

Studies for a Common Selection Software Environment in ATLAS: From the Level-2 Trigger to the Offline Reconstruction

S. Armstrong, J. T. Baines, C. P. Bee, M. Biglietti, A. Bogaerts, V. Boisvert, M. Bosman, S. Brandt, B. Caron, P. Casado, G. Cataldi, D. Cavalli, M. Cervetto, G. Comune, A. Corso-Radu, A. Di Mattia, M. Diaz Gomez, A. dos Anjos, J. Drohan, N. Ellis, M. Elsing, B. Epp, F. Etienne, S. Falciano, A. Farilla, S. George, V. Ghete, S. González, M. Grothe, A. Kaczmarska, K. Karr, A. Khomich, N. Konstantinidis, W. Krasny, W. Li, A. Lowe, L. Luminari, C. Meessen, A. G. Mello, G. Merino, P. Morettini, E. Moyses, A. Nairz, A. Negri, N. Nikitin, A. Nisati, C. Padilla, F. Parodi, V. Perez-Reale, J. L. Pinfold, P. Pinto, G. Polesello, Z. Qian, S. Resconi, S. Rosati, D. A. Scannicchio, C. Schiavi, T. Schoerner-Sadenius, E. Segura, J. M. de Seixas, T. Shears, S. Sivoklov, M. Smizanska, R. Soluk, C. Stanescu, S. Tapprogge, F. Touchard, V. Vercesi, A. Watson, T. Wengler, P. Werner, S. Wheeler, F. J. Wickens, W. Wiedenmann, M. Wielers, and H. Zobernig

Abstract—The ATLAS High Level Trigger’s (HLT) primary function of event selection will be accomplished with a Level-2 trigger farm and an event filter (EF) farm, both running software components developed in the ATLAS offline reconstruction framework. While this approach provides a unified software framework for event selection, it poses strict requirements on offline components critical for the Level-2 trigger. A Level-2 decision in ATLAS must typically be accomplished within 10 ms and with multiple event processing in concurrent threads. To address these constraints, prototypes have been developed that incorporate elements of the ATLAS data flow, high level trigger, and offline framework software. To realize a homogeneous software environment for offline components in the HLT, the Level-2 Steering Controller was developed. With electron/gamma- and muon-selection slices it has been shown that the required performance can be reached, if the offline components used are carefully designed and optimized for the application in the HLT.

Index Terms—Programming environments, software reusability, triggering.

I. INTRODUCTION

THE Large Hadron Collider (LHC) currently under construction at CERN will produce pp -collisions with a center of mass energy of $\sqrt{s} = 14$ TeV at a design luminosity of 10^{34} $\text{cm}^{-2}\text{s}^{-1}$. With a bunch crossing rate of 40 MHz and about 23 interactions per bunch crossing, it requires, however, highly selective trigger systems to reduce the expected 10^9 interactions per second to an acceptable rate of a few hundred Hz. ATLAS [1] is one of the two large general purpose experiments at the LHC and covers a widely diversified physics program [2], ranging from discovery physics to precision measurements of Standard Model parameters. ATLAS has an inner detector for precision tracking mounted inside a superconducting magnet with 2 T field strength. Outside the solenoid follow electromag-

netic and hadronic calorimeters. Muon identification is achieved with a high precision muon spectrometer. The total number of readout channels is about 10^8 .

Given the required large selectivity of the ATLAS trigger ($\approx 10^{-7}$) and the rare nature of the most interesting physics signatures at the LHC collider, it is essential to understand the efficiencies at each step of the event selection process. Sharing a large number of software components across all platforms from the trigger event selection software to the offline physics analysis and reconstruction environment helps in achieving this goal and allows for a common development and run environment.

II. THE ATLAS TRIGGER

The ATLAS trigger is based on three levels of online selection: Level-1, Level-2, and Event Filter (EF). The second and third level triggers, together known as the high level trigger (HLT) [3], [4], are implemented on PCs running the Linux operating system.

The Level-1 trigger [5] is implemented in custom hardware and will reduce the initial event rate to about 75 kHz. The Level-1 decision is based on data from the calorimeters and the muon detectors. For accepted events small localized regions in pseudo rapidity η and azimuthal angle ϕ centered on the high p_T objects identified by the Level-1 trigger are determined. Each region of interest (RoI) contains the type and the thresholds passed of the associated high p_T candidate objects.

The Level-2 trigger’s selection process is guided by the RoI information supplied by the Level-1 trigger and uses full granularity event data within a RoI from all detectors. In this way, only 2% of the full event data are needed for the decision process at Level-2, thus reducing the required aggregate bandwidth to serve the Level-2 trigger. The selection algorithms request data from the readout buffers (ROB) for specific detectors in a Level-1 defined RoI for each processing step. The data are held in the ROB until the Level-2 trigger accepts or rejects the event. The Level-2 event selection algorithms are controlled by the HLT selection framework and run inside

Manuscript received November 14, 2003; revised February 3, 2004.

The authors are with THE ATLAS HIGH LEVEL TRIGGER GROUP. Information is available online at: <http://atlas.web.cern.ch/Atlas/GROUPS/DAQTRIG/HLT/AUTHORLISTS/nss2003.pdf>.

Digital Object Identifier 10.1109/TNS.2004.829548

the Level-2 Processing Units [6], [7] (L2PU) in concurrent worker-threads, each processing one event. The multithreaded approach minimizes overheads from context-switching and avoids stalling the CPU when waiting for requested RoI data to arrive from the ROBs. Asynchronous services, like input of data and application monitoring, are executed in separate threads. This allows an efficient use of multi-CPU processors but requires also all software running in the L2PU to be thread-safe. The technical aspects of multithreading are handled by the dataflow software itself, including creation and deletion of threads and any locking mechanism that may be required. The Level-2 output rate is about 2 kHz with typical event decision times of 10 ms.

If an event is accepted by Level-2, the Event Builder collects all the event data fragments from the ROBs. The complete event is then made available to the EF for the final stage of trigger processing. Here, more complex algorithms provide a further rate reduction to about 200 Hz with typical event decision times of 1–2 s. While the Level-2 reconstructs localized regions, the baseline for the EF is a full offline-like event reconstruction guided by the *Level-2 Result*. It will also use more complete calibration, alignment and magnetic field data.

To achieve a fast rejection, the event processing in the HLT selection proceeds in steps for feature extraction and hypothesis decisions. At the end of each step, the step results are checked against abstract physics signatures defined in trigger menus.

III. HIGH LEVEL TRIGGER SELECTION SOFTWARE

The HLT selection framework [8] constitutes the run environment for the trigger algorithms. It is common to Level-2 and EF and is composed of four main components. The *HLT Steering* schedules the *HLT Algorithms* corresponding to the input seed, so that all necessary data for a trigger decision are produced. Information about event specific quantities is exchanged via components of the *event data model* (EDM). During event processing data are stored and accessed through a *data manager*. This allows to hide platform- and storage technology-specific details of event data access from the algorithms. The *HLT algorithms* either reconstruct new event quantities or check trigger hypotheses with previously computed event features.

As the main purpose of the HLT selection software is event selection, it has to run efficiently and reliably in the online environment. In addition, critical selection components must be transferable to the offline environment for development and testing purposes. By providing a common code base for the online and the offline software, the HLT guarantees the consistency of trigger performance evaluations. It also provides a “physicist-friendly” environment for trigger algorithm development. In addition, studies have already shown [9] that great cost savings can be obtained with the proper global optimization of the trigger. Having a single common framework, where the different trigger levels can be cross optimized, greatly facilitates these studies.

Since the EF provides an offline-like environment, the HLT selection software is naturally based on the ATLAS offline reconstruction and analysis environment ATHENA [10], which itself is based on the GAUDI [11] framework. This allows for the

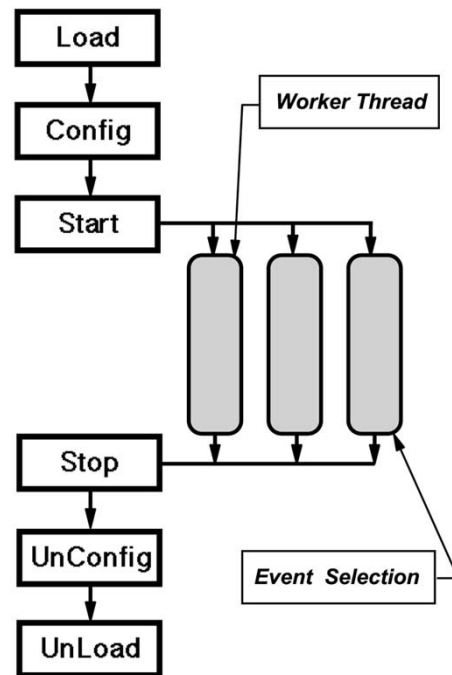


Fig. 1. The L2PU finite state machine. All algorithms and services are created at initialization time in the *Configure* step. The Level-2 selection code is executed in parallel in multiple worker-threads.

reuse of the storage manager, the EDM, the detector description and many algorithms, which are already developed by the offline community. Only the *HLT Steering* framework and certain algorithms remain as HLT specific developments. In the case of the Level-2 trigger, a similar ansatz is more difficult due to the multithreaded selection process and the more severe performance requirements. Even though Level-2 algorithms are specially developed to meet the tight timing requirements, they use the same EDM- and detector description objects present in the EF and offline. A transparent use of such components is possible and a common implementation of the HLT framework for both Level-2 and EF can be realized if the same interfaces are available at Level-2. This is provided by the *Level-2 steering controller* (SC).

IV. THE LEVEL-2 STEERING CONTROLLER

The SC [12] is the software component that interfaces the L2PU, the dataflow application which provides access to the event data stored in ROBs and the HLT event selection software. The purpose of the SC is threefold: to allow the L2PU to host and control the selection software; to allow the reuse of the same trigger algorithm steering software as in the EF; and to provide a mechanism for transmitting the results of Level-1 and Level-2 processing between the data acquisition system and the event selection software. All dataflow applications follow for control a state model implemented in form of *finite state machines* (FSM). The key to the SC design is to place this interface where the functionality of the dataflow and event selection frameworks can be cleanly separated. The location chosen is the FSM of the L2PU.

The SC provides the means for forwarding state changes from the dataflow software (Fig. 1) to the event selection software. An

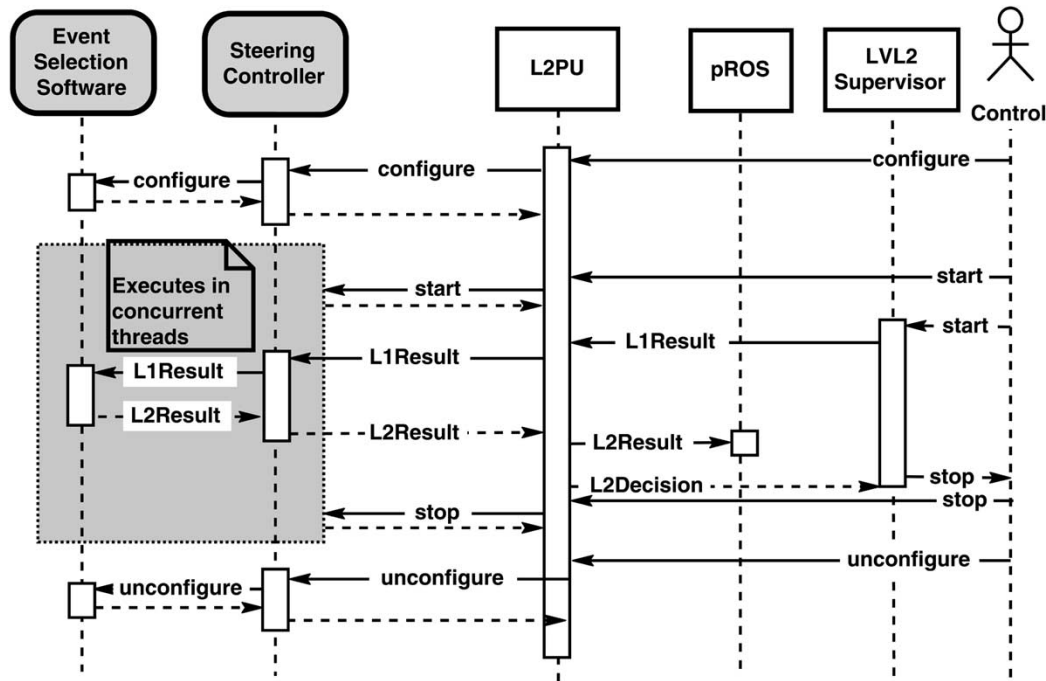


Fig. 2. The sequence of interactions of the SC with the L2PU and the Event Selection Software. Three states are shown: *configure*, *start*, and *stop*. The gray area shows the interactions which happen in multiple Worker Threads. The *pROS* is the component that forwards the *Level-2 Result* to the EF.

important aspect of this approach is that the Level-2 event data access is managed entirely by the dataflow. The SC then does not need to interact directly with the data input threads or other dataflow components.

Fig. 2 illustrates the sequence of interactions of the SC with the L2PU and the event selection software. The figure shows three states: *configure*, *start*, and *stop*. During the *configure* phase, configuration, and conditions data are obtained from external databases via an HLT-online interface. These data are then used to configure the selection software and all associated components.

After a *start* the SC receives an “execute event” directive with a *Level-1 Result* as argument. The result of event processing is directly returned as *Level-2 Result* to the L2PU. A *stop* command terminates algorithm execution and produces run summary information.

Since the trigger event selection software is being developed in the ATLAS offline framework, which is itself based on the GAUDI framework, the SC also has been designed to reuse the framework interfaces defined in GAUDI. In this way, there is a unified environment for event reconstruction and selection from the second level trigger to the offline analysis.

Since the event selection software executes in multiple worker threads, the SC must provide a thread-safe environment. At the same time, and in order to provide an easy-to-use framework for offline developers, the SC must hide all technical details of thread handling and locks. Thread safety has been implemented in the SC by using GAUDI’s name-based object and service bookkeeping system. Copies of components that need to be thread-safe are created for each worker thread with different labels. The labels incorporate the thread-ID of the worker thread, as obtained from the dataflow software. The number of threads created by the dataflow software is

transferred to the SC, which transparently creates the number of required copies. In this scheme, the same configuration can be used in the offline and in the Level-2 environments; the thread-ID collapses to null in the offline environment, as it is not needed there.

In contrast to the EF and offline environment, all Level-2 algorithm- and service instances need to be created and initialized in the configuration state of the L2PU FSM, since only in this phase a L2PU has access to external databases. Later in the event loop, only access to information stored locally on the processor is possible. The *configure* step is still executed in a single thread and only later after the transition to the *start* state of the FSM the thread specific copies of the algorithms and services are attached to the worker threads.

The implementation of the SC consists of three components.

- An interface class, which forward the L2PU state changes to the algorithm execution framework. Different interface implementations, e.g., for dataflow testing without algorithms, can be specified in the L2PU configuration and are loadable as shared libraries. This implementation also helps to minimize cross dependencies in the respective dataflow and offline repositories.
- The multithreaded algorithm execution environment. The necessary changes for multithreading support have been incorporated in the GAUDI base libraries and are available with the recent official releases of the GAUDI framework.
- A modified ATHENA event loop handler. Contrary to offline, the event loop is controlled by the dataflow software. The modified event loop handler makes the *Level-1 Result* available to the HLT selection software, executes the algorithms for a given event and forward the Level-2 decision to the L2PU interface class.

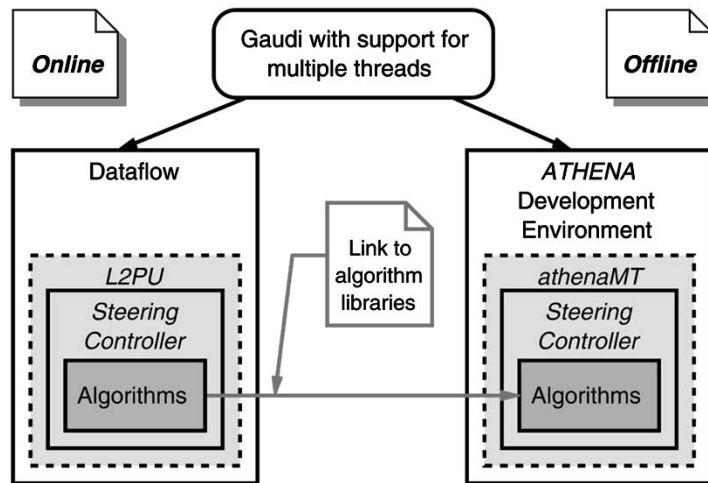


Fig. 3. A typical development and prototype setup. GAUDI with support for multiple threads is the basic algorithm execution environment. It is shared by the online and the offline environments. Algorithms are developed in the offline environment and tested with the L2PU emulator *athenaMT*. The same binary libraries containing the event selection algorithms are then used by the real L2PU running in a dataflow setup.

The implementation is complemented by special utility services, which connect, e.g., the GAUDI framework messaging to the corresponding dataflow implementations.

V. SOFTWARE DEVELOPMENT MODEL

HLT software developers follow a typical edit, compile and run cycle in the ATHENA offline environment, when creating new software components. For running an application, the ATHENA main program together with a job configuration file would be used: `athena <job-configuration>`.

Since the same interfaces are available in the EF and the L2PU environment the code developed in the offline environment can be directly downloaded in binary form to the processors. For Level-2 the developer has to follow however a set of simple coding rules [13] to produce thread safe code and to make its algorithms or services compatible with the automatic creation of multiple copies in the L2PU. Furthermore, it should be not necessary to use *locks* or *mutexes* to adapt the code to the multithreaded environment. To meet the timing requirements for the second level trigger, the number and type of available services is restricted to the necessary minimum, which is a subset of the services available in the EF and in offline. In this way, it should be always possible to move a Level-2 component to the EF and offline environment. The other direction, however, is only possible if the software component finds all its necessary resources in the restricted L2PU environment.

The dataflow software can be configured to run either as single or multinode system. The single-node system starts a Read-Out Subsystem (ROS) emulator, a Level-2 supervisor, and a L2PU in the same processing node, while the multinode system distributes these applications over various nodes. The setup of a complex dataflow system for application testing is in both cases a nontrivial task and most HLT developers lack also the necessary hardware resources. A multithreaded test application called *athenaMT* was therefore created, which presents internally the same run environment as a L2PU, but can be started as simple as the normal ATHENA main program: `athenaMT <number of worker-threads>`

TABLE I
RATES ON A DUAL-PROCESSOR 1.533-GHz AMD ATHLON MACHINE

Prototype configuration	Measured rate	Overhead/event
L2PU	21.7 kHz	46 μ s
L2PU+SC	17.0 kHz	59 μ s
L2PU+SC+Data Manager	15.3 kHz	65 μ s

`<job-configuration>`. *athenaMT* uses the SC and most of the dataflow components that are also used in a standard L2PU. It differs, however, in the supervision aspect of the L2PU and in the way detector raw data are made available to the processing unit. The application can be used on single- or multi-CPU machines. HLT developers need not to be familiar with detailed technical aspects of the dataflow software and are also shielded from changes in the dataflow part of the software. They can concentrate exclusively on the HLT software and are able to perform a large variety of useful tests, from thread safety to performance measurements, in a realistic L2PU environment. Fig. 3 shows the relation between an online dataflow setup and an offline development environment for Level-2 software. In this way the development effort for HLT and dataflow can be widely parallelized. It is clear, however, that the final certification of the HLT software has to be done on a large distributed system.

VI. PERFORMANCE MEASUREMENTS

After integrating the SC with the dataflow software, both performance and robustness tests were carried out on a dual-processor 1.533-GHz AMD Athlon machine. The SC ran for over 50 h with three threads with an early version of the selection software prototype. The prototype ran successfully on both single-CPU and double-CPU machines, showing it to be thread safe. A direct measurement of the SC overhead yielded 13 μ s per event. The overhead was estimated by comparing the number of events per second a L2PU can handle when running without or with the SC and executing a simple algorithm (first and second row in Table I). Table I shows also the obtained rate, when in addition the *Level-1 Result* is transferred to the *data manager*.

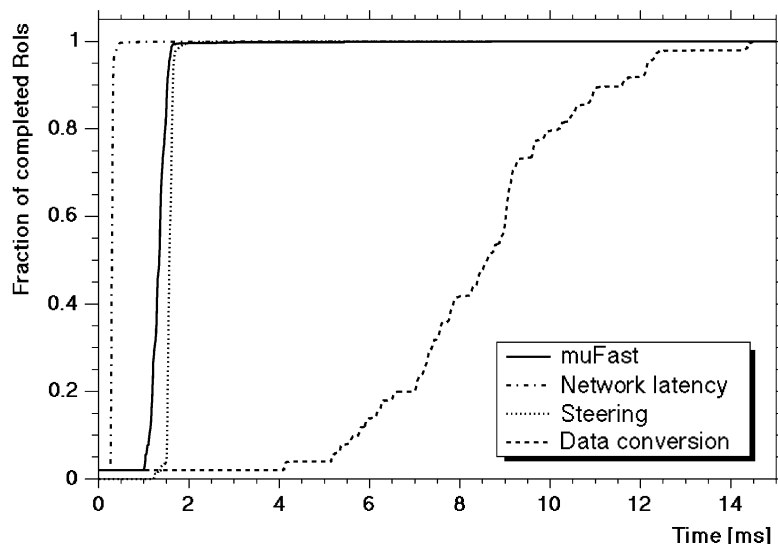


Fig. 4. The main contributions to the latency shown as curves of integrals for the μ selection slice. The slice used offline software components for data conversion, detector description, data manager, and steering. The contributions in order of decreasing importance are data preparation and conversion (Data conversion), HLT framework overheads (Steering), algorithmic processing (muFast), and network access time (Network latency). The curves show, e.g., that in this setup algorithmic processing is terminated for 95% of the events in less than 2 ms. The data sample consisted of 200 GeV muons at high luminosity. The muon background from the ATLAS cavern was boosted by a factor of 2, so that the results give a conservative estimate of the processing times.

The quoted overhead includes all GAUDI framework steps to schedule and execute algorithms and to execute the used base services. The measurements were based on runs of at least 100 000 events and almost perfect scaling of the measured latency was observed when the algorithm execution time was varied with a CPU burning loop from 0 to 8 ms in the different configurations, indicating that the overheads per event introduced by the SC are independent of the algorithm execution time.

More tests with increasing complexity were performed on a three-node dataflow system build with a dual-processor Intel XEON 2.2 GHz machine hosting the L2PU, a dual-processor 1.533-GHz AMD Athlon machine hosting the ROS emulator and a single-processor Intel Pentium 4 machine running the Level-2 Supervisor. The dual-processor machines were connected via a Gigabit Ethernet network. All machines were running the Redhat 7.3 Linux operating system.

In sequence, the offline *data manager* StoreGate [14] the *HLT steering* and a fast inner tracker feature extraction algorithm [15] were included in the setup. The algorithm used an early version of the *event data model* and the raw data conversion process from byte stream format to EDM classes. More complex tests were recently done with complete e/γ and μ selection slices. They contained the complete *HLT steering* framework with decoding of the *Level-1 Result*, scheduling of the feature extraction algorithms, and sending the results of the Level-2 processing to the EF. The *data manager*, the *event data model*, the detector description for the calorimeters, and the muon system and services for conversion from raw data byte stream format to high-level data containers were used from offline. The feature extraction algorithms were specially developed for Level-2 but used the aforementioned offline services. In the case of the e/γ slice, over 95% of the events were processed within 5 ms for a sample of di-jet events at low luminosity and a RoI size of $\Delta\eta \times \Delta\phi = 0.3 \times 0.3$ [4]. Fig. 4 shows similar results for the μ slice.

The overheads of the SC are negligible compared to the processing times of the HLT event selection software. For both selection slices the data preparation and conversion step dominates the processing times. It contains all contributions from unpacking of the raw data, the application of the complex detector description and calibration schemes to the creation of fully calibrated high level physics objects suited for event reconstruction and selection in the algorithms. The data preparation and conversion time is a direct function of the number of handled input raw data items and is therefore also sensitive, e.g., to the noise level in the calorimeters and the occupancy of the muon chambers due to cavern background. Appropriate data preprocessing and noise reduction schemes can help to reduce the time spend in this step [4]. The contribution from the *HLT steering* and the network latency are similar for both selection slices. Measurements on larger test beds [7] and detailed simulations of the full system have shown [4], that at most a contribution of 1 ms from the RoI data collection over the network is expected.

VII. EXPERIENCES

During the tests it was observed that the event throughput in a L2PU did not scale in the expected way with the number of worker-threads. This was due to the use of a common memory pool for container objects in the default memory allocation scheme of the *standard template library* (STL). The event processing model of Level-2 favors a scheme where every thread allocates its own memory pool. Such an allocation scheme is available in the STL. After carefully optimizing the code with it, the expected scaling behavior was observed. Such an optimization poses limitations on the offline components used and their external dependencies. For instance, external utility libraries may not be available with this allocation scheme. To avoid frequent memory allocation during event processing, all large containers that hold detector data were designed to

allocate memory only during initialization time and to reset their data for each new event.

The reconstruction of full events in offline favors a processing model where utility services and external data are retrieved on demand. In the case of Level-2 this model cannot be applied since during event processing only locally stored meta-data are accessible. All required data need to be known at configuration time and need to be prefetched by the processor. Furthermore, the creation of large configuration objects on demand may lead to time-outs during event processing. These restrictions required a redesign and alternative initialization methods for some offline components, especially for detector description and raw data conversion.

The use of offline components in multiple worker threads and the requirement to avoid *locks* in the event selection code limits certain design and implementation choices, e.g., the use of *singletons*. These restrictions were communicated to the HLT developers at an early design stage.

VIII. CONCLUSION

The presented implementation of the SC for the ATLAS Level-2 trigger enables the reuse of offline software components throughout the ATLAS High Level Triggers. It realizes a homogeneous software and development environment from the Level-2 trigger to offline. Realistic prototypes have shown that the required performances can be reached if the offline components are carefully optimized and designed for reuse in the triggers. This may limit architectural, design, and implementation choices that are otherwise available in a pure offline environment. An understanding of these restrictions is necessary for all contributing developers.

ACKNOWLEDGMENT

The authors would like to acknowledge the help and support of the ATLAS Data Acquisition Group and the ATLAS Offline and Detector software groups.

REFERENCES

- [1] "ATLAS: Technical Proposal for a General-Purpose Experiment at the Large Hadron Collider at CERN," ATLAS Collaboration, CERN/LHCC/94-43, 1994.
- [2] "ATLAS Detector and Physics Performance Tech. Design Rep.," ATLAS Collaboration, CERN/LHCC/99-014 and 99-015, 1999.
- [3] "ATLAS High-Level Triggers, DAQ and DCS Tech. Proposal," ATLAS Collaboration, CERN/LHCC/2000-017, 2000.
- [4] "ATLAS High-Level Trigger Data Acquisition and Controls, Tech. Design Rep.," ATLAS Collaboration, CERN/LHCC/2003-022, 2003.
- [5] "ATLAS Level-1 Trigger Tech. Design Rep.," ATLAS Collaboration, CERN/LHCC/98-014, 1998.
- [6] A. Bogaerts and F. Wickens, "LVL2 Processing Unit Application Design," ATL-DH-ES-0009, 2001. EDMS Note.
- [7] S. Armstrong *et al.*, "The second level trigger of the ATLAS experiment at CERN's LHC," in *Proc. Nuclear Science Symp.*, Portland, OR, Oct., 19–25 2003.
- [8] M. Elsing, Ed., "Analysis and Conceptual Design of the HLT Selection Software," PESA Software Group, ATL-DAQ-2002-013, 2002. ATLAS Internal Note.
- [9] J. T. M. Baines, S. González, R. K. Mommsen, A. Radu, T. G. Shears, S. Sivoklov, and M. Wielers, "First Study of the LVL2-EF Boundary in the High-pT e/γ High-Level Trigger," ATL-DAQ-2000-045, 2000. ATLAS Internal Note.
- [10] Athena: User Guide and Tutorial. [Online]. Available: <http://atlas.web.cern.ch/Atlas/GROUPS/SOFTWARE/OO/architecture/General/Tech.Doc/Manual/2.0.0-DRAFT/AthenaUserGuide.pdf>
- [11] Gaudi Project [Online]. Available: <http://cern.ch/proj-gaudi/>
- [12] S. González, A. Radu, and W. Wiedenmann, "Use of Gaudi in the LVL2 Trigger: The Steering Controller," ATL-DH-EN-0001, ATL-DAQ-2002-012, 2002. EDMS Note, and ATLAS Internal Note.
- [13] A. Bogaerts, A. dos Anjos, S. González, and W. Wiedenmann. Guidelines for Offline Preparation and Testing of LVL2 Code. [Online]. Available: <http://www.cern.ch/steve.armstrong/algorithms/guidelines>
- [14] P. Calafiura, H. Ma, M. Marino, S. Rajagopalan, and D. R. Quarrie. StoreGate: A data model for the ATLAS software architecture. presented at *Proc. Computing in High Energy and Nuclear Physics Conf.* [Online]. Available: <http://atlas.web.cern.ch/Atlas/GROUPS/SOFTWARE/OO/architecture/EventDataModel/Tech.Doc/StoreGate/CHEP01.pdf>
- [15] J. T. M. Baines, A. Bogaerts, M. Bosmann, A. dos Anjos, S. González, K. Karr, W. Li, G. Merino, J. Schlereth, and W. Wiedenmann, "Initial LVL2 Tests with the Si Tree Algorithm," ATL-DH-TN-0001, 2003. EDMS Note.