

Benchmarking Harp-DAAL: High Performance Hadoop on KNL Clusters

Langshi Chen¹ Bo Peng¹ Bingjing Zhang¹ Tony Liu¹ Yiming Zou¹ Lei Jiang¹
Robert Henschel² Craig Stewart² Zhang Zhang³ Emily McCallum³ Zahniser Tom³ Omer Jon³ Judy Qiu¹

¹School of Informatics and Computing, Indiana University

²UITS, Indiana University

³Intel Corporation

{lc37, pengb, zhangbj, xqiu, xl41}@indiana.edu

{yizou, jiang60, henschel, stewart}@iu.edu

{zhang.zhang, emily.l.mccallum, tom.zahniser, jon.omer}@intel.com

Abstract—Data analytics is undergoing a revolution in many scientific domains, demanding cost-effective parallel data analysis techniques. Traditional Java-based Big Data processing tools like Hadoop MapReduce are designed for commodity CPUs. In contrast, emerging manycore processors like Xeon Phi has an order of magnitude of computation power and memory bandwidth. To harness the computing capabilities, we propose a Harp-DAAL framework. We show that enhanced versions of MapReduce can be replaced by Harp, a Hadoop plug-in, that offers useful data abstractions for both of high-performance iterative computation and MPI-quality communication, and it can drive Intel’s native library DAAL. We select a subset of three machine learning algorithms and implement them within Harp-DAAL. Our scalability benchmarks run on Knights Landing (KNL) clusters and achieve up to 2.5 times speedup of performance to the HPC solution in NOMAD and 15 to 40 times faster than Java-based solutions in Spark. We further quantify the workloads on single node KNL with a performance breakdown at micro-architecture level.

Keywords-HPC, Xeon Phi, BigData

I. INTRODUCTION

In recent years, the volume and variety of data have been collected at enormous rates. The data comes from sources ranging from massive physics experiments and instruments that read our DNA to a multitude of sensors that monitor our environment. It also comes from our digitized libraries, our media streams, and our personal health monitors. Many of the primary software tools used to do the large-scale data analysis were born in the Cloud. Big data processing frameworks like Apache Hadoop are designed to run on large-scale commodity CPU clusters connected by Ethernet. While using large clusters of commodity servers is still the most cost-effective way to process petabytes or exabytes of data, the majority of data analytics jobs do not have tremendous workloads. Furthermore, machine learning algorithms become increasingly common and can easily fit into memory but require fine-grained parallelism for high performance. Modern scale-up servers like GPU and Xeon Phi provide substantial processing, memory, and I/O capabilities. Therefore small or middle-sized clusters for Big Data analytics become an attractive approach.

In this paper, we investigate and re-design optimized software stacks to effectively utilize scale-up servers in the Cloud for machine learning and data analytics applications. We conduct extensive benchmark on emerging Intel’s Many Integrated Core Architecture (MIC) based Xeon Phi Knights Landing (KNL). Our goal is to bridge the performance gap of Big Data tools and HPC systems. To the best of our knowledge, we are the first to run high-performance Hadoop for machine learning applications on KNL many-core clusters.

Here, we propose a hybrid machine learning framework named *Harp-DAAL*, which interfaces Harp¹, a highly efficient Hadoop based communication library, and *DAAL*², a native Data Analytics Acceleration Library from Intel. Harp is an open source project developed by Indiana University, which is a plug-in to the Apache Hadoop framework.

Harp has two distinctive functions: 1) Collective communication operations that are highly optimized for big data problems. 2) Efficient and innovative computation models for different machine learning problems. Original Harp project has its codebase written in Java, which is unable to leverage the capabilities of shared-memory HPC hardware. The central idea is to replace Java kernels by highly optimized native kernels or math libraries at the node level. We use Intel’s *DAAL* as the low-level kernels for Harp on HPC platforms. *DAAL* is a library that aims to provide the users with highly optimized building blocks for data analytics and machine learning applications. For each of its kernel, *DAAL* has three modes: 1) Batch Processing, 2) Online Processing, 3) Distributed Processing. The open source code of *DAAL* only provides MKL/TBB based kernels for intra-node computation while leaving the inter-node communication of the distributed processing mode to users. This motivates us to design and build the *Harp-DAAL* framework and perform further optimization as follows: 1) Intra-node, co-design and implement kernels to fully exploit the advantages of hardware architectures. 2) Inter-

¹<https://dsc-spidal.github.io/harp/>

²<https://github.com/01org/daal>

node, minimize the overhead of data conversion and data transfer to improve the scalability of codes. The rest of the paper is organized as below: In Section II, we give an overview of existing hybrid frameworks in the domain of data analytics and machine learning. Section III describes our *Harp-DAAL* interface, the techniques to reduce data conversion overhead; as well as three benchmark algorithms. In Section IV and Section V, we discuss the configuration of the experimentation and analyze the experimental results. Conclusion is included in Section VI.

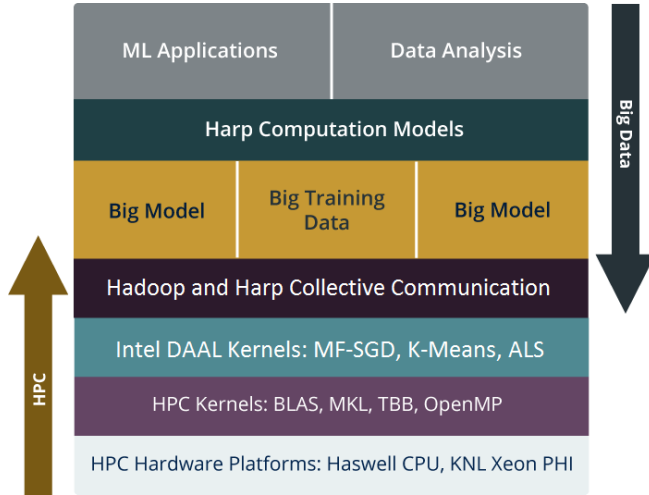


Figure 1. Harp-DAAL within HPC-BigData Stack

II. RELATED WORK

For years, scientific projects have tried to accelerate big data processing. Mars [1] combined the MapReduce framework with graphics processors and achieved 1.5 to 16 times performance improvements on a PC with G80 GPU compared to a CPU-based counterpart on web data. *MLlib* [2] is a machine learning library built upon Apache Spark, which includes a variety of machine learning applications implemented using the Scala programming language. Although it is easy to use, the performance is not satisfying, especially on HPC clusters. Awan et al. [3] have characterized the performance of Spark *MLlib* on a 24-core processor, and they conclude that Spark *MLlib* suffers from imbalanced workload on threads and poor usage of the memory system. Therefore, another approach is to build machine learning libraries in C/C++, which well suit the HPC clusters. Such efforts include *Petuum*, developed to run machine learning algorithms efficiently on any hardware [4]. However, the C/C++ based solution demands relatively high programming skills, in particular, the knowledge of multi-thread programming and architecture characteristics of multi-core/many-core processors. Other machine learning libraries build frameworks in a hybrid mode. Users only

need to know about high-level programming interfaces while letting the framework to invoke proper HPC kernels written in C/C++ or FORTRAN. Currently Spark *MLlib* can utilize the highly optimized HPC library like BLAS and LAPACK interfaced by Breeze and netlib-java. Torch³ is a scientific computing framework with broad support for machine learning algorithms from Facebook, which uses a script language called LuaJIT at the high level and the C/CUDA implementations on GPUs at the low level. Tensorflow⁴ is a deep learning framework developed by Google, which has a Python interface to write dataflow graphs, and low-level implementations on different hardware devices like CPU and GPU.

III. DESIGN AND IMPLEMENTATION

A. Data Conversion within Harp-DAAL

Harp-DAAL has a two-level structure. At the top level, it runs a group of Harp mappers, which extend Hadoop’s Java mapper. Unlike traditional MapReduce mapper, a Harp mapper holds data in main memory and invokes collective communication operations among different mappers.

At the bottom level, Harp-DAAL invokes native kernels written in C/C++, which uses multi-threading programming paradigms such as OpenMP, TBB, and so forth. Since there are many highly optimized native kernels developed by the HPC community over years, the invocation of HPC kernels can better leverage the hardware resource than the bottom level implementation of the original Harp applications, which uses Java threads to execute jobs in parallel. However, the invocation of DAAL kernels from Harp is a non-trivial work, and we must address two major problems: 1) Data conversion between Harp and DAAL, 2) Computation model of each application at both of inter-node and intra-node levels. Data conversion between Harp and DAAL is critical to applications like MF-SGD, where a massive model needs to be synchronized in each iteration. An inappropriate data conversion will significantly slow down the performance. Before discussing the data conversion, we first introduce the data structure of Harp library and Intel’s DAAL framework.

1) *Data Structures of Harp and Intel’s DAAL*: Harp has a three-level hierarchy of data structures. At the top level is the *Table* class, and a *Table* contains a certain number of *Partitions*. Each *Partition* consists of a *partition id* and a data container. Data containers could wrap up Java objects or primitive arrays, where the data is stored on heap memory managed by JVM. This design has two consequences: 1) The data stored within a *Table* is scattered into different partitions, whose memory addresses are not contiguous, while many native kernels with optimized memory access require contiguous memory allocation. 2) The memory addresses on Java heap can not be directly passed to native kernels

³<http://torch.ch>

⁴<https://www.tensorflow.org>

because JVM may change the objects’ physical addresses during Garbage Collection. DAAL library consists of two modules.

- *Native Kernels*, the implementation of algorithms and data structures in C/C++
- *Interface*, the API written in Java and Python to access *Native Kernels*

The native kernels are in charge of the computational work while the Java/Python APIs allow users to construct their applications without knowing the low-level implementations. In Harp-DAAL, data conversion happens between the Harp Java codes and the DAAL Java/Python APIs. Thus, the dataflow of the Harp-DAAL framework is from Harp Java codes to DAAL Java APIs and finally accessed by DAAL native kernels. For instance, the *HomogenNumericTable* provides users of two ways to store data.

- *ArrayImplm*, store the data at Java heap side
- *ByteBufferImplm*, store the data on the native side, use *DirectByteBuffer* to read and write data between Java and native side.

In the way of *ArrayImplm*, DAAL Java APIs hold data at Java heap space, and it creates a data structure called *JavaNumericTable* for the native kernels to access these data. Whenever a native kernel requests the data in *JavaNumericTable*, it will callback to a member function of *HomogenNumericTable* at Java side to trigger the data transfer via JNI functions and *DirectByteBuffer* class. Most of the DAAL’s data structures support this way of storing data. In the way of *ByteBufferImplm*, the Java APIs allocates an empty array in the native memory space, and it is the user’s responsibility to read and write data between Java heap memory and allocated native memory. In current DAAL release version 2017, only *HomogenNumericTable* supports this way of data storage. For native kernels, they can only manipulate data stored in the native memory space, and they use an auxiliary data structure named *MicroTable* to retrieve specified rows or columns from a *NumericTable*. We design two data conversion operations in Harp-DAAL. The first operation named *JavaBulkCopy* starts from Harp side and copies data from a Harp *Table* to a *HomogenNumericTable* via the *ByteBufferImplm*. The second operation named *NativeDiscreteCopy* starts from a native kernel and transfers data from DAAL Java APIs to native kernels via the *ArrayImplm*. Figure 2 (a) and (b) explain the workflow of the two data conversion operations.

- 2) *JavaBulkCopy*: There are two steps in *JavaBulkCopy*
 - Create a large *DirectByteBuffer*, and copy data from a Harp *Table* into the *DirectByteBuffer*
 - Write the content of *DirectByteBuffer* into a *HomogenNumericTable*.

For step one, we use *java.lang.Thread* class to do a parallel data copy from the non-contiguous JVM heap memory of a Harp *Table* to a contiguous memory block within

DirectByteBuffer. In step two, a member function named *releaseBlockOfRows* from *HomogenNumericTable* does a bulk data copy from *DirectByteBuffer* to the native memory allocated in a *HomogenNumericTable*. The advantage of *JavaBulkCopy* is that the data within *HomogenNumericTable* is contiguous, which favors many of DAAL’s algorithms such as K-means, and it usually results in an efficient cache usage on multi-core and many-core processors. However, it has a two-fold downside. Firstly, the maximal size of a single *DirectByteBuffer* is limited by 2GB. Secondly, the scheduling of Java threads is not as efficient as that of OpenMP and TBB, which increases the overhead of parallel data copy.

3) *NativeDiscreteCopy*: *NativeDiscreteCopy* could be applied to all sorts of *NumericTable* within DAAL, where the data is stored in Java heap memory and DAAL uses a C++ class named *JavaNumericTable* to expose the data to native kernels. In a *NativeDiscreteCopy* process, a thread from a native kernel will be attached to a C++ pointer, which is a member of *JavaNumericTable* and points to a JVM object via JNI. The thread can then call Java functions of *NumericTable* to copy data from Java heap memory back to native memory by using *DirectByteBuffer*. Since the JVM pointer allows concurrent accesses from different threads, we can use OpenMP or TBB threads to copy data from Java heap memory to native memory in parallel (See Figure 2 (b)). The advantage of using *NativeDiscreteCopy* is to reduce the size of *DirectByteBuffer*, because each thread can reuse the assigned buffer to copy different data. The downside of *NativeDiscreteCopy* is exposure of the JNI interface and low-level implementations to users. An ill-implemented *NativeDiscreteCopy* will cause issues such as memory leak. Deciding which memory copy operation to use depends on the data structures and model synchronization operations. We will discuss the details in Section III-B for each application.

B. Benchmark Algorithms

We select three typical learning algorithms to evaluate our framework: 1) K-means Clustering (K-means), a computation-bounded algorithm; 2) Matrix Factorization by Stochastic Gradient Descent (MF-SGD), a computation-bounded and communication-bounded algorithm; 3) Alternating least squares (ALS), a communication-bounded algorithm.

1) *K-means Clustering*: K-means is a widely used clustering algorithm in machine learning community. It computes the distance between each training point to every centroid, re-assigns the training point to the new cluster and re-computes the new centroid of each cluster at each iteration. Harp-DAAL-Kmeans is built upon Harp’s original K-means parallel implementation, with a computation complexity for each iteration as $\mathcal{O}(|\Omega|KM)$, where Ω is the set of training samples, K is the feature dimension of a training point, and

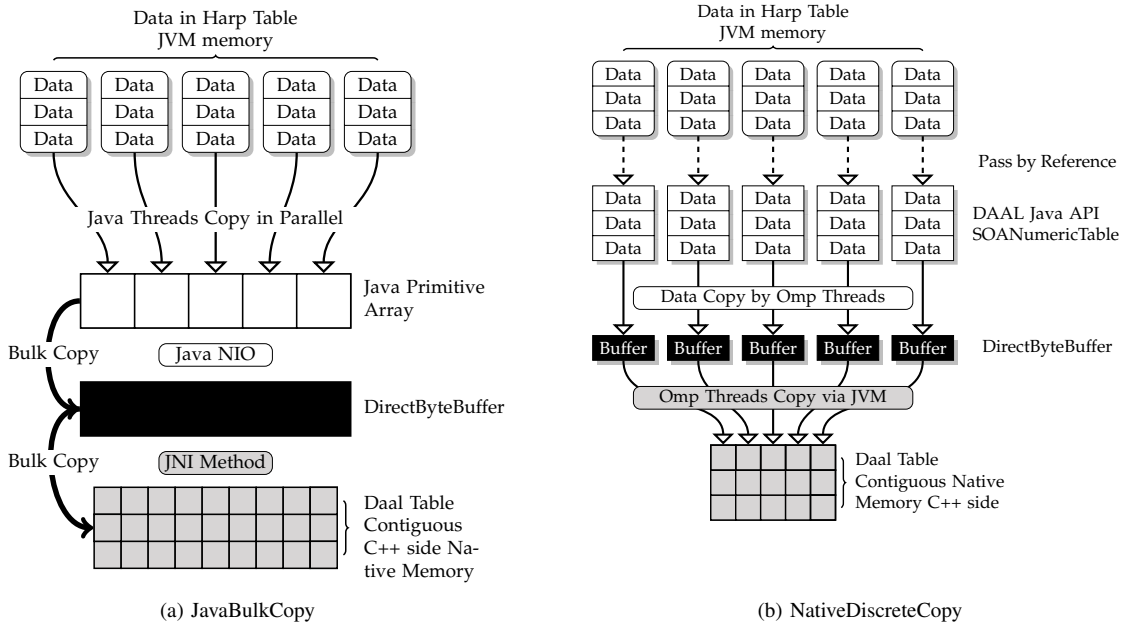


Figure 2. (a) JavaBulkCopy; (b) NativeDiscreteCopy

M is the number of centroids. Harp-DAAL-Kmeans uses the regroup-allgather operation [5] to synchronize model, i.e. centroids, among each mapper. Harp-DAAL-Kmeans uses DAALs K-means kernel, where the computation of point-centroid distance is implemented by BLAS-level 3 matrix-matrix operations. This optimization significantly increases the computation intensity, resulting in highly vectorized codes against the original Harp-Kmeans.

2) *MF-SGD*: Matrix Factorization based on Stochastic Gradient Descent (MF-SGD) is commonly used in recommender systems [6], where it aims to factorize a sparse matrix into two dense model matrices W and H . The computation complexity for each iteration is $\mathcal{O}(|\Omega|K)$, where Ω is the set of training samples, and K is the feature dimension of a training point. Previous work such as [7] concentrates on the single-node shared memory optimization for MF-SGD. For distributed memory system, we already implemented a pure Java version within the Harp framework [5], and we reimplement it by invoking native kernels in our hybrid Harp-DAAL framework. The implementation of MF-SGD consists of two levels. At the inter-node level, we use a rotation operation to synchronize model distributed among nodes [8]. At the intra-node level, we choose an asynchronous operation to update model manipulated by different threads. By eliminating the lock and waits at thread level, we relieve the load balancing problems caused by uneven distribution of training points in each row and column.

3) *ALS*: Alternating least squares (ALS) is another frequently used algorithm to decompose rating matrices in recommender systems. The algorithm has a computation

complexity for each iteration as $\mathcal{O}(|\Omega|K^2 + (m+n)K^3)$. Ω is the set of training samples, K is the feature dimension, m is the row number of the rating matrix, and n is the column number of the rating matrix. Unlike MF-SGD, ALS alternatively computes model W and H independently of each other. The implementation of ALS in our Harp-DAAL framework chooses the regroup-allgather operation to interface the DAAL-ALS kernels based on the work of Zou et al. [9]

IV. EXPERIMENTATION

In the experiments, we compare the performance of the following six implementations:

- Harp-DAAL-Kmeans
- Spark-Kmeans
- Harp-DAAL-SGD
- NOMAD-SGD
- Harp-DAAL-ALS
- Spark-ALS

Apart from the three applications from Harp-DAAL, we use Spark-Kmeans and Spark-ALS from Apache Spark in Section II, both are written in Java. NOMAD-SGD is a distributed MF-SGD developed by Yun et al. [10], which is written in C/C++ and uses MPI in inter-node communication.

A. Hardware Platform

We conduct experiments on a cluster of Intel’s Xeon Phi processor 7250 codenamed Knights Landing (KNL). Table I gives the specification of one KNL node. Compared to Intel’s

Table I
SPECIFICATION OF XEON PHI 7250 KNL

Cores		Memory		Node Spec		Misc Spec	
Cores	68	DDR4	190 GB	Network	Omni-path	Instruction Set	64 bit
Base Freq	1.4GHz	MCDRAM	16 GB	Peak Port Band	100 Gbps	IS Extension	AVX512
L1 Cache	2 MB	DDR4-Band	90 Gbps	Socket	1	Max Threads	271
L2 Cache	34 MB	MCDRAM-Band	400 Gbps	Disk	1 TB	VPU	136

Table II
DATASETS USE IN K-MEANS, MF-SGD, AND ALS

Dataset	Kmeans-Single	Kmeans-Multi	Movielens	Netflix	Yahoomusic	Enwiki	Hugewiki
#Training	5000000	20000000	9301274	99072112	252800275	609700674	3074875354
#Test	none	none	698780	1408395	4003960	12437156	365998592
#centroid	10000	100000	none	none	none	none	none
Dim	100	100	40	40	100	100	1000
λ	none	none	0.05	0.05	1	0.01	0.01
γ	none	none	0.003	0.002	0.0001	0.001	0.004

Xeon processor family, KNL has three advantages: 1) A high number of physical cores, 2) Up to 136 AVX-512 Vector Processing Units (VPU). Each VPU could simultaneously compute 8 double or 16 float operations in parallel by enabling Intel’s AVX-512 instruction set extension. 3) On-chip high bandwidth memory named MCDRAM, whose bandwidth reaches 400Gbps. Therefore, KNL favors applications that maximally leverage the intra-node threads parallelism, codes vectorization, and memory bandwidth usage. To compare the utilization of KNL’s features, we employ Intel VTune Amplifier to do the micro-benchmark profiling.

B. Dataset

A variety of datasets are used to examine the performance of implementations. Table II describes the details including sizes and parameters. For K-means, we use synthetic datasets, a small one for single node test and a large one for multi-node scaling test. For MF-SGD and ALS, we use same datasets as in related work.

V. RESULTS AND ANALYSIS

A. Single Node Performance

We first evaluate the performance of Harp-DAAL framework on a single KNL node. We use 64 cores out of the total 68 cores because the self-booted KNL node requires several cores to run the OS system. Each physical core runs one thread, which implies a total of 64 threads for the applications. The metric includes the execution time per training iteration for K-means, MF-SGD and ALS, respectively. The time is averaged over a certain number of iterations to avoid potential cold start phenomenon.

In Figure 3, we find that Harp-DAAL achieves the best performance among the three frameworks. For K-means and ALS, it runs 20x to 50x faster than Spark, this result is

not surprising because both of K-means and ALS contains matrix-matrix computations that benefit from the optimized MKL native kernels invoked by Harp-DAAL. For MF-SGD, we still achieve comparable performance to NOMAD on datasets Movielens, Netflix, and Enwiki. On Yahoomusic, we even have a 2x speedup, which proves that the intra-node optimization brought by DAAL kernels achieves the same performance level of the state-of-the-art work.

B. Multi-node Performance

To compare multi-node performance of the applications, we use 60 out of 64 threads per node for K-means and ALS. For MF-SGD, because extra threads may be used to do communication and overlap with the computation in a pipeline. The scalability for K-means and MF-SGD is not measured against one node, due to the datasets are too large to fit into the memory of single node. Instead, we assume that the scalability from one node to 10 nodes are linear, and measure the strong scalability over 10 nodes. Figure 6 decomposes the execution time of Harp-DAAL applications on multiple nodes into three components: 1) Computation time on local nodes, 2) Data conversion time on local nodes, and 3) Data communication time among remote nodes. For K-means, the computation time is dominant, taking more than 80% of the total execution time. When it scales from 10 nodes to 20 nodes, the communication ratio decreases because the model volume on each node decreases, and Harp regroup-allgather operation favors small communication data. The communication ratio slightly increases from 20 nodes to 30 nodes, which is due to the insufficient computation work on each node that causes a reduced computation ratio. In Figure 6 (a), Harp-DAAL-Kmeans runs 15x to 40x faster than Spark Kmeans, and shows better scalability from 10 nodes to 20. Beyond 20 nodes, due to low computation workload on local nodes, the scalability of Harp-DAAL-Kmeans drops while

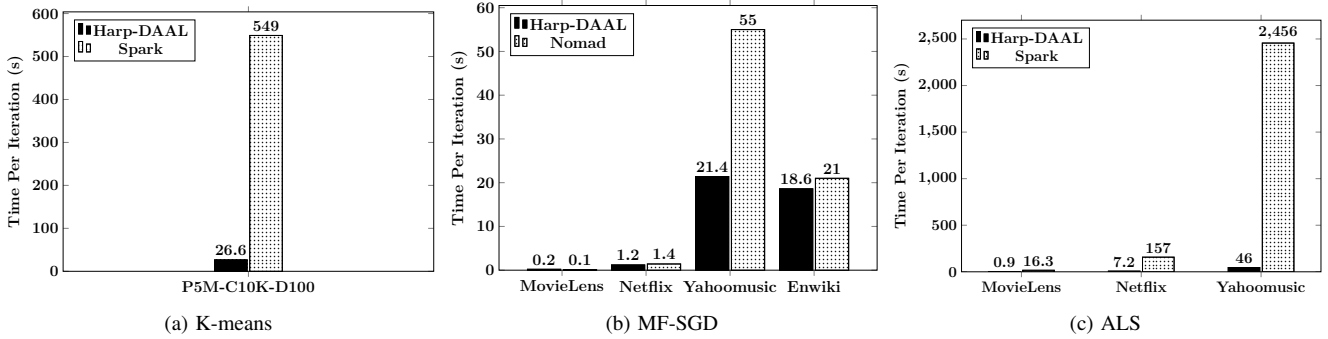


Figure 3. The execution time (per iteration) on a single KNL node. Each subgraph compares two implementations on the same algorithm with different datasets. (a) K-means dataset: 5 million points, 10 thousand centroids, 1000 feature dimension; (b) MF-SGD dataset: 1) MovieLens 9 million points, 40 feature dimension, 2) Netflix 100 million points, 40 feature dimension, 3) Yahooomusic 250 million points, 100 feature dimension, 4) Enwiki 600 million points, 100 feature dimension; (c) ALS dataset : 1) MovieLens 9 million points, 40 feature dimension, 2) Netflix 100 million points, 40 feature dimension, Yahooomusic 250 million points, 100 feature dimension

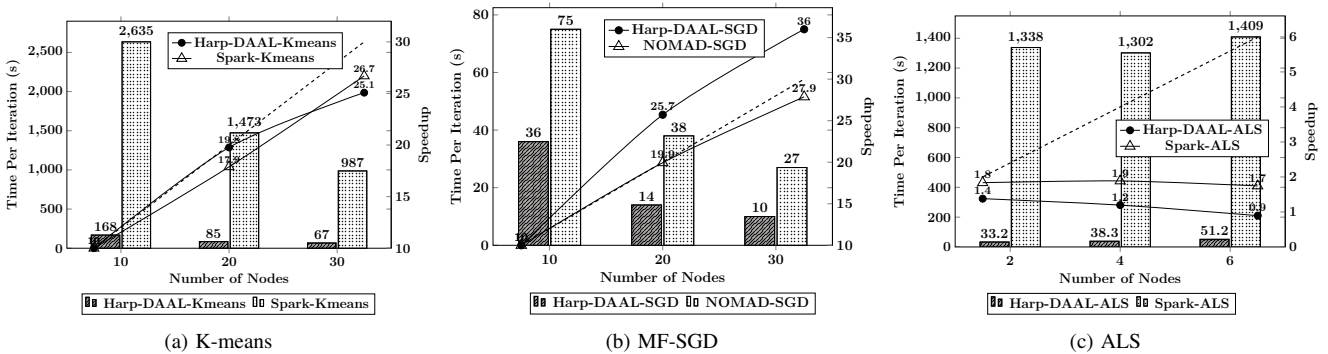


Figure 4. Strong Scaling on multiple KNL nodes (a) K-means dataset:20 million points, 100 thousand centroids, 100 feature dimension; (b) MF-SGD dataset, Hugelwiki, with 3 billion training points, 1000 feature dimension; (c) ALS dataset, Yahooomusic, with 250 million training points, 100 feature dimension

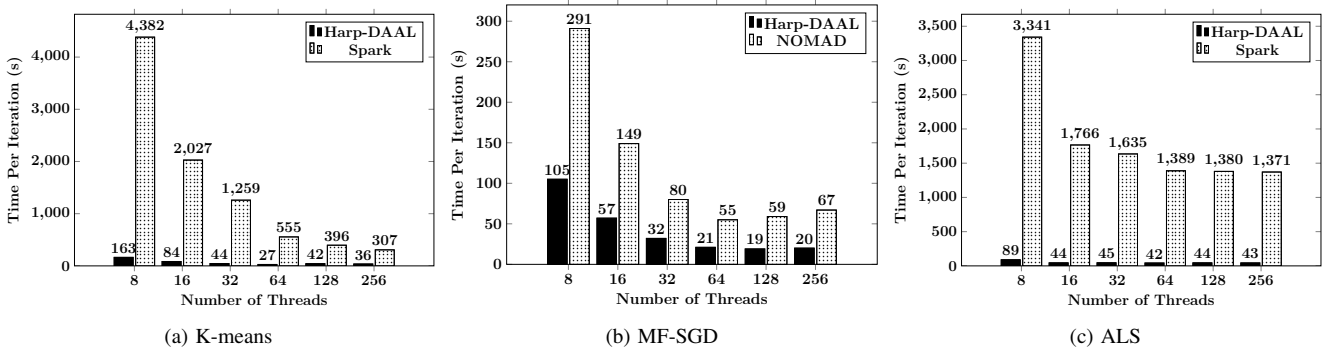


Figure 5. Strong Scaling on Threads of a single KNL (a) K-means dataset: 5 million points, 10 thousand centroids, 100 feature dimension; (b) MF-SGD dataset: Yahooomusic with 250 million training points, 1000 feature dimension; (c) ALS dataset, Yahooomusic, with 250 million training points, 100 feature dimension

Spark-Kmeans still has substantial computation workload. It suggests that Spark-Kmeans is more computation-bounded than Harp-DAAL-Kmeans because strong scalability reflects how an implementation is bounded by local computations. Since Harp-DAAL-Kmeans invokes fast MKL kernels at the

low level, it is much less bounded by computation time.

For MF-SGD, Figure 6 (b) shows that Harp-DAAL-SGD runs 2.5x faster than NOMAD-SGD, and it even achieves super-linear on 20 nodes and 30 nodes, which is also better than the scalability of NOMAD-SGD. The reason why Harp-

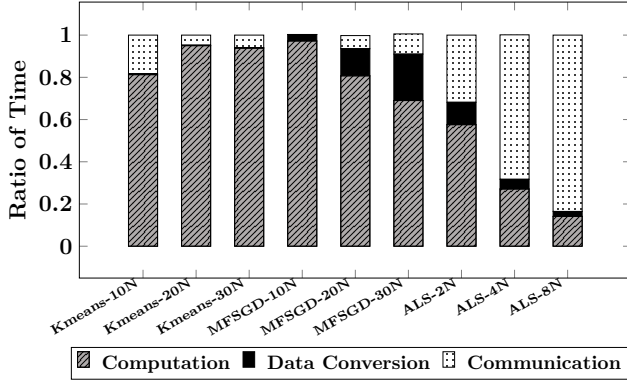


Figure 6. Breakdown of execution time of Harp-DAAL Applications on multiple KNL nodes 1) K-means dataset:20 million points, 100 thousand centroids, 100 feature dimension; 2) MF-SGD dataset, Hugewiki, with 3 billion training points, 1000 feature dimension; 3) ALS dataset, Yahoo music, with 250 million training points, 100 feature dimension

DAAL-SGD has a super-linear speedups has two-folds: 1) Harp-DAAL-SGD has its native kernels implemented by OpenMP, and a small number of training data may significantly reduce the scheduling overheads. 2) The native kernel for intra-node computation uses asynchronous shared-memory data access operations, which favors sparse training points that has less conflicts in accessing the same data.

For ALS, Figure 6 shows that the communication time already takes up more than 50% of the execution time on four nodes. It means that ALS is not bounded by local computation, and therefore both of Harp-DAAL-ALS and Spark-ALS have bad strong scalabilities in Figure 6 (c). However, Harp-DAAL-ALS still has 25x to 40x speedups to Spark-ALS because it invokes highly efficient MKL kernels.

C. Micro-Benchmark

1) Performance Breakdown on a KNL Single Node:

The execution time breakdown of all benchmarks in both frameworks is shown in Figure 7. Compared to another C++-based SGD implementation, Nomad-SGD, our SGD with Harp-DAAL significantly reduces execution time. Because Harp-DAAL-SGD is able to fully take advantage of AVX-512 to reduce retiring instruction number by packaging multiple floating point operations into one SIMD instruction. One AVX-512 instruction may cause multiple simultaneous L1 cache accesses, so that Harp-DAAL-SGD also improves the L1 cache bandwidth utilization. Compared to our C++ and Java hybrid framework, Harp-DAAL, the pure Java framework Spark inflates the executed instruction number by at least 10 times on K-means and ALS. Therefore, the retiring of benchmarks with Harp-DAAL only takes at most 10% of that with Spark. Spark is developed upon only Java virtual machine which can barely benefit from AVX-512. Therefore, its non-vectorized code only generates memory

requests with poor temporal locality hardly saturating the large bandwidth supplied by the MCDRAM. The serialized memory accesses of Spark framework substantially prolong the time stall along the memory hierarchy.

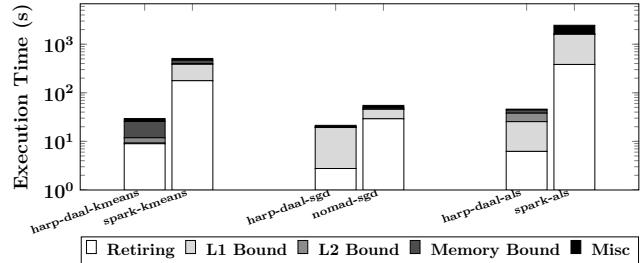


Figure 7. Breakdown of execution time on a single KNL node. X Bound indicates the time stall caused by X-event. Retiring means the execution time of instructions. Misc represents the time stall triggered by instruction cache miss, branch prediction miss, TLB miss and all other micro-architecture events.

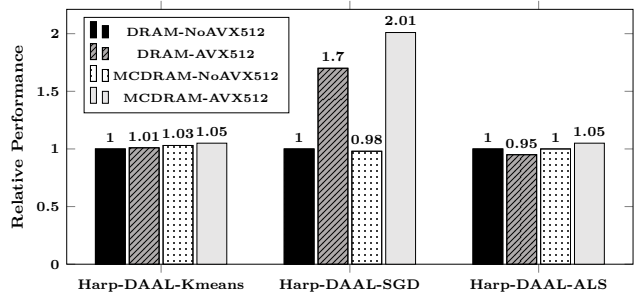


Figure 8. Relative Performance by enabling AVX-512 and MCDRAM. The baseline performance is the configuration without AVX-512 and MCDRAM, which is set to 1. The relative performance of other configuration is their acceleration of execution time compared with the baseline performance

2) *Thread Scaling*: Figure 5 shows the performance of K-means, MF-SGD and ALS with varying numbers of threads on one single node. With an increasing number of threads, total computing power of KNL boosts. However, communications between cores intensify and cache capacity per thread also drops significantly. Therefore, more threads do not necessarily indicate shorter execution time on a KNL node. K-means, MF-SGD and ALS encapsulated by Harp-DAAL achieve the best performance with 64, 128 and 64 threads, respectively. Since Spark cannot fully utilize AVX-512, both K-means and ALS with Spark have to execute more instructions to do the same job. Both prefer 256-thread configuration to relieve their bottleneck on instruction retiring.

3) *AVX-512 and MCDRAM*: Figure 8 exhibits the performance improvement achieved by AVX-512 and MCDRAM

in our Harp-DAAL framework. Instead of configuring MCDRAM as a hardware managed cache, we deployed it as a parallel component to DDR4 in main memory system. Through `numactl` command, we can use either DDR4-based main memory or MCDRAM-based main memory. All bars are normalized to the scheme compiled with disabling vectorization and with DDR4-based main memory. By enabling AVX-512, Kmeans and ALS do not reduce execution time significantly. This is because, although we compiled them without vectorization, kernels in Kmeans and ALS invoke functions from Intel Math Kernel Library which are fully optimized with AVX-512. AVX-512 improves the performance of SGD by 70%, since VPUs on KNL compute matrix multiplications with much higher throughput. Only enabling MCDRAM does not obviously boost the performance of all three benchmarks. K-means, non-vectorized SGD and ALS have relatively high L2 cache hit rate, so that they cannot benefit from large memory bandwidth provided by MCDRAM. By enabling both AVX-512 and MCDRAM, SGD improves the performance by 101%, since AVX-512 instructions may generate multiple memory accesses in one cycle and the MCDRAM memory bandwidth can be fully utilized.

VI. CONCLUSION

We design and implement Harp-DAAL that enables Hadoop on cloud servers with manycore KNL processors. Many machine learning applications can be implemented with MapReduce-like interfaces with significantly boosted performance by scaling up. Through evaluating computation and communication-bounded applications, we show that Harp-DAAL combines advanced communication operations from Harp and high performance computation kernels from DAAL. Our framework achieves 15x to 40x speedups over Spark-Kmeans and 25x to 40x speedups to Spark-ALS. Compared to NOMAD-SGD, a state-of-the-art C/C++ implementation of the MF-SGD application, we still get higher performance by a factor of 2.5. An interesting future direction will be to compare with other hardware including Haswell and GPU, and develop high performance machine learning libraries. The code and documentation of Harp-DAAL framework can be found at <https://dcspsidal.github.io/harp/>.

ACKNOWLEDGMENT

We gratefully acknowledge generous support from the Intel Parallel Computing Center (IPCC) grant, NSF OCI-114932 (Career: Programming Environments and Runtime for Data Enabled Science), CIF-DIBBS 143054: Middleware and High Performance Analytics Libraries for Scalable Data Science. We appreciate the support from IU PHI, FutureSystems team and ISE Modelling and Simulation Lab.

REFERENCES

- [1] B. He, W. Fang, Q. Luo, N. K. Govindaraju, and T. Wang, "Mars: a MapReduce framework on graphics processors," *International Conference on Parallel Architectures and Compilation Techniques*, pp. 260–269, 2008. [Online]. Available: <http://portal.acm.org/citation.cfm?doid=1454115.1454152>
- [2] X. Meng, J. Bradley, S. Street, S. Francisco, E. Sparks, U. C. Berkeley, S. Hall, S. Street, S. Francisco, D. Xin, R. Xin, M. J. Franklin, U. C. Berkeley, and S. Hall, "MLlib : Machine Learning in Apache Spark," *The Journal of Machine Learning Research*, vol. 17, pp. 1235–1241, 2016.
- [3] A. J. Awan, M. Brorsson, V. Vlassov, and E. Ayguade, "Performance Characterization of In-Memory Data Analytics on a Modern Cloud Server," *Proceedings - 2015 IEEE 5th International Conference on Big Data and Cloud Computing, BDCLOUD 2015*, pp. 1–8, 2015.
- [4] E. P. Xing, Q. Ho, W. Dai, J. K. Kim, J. Wei, S. Lee, X. Zheng, P. Xie, A. Kumar, and Y. Yu, "Petuum: A New Platform for Distributed Machine Learning on Big Data," *Kdd*, vol. XX, no. X, p. 15, 2015. [Online]. Available: http://www.cs.cmu.edu/~seunghak/petuum/_kdd15.pdf
- [5] B. Zhang, Y. Ruan, and J. Qiu, "Harp: Collective communication on Hadoop," *Proceedings - 2015 IEEE International Conference on Cloud Engineering, IC2E 2015*, pp. 228–233, 2015.
- [6] Y. Koren, R. Bell, and C. Volinsky, "Matrix Factorization Techniques for Recommender Systems," *Computer*, vol. 42, no. 8, pp. 42–49, 2009.
- [7] W.-S. Chin, Y. Zhuang, Y.-C. Juan, and C.-J. Lin, "A Fast Parallel Stochastic Gradient Method for Matrix Factorization in Shared Memory Systems," *ACM Transactions on Intelligent Systems and Technology*, vol. 6, no. 1, pp. 1–24, 2015. [Online]. Available: <http://dl.acm.org/citation.cfm?doid=2745393.2668133>
- [8] B. Zhang, B. Peng, and J. Qiu, "Model-centric computation abstractions in machine learning applications," *Proceedings of the 3rd ACM SIGMOD Workshop on Algorithms and Systems for MapReduce and Beyond - BeyondMR '16*, vol. 09, pp. 1–4, 2016. [Online]. Available: <http://doi.acm.org/10.1145/2926534.2926539>
- [9] R. P. Yunhong Zhou, Dennis Wilkinson, Robert Schreiber, "Large-Scale Parallel Collaborative Filtering for the Netix Prize," *Lecture Notes in Computer Science*, vol. 5034, pp. 337–347, 2010. [Online]. Available: <http://link.springer.com/content/pdf/10.1007/978-3-642-14355-7.pdf>
- [10] H. Yun, H.-F. Yu, C.-J. Hsieh, S. V. N. Vishwanathan, and I. Dhillon, "NOMAD: Non-locking, stochastic Multi-machine algorithm for Asynchronous and Decentralized matrix completion," *Pvldb*, vol. 7, no. 11, pp. 975–986, 2014. [Online]. Available: <http://www.vldb.org/pvldb/vol7/p975-yun.pdf>