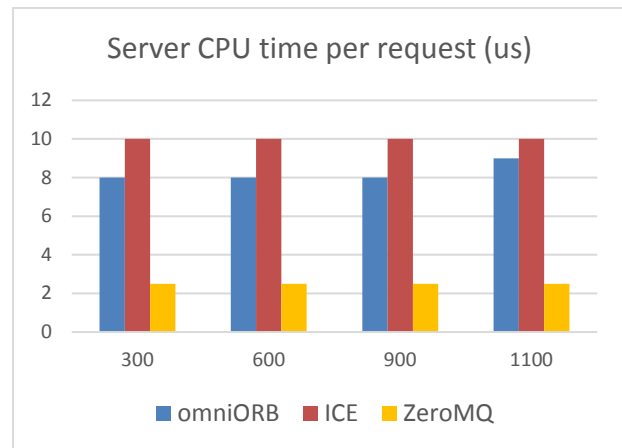


Figure 7. Average CPU time for a single request execution (measured on the server).

194 6. Conclusion

195 The CORBA standard has a long and successful history, but now the interest of the IPC software
196 development community has been shifting away from universal frameworks to relatively low-level
197 messaging systems. Many such systems are so simple that they are distributed in the form of a library
198 with a simple yet very powerful API. The new implementation of the ATLAS TDAQ IPC software
199 uses such a library called ZeroMQ in conjunction with a modular software architecture approach. This
200 approach gives an attractive alternative to the usage of a traditional high-level object-oriented
201 communication framework by offering a high performance, simple and flexible solution for
202 communication system development. The new IPC software offers a RPC-style API for implementing
203 simple communications, but at the same time makes it possible to customize data serialization along
204 with providing access to the low-level communication library API for performance optimization.

205 7. References

- 206 [1] ATLAS Collaboration 2003 ATLAS high-level trigger data-acquisition and controls *Technical*
207 *Design Report* ATLAS-TDR-016 CERN-LHCC-2003-022.
- 208 [2] ATLAS Collaboration 2008 The ATLAS experiment at the CERN Large Hadron Collider
209 *Journal of Instrumentation* JINST 3 S08003.
- 210 [3] Common Object Request Broker Architecture, www.corba.org
- 211 [4] JacORB, www.jacorb.org
- 212 [5] omniORB, omniorb.sourceforge.net
- 213 [6] Distributed Messaging, www.zeromq.org
- 214 [7] A Java serialization/deserialization library, <https://github.com/google/gson>
- 215 [8] Java Reflection API tutorial, <https://docs.oracle.com/javase/tutorial/reflect/>
- 216 [9] Ice - Comprehensive RPC Framework, <https://zeroc.com/products/ice>
- 217 [10] Intel Xeon Processor E564, [http://ark.intel.com/de/products/48768/Intel-Xeon-Processor-](http://ark.intel.com/de/products/48768/Intel-Xeon-Processor-E5645-12M-Cache-2_40-GHz-5_86-GTs-Intel-QPI)
218 [E5645-12M-Cache-2_40-GHz-5_86-GTs-Intel-QPI](http://ark.intel.com/de/products/48768/Intel-Xeon-Processor-E5645-12M-Cache-2_40-GHz-5_86-GTs-Intel-QPI)
- 219 [11] Intel Xeon Processor E5420, [http://ark.intel.com/products/33927/Intel-Xeon-Processor-E5420-](http://ark.intel.com/products/33927/Intel-Xeon-Processor-E5420-12M-Cache-2_50-GHz-1333-MHz-FSB)
220 [12M-Cache-2_50-GHz-1333-MHz-FSB](http://ark.intel.com/products/33927/Intel-Xeon-Processor-E5420-12M-Cache-2_50-GHz-1333-MHz-FSB)