

1 Multi-threaded ATLAS Simulation on Intel Knights 2 Landing Processors

3 **Steven Farrell¹, Paolo Calafiura¹, Charles Leggett¹, Vakhtang
4 Tsulaia¹, Andrea Dotti², on behalf of the ATLAS Collaboration**

5 ¹Lawrence Berkeley National Laboratory, ²SLAC National Accelerator Laboratory

6 E-mail: SFarrell@lbl.gov

7 **Abstract.** The Knights Landing (KNL) release of the Intel Many Integrated Core (MIC)
8 Xeon Phi line of processors is a potential game changer for HEP computing. With 72 cores and
9 deep vector registers, the KNL cards promise significant performance benefits for highly-parallel,
10 compute-heavy applications. Cori, the newest supercomputer at the National Energy Research
11 Scientific Computing Center (NERSC), was delivered to its users in two phases with the first
12 phase online at the end of 2015 and the second phase now online at the end of 2016. Cori
13 Phase 2 is based on the KNL architecture and contains over 9000 compute nodes with 96GB
14 DDR4 memory. ATLAS simulation with the multithreaded Athena Framework (AthenaMT)
15 is a good potential use-case for the KNL architecture and supercomputers like Cori. ATLAS
16 simulation jobs have a high ratio of CPU computation to disk I/O and have been shown to scale
17 well in multi-threading and across many nodes. In this paper we will give an overview of the
18 ATLAS simulation application with details on its multi-threaded design. Then, we will present
19 a performance analysis of the application on KNL devices and compare it to a traditional x86
20 platform to demonstrate the capabilities of the architecture and evaluate the benefits of utilizing
21 KNL platforms like Cori for ATLAS production.

22 1. Introduction

23 In the multi-core computing era, processor chip trends such as increasing core multiplicity,
24 decreasing memory per core, and increasing importance of vector processing are changing the
25 way scientific software developers write efficient, scalable code. Modern computing devices such
26 as Intel's Xeon Phi line of many-core processors are good examples of what will be used more
27 frequently in high performance computing facilities. These devices are best utilized with highly-
28 parallel applications, so scientific computing models must adapt for greater concurrency and
29 intelligent usage of memory resources.

30 High energy physics (HEP) experiments such as ATLAS[1] are no exception to this paradigm
31 shift. Particle collision data is typically trivially parallelizable, but production software such as
32 the Athena framework[2] have historically been written for sequential processing. In order to
33 ensure that ATLAS can efficiently utilize modern computing devices and devices of the future,
34 a large campaign is underway to adopt a multi-threading concurrency model for parallelism and
35 efficient use of memory resources[3][4]. ATLAS simulation is the most advanced use-case for
36 multi-threading, with a nearly complete configuration working and performing well on traditional
37 Intel Xeon devices.



38 In this paper we will share results and experience preparing the ATLAS simulation software
39 for the Knights Landing generation of Intel Xeon Phi processors. Section 2 gives a brief overview
40 of the hardware used. Section 3 details the multi-threaded ATLAS simulation application.
41 Performance results on Xeon and Xeon Phi machines are then given in Section 4. Ideas for
42 future work and conclusions are given in Section 5.

43 **2. Intel Xeon Phi processors**

44 Current state-of-the-art processors for high-performance computing offer a wide array of
45 capabilities and challenges. Devices such as FPGAs and general-purpose GPUs offer a high
46 degree of parallelism with low power consumption for effective throughput. However, both of
47 these devices use highly specialized programming models and have challenging constraints on
48 memory capacity and data bandwidth. In response, Intel has been pursuing an alternate model
49 that promises high performance with ease of use: the Xeon Phi product line.

50 Intel Xeon Phi processors are built with Intel’s Many-Integrated-Core (MIC) architecture.
51 General features of the product include high core multiplicity, deep vector registers, and low
52 power consumption (relative to Xeon devices). Xeon Phi chips run a Linux OS, making them
53 substantially easier to use than FPGAs and GPUs.

54 The current (2nd) generation of Xeon Phi processors is codenamed Knights Landing (KNL).
55 KNL chips are the first release of the product line to offer full x86 binary compatibility and the
56 first which can be installed as a host device or as a coprocessor. They are available with up to 72
57 Airmont cores and 4-way hardware threads, giving a maximum of 288 threads of execution. For
58 SIMD parallelism, KNL devices have two 512-bit vector units per core and support AVX-512
59 instructions. Finally, the KNL generation introduces a deeper memory hierarchy compared to
60 previous releases, providing both traditional DDR4 RAM as well as 8-16 GB of on-package, high-
61 bandwidth MCDRAM. The MCDRAM can be utilized as an additional addressable memory
62 space (“flat” mode), as a transparent cache (“cache” mode), or as a mixture of both (“hybrid”
63 mode).

64 Xeon Phi processors are well suited for high performance computing facilities. A number
65 of planned supercomputers will be based on Xeon Phi processors. At NERSC, the Cori
66 supercomputer will have 9,300 KNL nodes with 68 cores each (2.5 million possible threads
67 of execution). The Theta system at Argonne National Lab will have over 2,500 KNL nodes
68 as well and will be a stepping-stone machine for the massive future Aurora system. Aurora
69 is planned for 2018 to have over 50 thousand nodes equipped with 3rd-generation Xeon Phi
70 (codenamed Knights Hill) processors.

71 **3. Multi-threaded ATLAS simulation**

72 The ATLAS simulation application (G4Atlas) is used to produce simulated ATLAS data in the
73 Athena production framework[6]. It has been used extensively in the ATLAS experiment for
74 many years for data analysis. It uses the Geant4[7] particle simulation toolkit to model physics
75 processes and detector response. Production is traditionally performed with sequential jobs or
76 multi-process jobs in the AthenaMP framework[8]. In the latter case, worker processes are forked
77 from the main process after initialization of the job and before the event loop. This procedure
78 allows worker processes to implicitly share some memory pages via the Linux copy-on-write
79 mechanism.

80 An effort is currently underway to migrate the ATLAS simulation application to a
81 multi-threading processing model (G4AtlasMT) in the AthenaMT (Multi-threaded Athena)
82 framework. AthenaMT, which is based on the Gaudi concurrent framework, uses Intel Threading
83 Building Blocks (TBB) for task-based parallelism. It schedules algorithms to operate on event
84 data as tasks to run concurrently on different threads. This model allows both inter-event and
85 intra-event parallelism. The simulation application uses few algorithms, however, with most of

86 the computation work happening in one algorithm (G4AtlasAlg) which simply invokes Geant4.
 87 The result is that the G4Atlas runs effectively with only inter-event parallelism. Memory
 88 savings are achieved by sharing physics and geometry tables across threads within Geant4.
 89 An illustration of the AthenaMT algorithm scheduling model is shown in Figure 1.

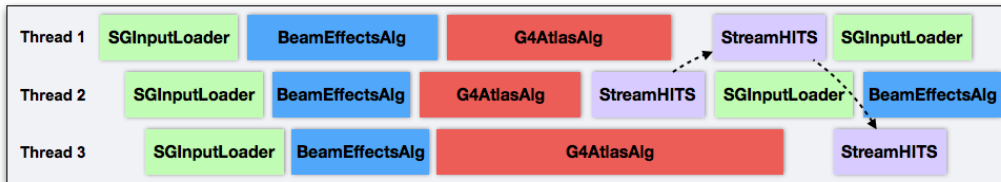


Figure 1. Illustration of worker thread processing in ATLAS multi-threaded simulation. SGInputLoader preloads some data from the input file to kickoff the event data flow. BeamEffectsAlg applies beam corrections and smearing to the input generated event. G4AtlasAlg is the main simulation algorithm which invokes Geant4. StreamHITS is the output stream algorithm which writes hit collections to the output file. StreamHITS is not cloned for concurrent processing. One instance serves all worker threads. Algorithm sizes are not shown to scale.

90 ATLAS simulation is potentially a good use-case for Xeon Phi processors. Relative to other
 91 ATLAS production workloads, simulation is CPU-heavy and uses little I/O. Not coincidentally,
 92 these are the same reasons that simulation is the primary ATLAS workload for supercomputers.
 93 The support for multi-threading is expected to be a powerful advantage in running effectively in
 94 the constrained memory environment of Xeon Phi cards. However, some challenges are expected
 95 as well. It is well known that vectorization is essential for effective utilization of KNL processors,
 96 but ATLAS simulation code does not vectorize well. Also, the highly object-oriented nature of
 97 ATLAS and Geant4 code tends to result in large code size and poor memory access patterns,
 98 which could hurt performance on KNL devices.

99 4. Performance measurements

100 The runtime performance of G4AtlasMT was measured on both Xeon and Xeon Phi machines.
 101 For the Xeon measurements, both a 16-core Ivy Bridge machine (E5-2650 v2 @ 2.60GHz) and
 102 a Cori Phase 1 Haswell node (E5-2698 v3 @ 2.30GHz) were used. The Xeon Phi measurements
 103 were taken on a KNL testbed (7210 @ 1.30GHz) for Cori Phase 2.

104 The important performance metrics are the event throughput and the memory consumption
 105 (RSS), and the scaling of these metrics with the number of worker threads. Figures 2 and 3
 106 show the measurements for the simulation of a $Z \rightarrow \tau\tau$ sample. On the Xeon, the throughput
 107 scales perfectly up to the physical number of cores on the machine (16), and small gain is seen in
 108 the hyper-threading regime. The memory consumption shows a nice gradual scaling with each
 109 additional worker thread adding only about 70 MB. On the Xeon Phi, good scaling is again
 110 seen up to the number of physical cores on the device (64), with substantial throughput gains
 111 seen in hyper-threading all the way up to the maximum 256 threads. As with the Xeon, the
 112 memory consumption on the Xeon Phi is gradual and linear, reaching about 14 GB when the
 113 device is fully loaded. For the sake of comparison, the scaling results for purely multi-process
 114 jobs are shown in Figure 4. The per-worker contribution to the memory consumption is about
 115 five times larger in multi-process jobs compared to multi-threaded jobs, a substantial reduction
 116 in memory footprint.

117 To test the scaling of G4AtlasMT in more extreme configurations, a single-muon particle gun
 118 sample was used. Whereas the $Z \rightarrow \tau\tau$ sample is representative of typical ATLAS simulation

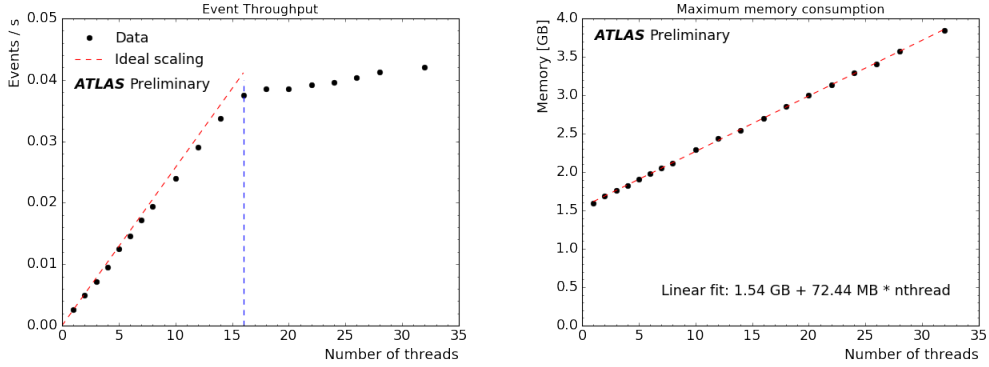


Figure 2. Event processing throughput (left) and memory consumption (right) on Intel Ivy Bridge Xeon for multi-threaded jobs with a $Z \rightarrow \tau\tau$ sample. The number of events processed is scaled as 50 times the number of threads [9].

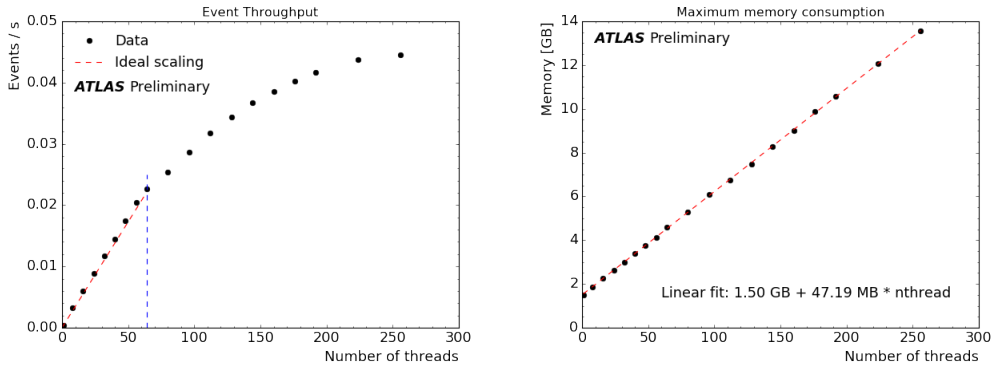


Figure 3. Event processing throughput (left) and memory consumption (right) on Intel KNL Xeon Phi for multi-threaded jobs with a $Z \rightarrow \tau\tau$ sample. The number of events processed is scaled as 10 times the number of threads [9].

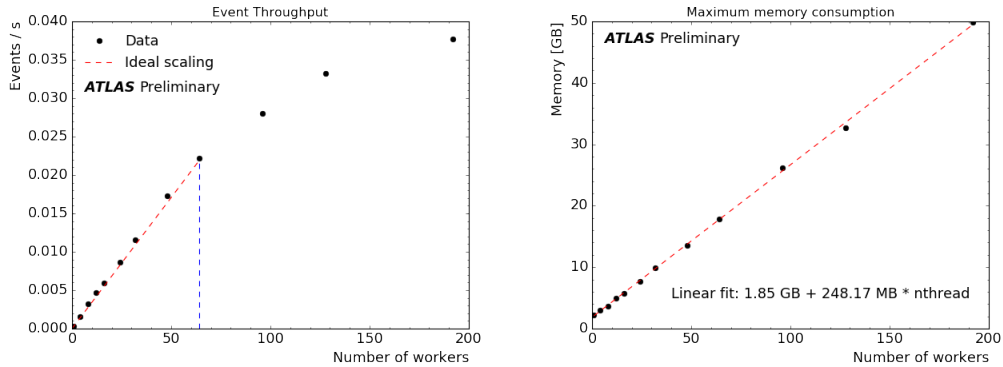


Figure 4. Event processing throughput (left) and memory consumption (right) on Intel KNL Xeon Phi for multi-process jobs with a $Z \rightarrow \tau\tau$ sample. The number of events processed is scaled as 10 times the number of threads [9].

119 production jobs and may take around 5 min per event, the single-muon sample typically takes
 120 less than one second to simulate one event. This applies more pressure to the scheduling system
 121 and other pieces of the framework infrastructure. Figures 5 and 6 show the results for the Xeon
 122 and the Xeon Phi, respectively. In this case, the throughput scales poorly above 180 threads.
 123 The source of the poor scaling was discovered to be the bottleneck in the sequential output
 124 stream which writes the simulated hit collections to the output file.

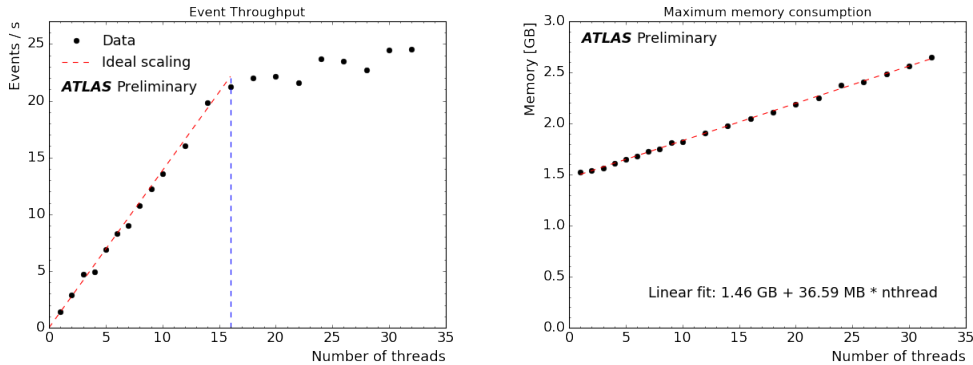


Figure 5. Event processing throughput (left) and memory consumption (right) on Intel Ivy Bridge Xeon with a single-muon sample. The number of events processed is scaled as 1000 times the number of threads [9].

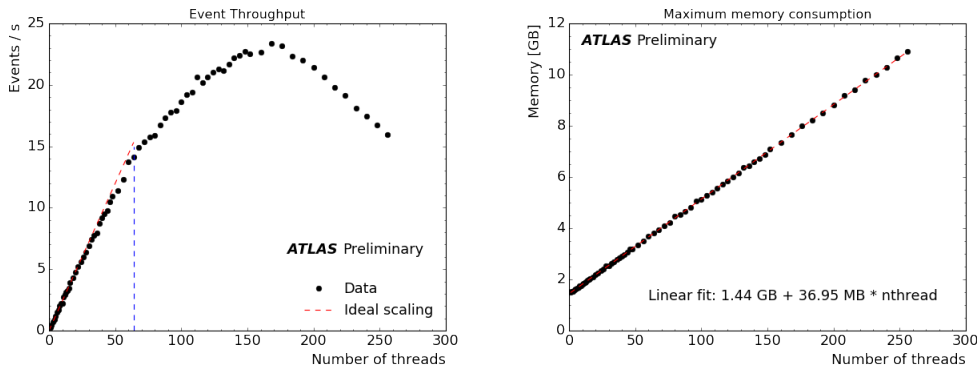


Figure 6. Event processing throughput (left) and memory consumption (right) on Intel KNL Xeon Phi with a single-muon sample. The number of events processed is scaled as 1000 times the number of threads. The sharp decrease in throughput starting around 180 threads is due to a bottleneck in the output serialization layer [9].

125 Despite good scaling results on the KNL, the absolute event throughput is not impressive.
 126 Table 1 summarizes and compares the measured event throughput for a single worker thread
 127 and for a fully-loaded device. The maximal throughput achieved on the KNL with the $Z \rightarrow \tau\tau$
 128 sample is only slightly higher than the 16-core Ivy Bridge. A fairer comparison would be a 32-
 129 core Haswell processor, which should have substantially higher throughput. The single-thread
 130 performance on KNL is observed to be about 6-7 times slower than the Xeon. While some
 131 slowdown is expected due to the reduced clock-rate and sophistication of the Airmont cores,
 132 this large difference warrants further investigation.

133 To further understand the performance characteristics on KNL, Intel VTune Amplifier was
 134 used to collect and summarize various metrics based on hardware counters. Table 2 shows some
 135 of the interesting metrics reported by VTune when comparing G4AtlasMT on a Haswell to the
 136 KNL. The clocks-per-instruction rate on Haswell is fairly reasonable, but on KNL an average
 137 of three clock cycles are needed to execute every instruction. In addition, VTune reports that
 138 the application is highly front-end bound, meaning that the processors are frequently unable
 139 to load instructions fast enough to fill the execution pipeline. Finally, we see that the rate of
 140 instruction cache misses is nearly 1 on KNL. Such results can be due to poor code layout and
 141 large code size.

Table 1. Throughput summary table for an Ivy Bridge Xeon and a KNL Xeon Phi. Results are shown for $Z \rightarrow \tau\tau$ and single-muon samples and are split for the case of a single worker thread and a fully-loaded device (or best performing configuration). Ratios of the Xeon Phi to Xeon throughput are shown in the KNL speedup column.

Sample	Threads	Throughput [events/s]		KNL speedup
		Ivy Bridge	KNL	
$Z \rightarrow \tau\tau$	single	0.00257	0.000345	0.134
	full	0.0421	0.0445	1.06
Single μ	single	1.38	0.239	0.173
	full	24.6	23.2	0.943

Table 2. Profiling metrics obtained with VTune Amplifier. A single worker thread was used to process a $Z \rightarrow \mu\mu$ sample.

Architecture	CPI rate	Front-end bound	ICache misses	Bad speculation	Back-end bound
KNL	3.0	60.2%	0.96	2.4%	18.6%
Haswell	0.9	31.5%	0.09	11.7%	27.6%

142 5. Conclusion

143 It has been shown that multi-threaded ATLAS simulation can run on Knights Landing Xeon Phi
 144 machines. Good scaling is observed in typical production samples in terms of event throughput
 145 and in memory consumption. Multi-threading allows for substantial decreases in the memory
 146 footprint of jobs relative to multi-process jobs.

147 More work is needed to understand and improve the performance on KNL in order to use this
 148 architecture effectively. The current performance achieved is comparable to a 16-core Ivy Bridge
 149 Xeon, which falls short of the full potential of KNL processors. Since the profiling studies thus
 150 far have pointed to issues with large code size and poor code layout, steps should be taken to try
 151 and mitigate these problems. Some things to try include pruning unused or unnecessary pieces
 152 of code out of the shared libraries, improving code inlining, using statically linked libraries for
 153 problematic parts of the builds (e.g. Geant4), and using profiler guided optimization to improve
 154 the binaries.

References

- 155 [1] ATLAS Collaboration, 2008 “The ATLAS Experiment at the CERN Large Hadron Collider,” JINST **3**, S08003.
156 doi:10.1088/1748-0221/3/08/S08003
- 157 [2] Calafiura P, Lavrijsen W, Leggett C, Marino M, Quarrie D 2004 “The athena control framework in production,
158 new developments and lessons learned” *Interlaken, Computing in high energy physics and nuclear physics*
159 456-458
- 160 [3] Calafiura P, Lampl W, Leggett C, Malon D, Stewart G A, Wynne B, 2015 “Development of a Next
161 Generation Concurrent Framework for the ATLAS Experiment,” J. Phys. Conf. Ser. **664**, no. 7, 072031.
162 doi:10.1088/1742-6596/664/7/072031
- 163 [4] Stewart G A *et al.*, 2016 “Multi-threaded software framework development for the ATLAS experiment,” J.
164 Phys. Conf. Ser. **762**, no. 1, 012024. doi:10.1088/1742-6596/762/1/012024
- 165 [5] Clemencic M, Hegner B, Mato P, Piparo D, 2014 “Introducing concurrency in the Gaudi data processing
166 framework,” J. Phys. Conf. Ser. **513**, no. 2, 022013. doi:10.1088/1742-6596/513/2/022013
- 167 [6] Aad G *et al.*, 2010 “The ATLAS Simulation Infrastructure,” Eur. Phys. J. C **70** 823. doi:10.1140/epjc/s10052-
168 010-1429-9
- 169 [7] Agostinelli S *et al.*, 2003 “Geant4—a simulation toolkit,” Nucl. Instrum. Meth. A **506** 250-303.
170 doi:10.1016/S0168-9002(03)01368-8
- 171 [8] Calafiura P, Leggett C, Seuster R, Tsulaia V, Gemmeren P V, 2015 “Running ATLAS workloads within
172 massively parallel distributed applications using Athena Multi-Process framework (AthenaMP)” J. Phys.
173 Conf. Ser. **664** no. 7, 072050. doi:10.1088/1742-6596/664/7/072050
- 174 [9] Farrell S, Dotti A, Calafiura P, Leggett C, Tsulaia V, 2016 “Multi-threaded ATLAS Simulation on Intel
175 Knights Landing Processors” ATL-SOFT-SLIDE-2016-739. <https://cds.cern.ch/record/2220833>
176