

Performance Characterization of Multi-threaded Graph Processing Applications on Many-Integrated-Core Architecture

Lei Jiang Langshi Chen Judy Qiu

School of Informatics, Computing, and Engineering, Indiana University Bloomington
 {jiang60, lc37, xqiu}@indiana.edu

Abstract—In the age of Big Data, parallel graph processing has been a critical technique to analyze and understand connected data. Meanwhile, Moore’s Law continues by integrating more cores into a single chip in the deep-nano regime. Many-Integrated-Core (MIC) processors emerge as a promising solution to process large graphs. In this paper, we empirically evaluate various computing platforms including an Intel Xeon E5 CPU, an Nvidia Tesla P40 GPU and a Xeon Phi 7210 MIC processor codenamed Knights Landing (KNL) in the domain of parallel graph processing. We show that the KNL gains encouraging performance and power efficiency when processing graphs, so that it can become an auspicious alternative to traditional CPUs and GPUs. We further characterize the impact of KNL architectural enhancements on the performance of a state-of-the-art graph framework. We have four key observations: ① Different graph applications require distinctive numbers of threads to reach the peak performance. For the same application, various datasets need even different numbers of threads to achieve the best performance. ② Not all graph applications actually benefit from high bandwidth MCDRAMs, while some of them favor low latency DDR4 DRAMs. ③ Vector processing units executing AVX512 SIMD instructions on KNLs are underutilized when running the state-of-the-art graph framework. ④ The sub-NUMA cache clustering mode offering the lowest local memory access latency hurts the performance of graph benchmarks that are lack of NUMA awareness. At last, we suggest future works including system auto-tuning tools and graph framework optimizations to fully exploit the potential of KNL for parallel graph processing.

I. INTRODUCTION

Big Data explodes exponentially in the form of large-scale graphs. Graph processing has been an important technique to compute, analyze and visualize connected data. Real-world large-scale graphs [20], [27] include social networks, transportation networks, citation graphs and cognitive networks, which typically have millions of vertices and millions to billions of edges. Because of the huge graph size, it is natural for both industry and academia to develop parallel graph frameworks to accelerate graph processing on various hardware platforms. A large number of CPU frameworks, e.g., Giraph [24], Pregel [17] and GraphLab [16], have emerged for scalable in-memory graph processing. Recent works propose GPU frameworks such as CuSha [9], Medusa [28], LonestarGPU [4] and Gunrock [27] to perform high throughput graph processing on GPUs.

In CPU frameworks, graphs are partitioned, distributed and processed among multiple nodes [16], [17], [24] by message passing schemes or computed locally on a shared memory node [20]. Distributed graph processing is notorious for frequent communications between computations on each vertex or edge [20]. In a large-scale cluster, frequent communications are translated to a huge volume of messages across multiple nodes [16], [17], [24], seriously limiting the performance of graph processing. Even a single Mac-Mini SSD-based laptop potentially can outperform a medium-scale cluster for graph processing [12]. On a shared memory multi-core CPU node, communications are interpreted to loads and stores in memory hierarchy [20]. The core efficiency of a shared memory node is averagely $100\times$ higher than that in a cluster [22]. However, compared to GPUs, the graph computing throughput on CPUs is still constrained by the limited number of cores.

GPU frameworks capture graph processing parallelism by mapping and computing vertices or edges on thousands of tiny GPU cores [4], [27]. Although graph applications exhibit irregular memory access patterns, frequent thread synchronizations and data dependent control flows [4], GPUs still substantially boost the performance of graph processing over CPUs. However, the performance of graph applications on GPUs is sensitive to graph topologies. When traversing graphs with large diameter and small degree, GPUs exhibit poor performance due to the lack of traversal parallelism [27]. Moreover, in some commercial applications, simple operations, such as adding or deleting an edge in a graph, may cost a significant portion of the application execution time [20], but it is difficult for GPUs to support such basic operations.

Recently, Intel releases Xeon Phi MIC processors as an alternative to traditional CPUs and GPUs for the high performance computing market. A MIC processor delivers the GPU-comparable computing throughput and CPU-like sequential execution power. Compared to GPUs, a MIC processor can easily support basic graph operations, e.g., inserting/deleting a node/edge, since it is fully compatible with the X86 instruction set. In this paper, we focus on the performance characterization of multi-threaded graph applications on an Intel Xeon Phi (2nd generation) processor - KNL [21]. Our contributions are summarized as follows:

- We empirically evaluate a state-of-the-art graph framework on three hardware platforms: an Intel Xeon E5 CPU, an Nvidia Tesla P40 GPU and a KNL MIC processor. We show that the KNL demonstrates encouraging performance and power efficiency when running multi-threaded graph applications.
- We further characterize performance details of multi-threaded graph benchmarks on KNL. We measure and analyze the impact of KNL architectural enhancements such as many out-of-order cores, simultaneous multi-threading, cache clustering modes, vectorization processing units and 3D stacked MCDRAM on parallel graph processing. We have four key observations: ❶ Different graph applications require distinctive numbers of threads to achieve their best performance. For the same benchmark, various datasets ask for different numbers of threads to attain the shortest execution time. ❷ Not all graph applications actually benefit from high bandwidth MCDRAMs, while some of them favor low latency DDR4 DRAMs. ❸ VPU's executing AVX512 SIMD instructions on KNLs are underutilized when processing graphs. ❹ The sub-NUMA cache clustering mode offering the lowest local memory access latency hurts the performance of graph benchmarks that are lack of NUMA awareness.
- For future work, we suggest possible system auto-tuning tools and framework optimizations to fully exploit the potential of KNL for multi-threaded graph processing applications.

II. BACKGROUND

A. Graph Processing Frameworks and Benchmark Suite

The rapid and wide deployment of graph analytics in real-world diversifies graph applications. A lot of applications such as breadth first search incorporate graph traversals, while other benchmarks such as page rank involve intensive computations on vertex properties. Though low-level hardwired implementations [1], [3], [7], [8] have demonstrated high computing efficiency on both CPUs and GPUs, programmers need high-level programmable frameworks to implement a wide variety of complex graph applications to solve real-world problems. Therefore, graph processing heavily depends on specific frameworks composed of data structures, programming models and basic graph primitives to fulfill various functionalities. Previous research efforts propose many graph frameworks on both CPUs [16], [20] and GPUs [4], [9], [20], [28] to hide complicated details of operating on graphs and provide basic graph primitives. Most workloads spend 50% ~ 76% of execution time within graph frameworks [20].

The graph applications we studied in this paper are from a state-of-the-art graph benchmark suite, *graphBIG* [20]. The suite is implemented with IBM *System-G* framework that is a comprehensive graph library used by many industrial

solutions [23]. The framework adopts the vertex centric programming model, where a vertex is a basic element of a graph. Properties and outgoing edges of a vertex are attached to the same vertex data structure. The data structure of all vertices is stored in an adjacency list, and outgoing edges inside a vertex data structure also form an adjacency list of edges. The *major* reason to choose *graphBIG* is that each benchmark in this suite has a CPU OpenMP version and a GPU CUDA version sharing the same data structure and vertex centric programming model provided by *System-G*. In this way, we can minimize the impact of differences between CPU and GPU graph processing frameworks, and focus on only differences between various hardware computing platforms. On the contrary, other graph suites [1], [3], [4] implement graph applications without using a general framework or by using different frameworks on CPUs and GPUs respectively. Particularly, as one of the best-known graph benchmark suites, *Graph500* [19] provides only limited number of CPU workloads, since it is designed for only CPU-based system performance ranking.

B. Target Graph Benchmarks

Graph applications exhibit diversified program behaviors. Some workloads, e.g., breadth first search, are traversal-based. Only a subset of vertices is typically active at a given point during the execution of this type of applications. They often introduce many memory accesses but limited arithmetic operations. Their irregular memory behavior results in extremely poor locality in memory hierarchy. Some graph benchmarks operating on rich vertex properties, e.g., page rank and triangle count. They incorporate heavy arithmetic computations on vertex properties and intensive memory accesses leading to hybrid workload behaviors. Most vertices are active in all stages of their executions. We selected, compared and studied six common graph benchmarks in this paper. We introduce their details as follows:

- **Breadth First Search** traverses a graph by starting from a root vertex and exploring its neighbors first.
- **K-Core** finds a maximal connected sub-graph where all vertices have a $\geq K$ degree. It follows Matula & Beck's algorithm [20].
- **Single Source Shortest Path** calculates the minimum cost paths from a root vertex to each vertex in a given graph by Dijkstra's algorithm [20].
- **Graph Coloring** partitions a graph into independent sets. One set contains vertices sharing the same color. It adopts Luby-Jones' algorithm [20].
- **Page Rank** uses the probability distribution to compute page rankings. The rank of each vertex is computed by $PR(u) = \sum_{N_u} \frac{PR(v)}{E(v)}$, where the *PR* of a vertex (*u*) is decided by the *PR* of each vertex (*v*) in its neighbor set (N_u) divided by its outgoing edge number $E(v)$.
- **Triangle Count** measures the number of triangles that are formed in a graph when three vertices are connected to each other (Schank's algorithm [20]).

C. Graph Topology

The performance of graph applications is heavily influenced by the topology of graph datasets. We studied how topologies impact the benchmark performance via following metrics. The *eccentricity* $\epsilon(v)$ of a vertex v in a given connected graph G is the maximum graph distance between v and any other vertex u of G . The *diameter* d of a graph is the maximum eccentricity of any vertex in the graph ($d = \max_{v \in V} \epsilon(v)$). The *Fiedler eigenvalue* of the graph ($F(G)$) is the second smallest eigenvalue of the Laplacian matrix of G . The magnitude of $F(G)$ exhibits how well G is connected. For traversal-based applications such as breadth first search, traversal iteration number is proportional to the eccentricity and Fiedler eigenvalue. The *vertex degree* indicates the number of edges connected to a vertex. The average vertex degree and the vertex degree distribution of a graph reflect the amount of parallelism.

D. Intel Xeon Phi Architecture

1) *Overall Architecture*: Xeon Phi MIC processors are fully compatible with the x86 instruction set and hence support standard parallel shared memory programming tools, models and libraries such as OpenMP. The KNL [21] is the 2nd generation architecture (14nm) and has been adopted by several data centers such as the national energy research scientific computing center [2].

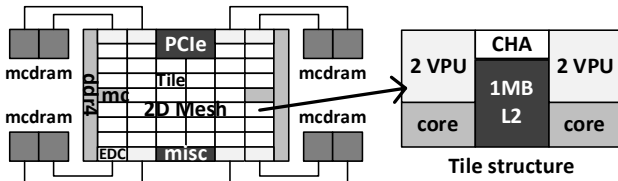


Figure 1. KNL (MC: DDR4 controller; EDC: MCDRAM controller).

The detailed KNL architecture is shown in Figure 1. It is built by up to 72 Atom (Silvermont) cores, each of which is based on a low operating frequency 2-wide issued out-of-order (OoO) micro-architecture supporting four concurrent threads, a.k.a, simultaneous multithreading (SMT). Additionally, every core has two vector processing units (VPUs) that support SIMD instructions such as AVX2 and AVX512 (a new 512-bit advanced vector extension of SIMD instructions for the x86 instruction set). A VPU can execute up to 16 single precision operations or 8 double precision floating point operations in each cycle. Two cores form a tile sharing a 1MB 16-way L2 cache and a caching home agent (CHA), which is a distributed tag directory for cache coherence. All tiles are connected by a 2D Mesh network-on-chip (NoC).

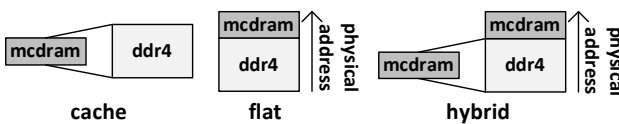


Figure 2. Configurations on MCDRAM.

2) *MCDRAM*: The KNL main memory system supports up to 384GB of DDR4 DRAM and 8~16GB of 3D stacked MCDRAM. The MCDRAM significantly boosts the memory bandwidth, but the access latency of MCDRAM is actually longer than that of DDR4 DRAM [18]. As Figure 2 shows, MCDRAM can be configured as a parallel memory component to DDR4 DRAM in main memory (*flat* mode), a hardware-managed L3 cache (*cache* mode) or both (*hybrid* mode). The *flat* mode offers high bandwidth by MCDRAM and low latency by DDR4 DRAM. However, programmers have to track the location of each data element and manage software complexity. On the contrary, the *cache* mode manages MCDRAM as a software-transparent direct-mapped cache. The *hybrid* mode combines the other two modes.

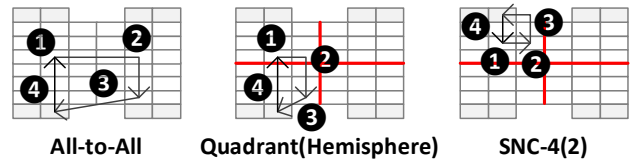


Figure 3. Configurations on Cache Clustering.

3) *Cache Clustering Mode*: Since all KNL tiles are connected by a Mesh NoC where each vertical/horizontal link is a bidirectional ring, all L2 caches are maintained coherent by the MESIF protocol. To enforce cache coherence, KNL has a distributed cache tag directory organized as a set of per-tile tag directories (shown as CHA in Figure 1) that record the state and location of all memory lines. KNL identifies which tag directory records a memory address by a hash function. As Figure 3 shows, the KNL cache can be operated in five modes including *All-to-All*, *Hemisphere*, *Quadrant*, *SNC-2* and *SNC-4*. When ① a core confronts a L2 miss, ② it sends a request to look up a tag directory. ③ The directory finds a miss and transfers this request to a memory controller. ④ At last, the memory controller fetches data from main memory and returns it to the core enduring the L2 miss. In the *All-to-All* mode, memory addresses are uniformly hashed across all tag directories. During a L2 miss, the core may send a request to a tag directory physically located in the farthest quadrant from the core. After the directory finds a miss, it may also send this request to a memory controller located in a third quadrant. Therefore, a L2 miss may have to traverse the entire Mesh to read one line from main memory. The *Quadrant (Hemisphere)* mode divides a KNL chip into four (two) parts. During a L2 miss, the core still needs to send a request to every on-chip tag directory, but the data associated to the target tag directory must be in the same part that the tag directory is located. Memory accesses from a tag directory are managed by its local memory or cache controller. The *Hemisphere* and *Quadrant* modes are managed by hardware and transparent to Operating System (OS). In contrast, the sub-NUMA-clustering (*SNC-2/4*) also separates the chip into two or four parts and exposes each as a separate NUMA domain to the OS. During a L2 miss, the

core only needs to send a request to its local tag directory that also transfers this request to the local memory controller if there is a directory miss. Therefore, the *SNC-4* mode has the lowest local memory access latency, but longer memory latency than that of the *Quadrant* mode when a request crosses NUMA boundaries. This is because requests traveling to another NUMA region have to be managed by the OS or applications themselves.

III. RELATED WORK

A recent graph characterization work on CPUs [3] identifies that eight fat OoO Ivy Bridge cores fail to fully utilize off-chip memory bandwidth when processing graphs, since the small instruction window cannot produce enough outstanding memory requests. In contrast, the KNL MIC processor can saturate main memory bandwidth by > 64 small OoO cores, each of which supports SMT and has two VPU. More research efforts [4] analyze GPU micro-architectural details to identify bottlenecks of graph benchmarks. In this paper, we pinpoint inefficiencies on the KNL MIC architecture when running multi-threaded graph applications. Previous physical-machine-based works [6], [8] find that the first generation Xeon Phi, Knight Corner, has poor performance for graph applications due to its feeble cores. Although the latest simulator-based work [1] characterizes the performance of multi-threaded graph applications on a MIC processor composed of 256 single-issue in-order cores, it fails to consider the architectural enhancements offered by KNLs, e.g., OoO cores, SMT, cache clustering modes, VPUs and MCDRAM. To our best knowledge, this work is the first to characterize and analyze the performance of multi-threaded graph applications on the KNL MIC architecture.

IV. EXPERIMENTAL METHODOLOGY

Graph benchmarks. In this paper, we adopted and studied six benchmarks including breadth first search (BF), single source shortest path (SS), graph coloring (GC), k-core (KC), triangle count (TC) and page rank (PR) from graphBIG [20]. Each benchmark contains a CPU OpenMP version and a GPU CUDA version. Due to the absence of PR GPU code in graphBIG, we created a vertex-centric GPU PR implementation based on a previous GPU PR benchmark [4] and the System-G framework. We compiled codes on CPU by `icc (17.0.0)` with `-O3, -MIC-AVX-512` and Intel OpenMP library. Through the command `numactl`, we chose to run CPU benchmarks in DDR4 DRAM under the flat mode. On GPU, we compiled programs by `nvcc (V8.0.61)` with CUDA-8.0. Since the size of large graph datasets exceeds the maximum physical GPU memory capacity, we used the `cudaMallocManaged` function to allocate data structures into a CUDA-managed memory space where both CPUs and GPUs share a single coherent address space. In this way, CUDA libraries transparently manage GPU memory access, data locality and data migrations between CPU and GPU. The OS we used is Ubuntu-16.04.

Table I
BENCHMARK DATASETS.

Dataset	Abbr.	Vertices	Edges	\overline{Degree}	\overline{Iter}_{BF}
roadNet_CA	road	1.97M	5.53M	2.81	665
roadNet_TX		1.38M	3.84M	2.27	739
soc-Slashdot0811	social	0.08M	0.9M	11.7	18
ego-Gplus		0.11M	13.6M	127.1	6
delaunay_n18	delaunay	0.26M	1.57M	6.0	223
delaunay_n19		0.52M	3.15M	6.0	296
kron_g500-logn17	kron	0.11M	10.23M	78.0	3
kron_g500-logn18		0.21M	21.17M	80.7	3
twitter7	large	17M	0.48B	4.6	23
com-friendster		65.6M	1.81B	27.6	15

Graph datasets. The graph inputs for our simulated benchmarks are summarized in Table I, where \overline{Degree} denotes the average vertex degree of the graph and \overline{Iter}_{BF} describes the average BF iteration number computed by searching from 10K random root vertices. *roadNet_XX* (road) are real-world road networks where most vertices have an outgoing degree below 4. *soc-Slashdot0811* and *ego-Gplus* (social) are two social networks typically having a scale-free vertex degree distribution and a small diameter. *delaunay* (delaunay) datasets are Delaunay triangulations of random points in the plane that also have extremely small outgoing degrees. Like social networks, *kron* (kron) datasets have large average vertex outgoing degree and can be traversed by only several BF iterations. We also included two large graph datasets (large), *twitter7* and *com-friendster*, each of which contains millions of vertices and billions of edges. The size of graph datasets is mainly decided by the number of edges, because the average degree of graphs we chose is ≥ 2 . We selected graph datasets from the Stanford SNAP [13]. We also used the synthetic graph generator, PaRMAT [10], in dataset sensitivity studies.

Table II
MACHINE CONFIGURATIONS.

Hardware	Description
Intel Xeon E5-4655v4	launched in Q2'16, 8-core@2.5GHz, 2T/core, 3.75MB L3/core, 512GB 68GB/s DDR4 DRAM thermal design power (TDP) 135W
Nvidia GTX1070	launched in Q2'16, 1920-cudaCore@1.5GHz, peak SP perf: 5783 GFLOPS, 8GB 256GB/s GDDR5, PCIe 3.0 \times 16 to CPU, TDP 150W
Nvidia Tesla P40	launched in Q3'16, 3840-cudaCore@1.3GHz, peak SP perf: 10007 GFLOPS, 24GB 346GB/s GDDR5, PCIe 3.0 \times 16 to CPU, TDP 250W
Intel Xeon Phi 7210 (KNL)	launched in Q2'16, 64-core@1.3GHz, 4T/core, 0.5MB L2/core, peak SP perf: 5324 GFLOPS, 8 channels 16GB 400GB/s 3D MCDRAM, 512GB 102GB/s DDR4 DRAM, TDP 230W
Disk	4TB SSD NAND-Flash

Hardware platforms. We chose and compared three server-level hardware platforms including a Xeon E5 CPU, a Tesla P40 GPU and a Xeon Phi 7210 (KNL) MIC processor. Since, compared to KNL, Tesla P40 achieves almost double peak single point float (SPF) throughput, we also included an Nvidia GTX1070 GPU having similar peak SPF throughput in our experiments. The machine configurations are shown in Table II. Both Xeon E5 and KNL have a 512GB DDR4 main

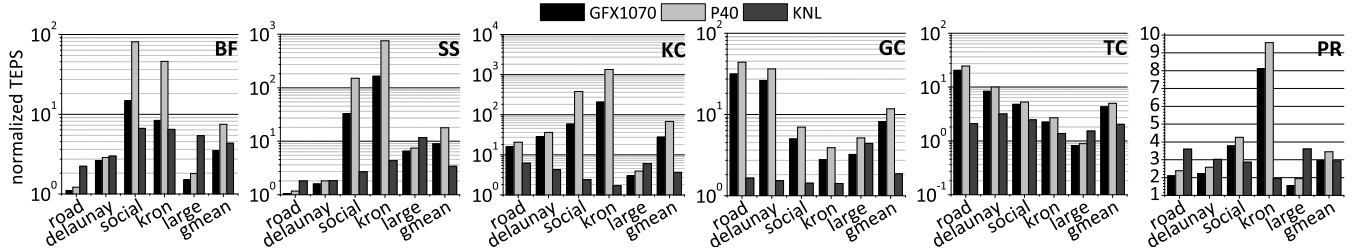


Figure 4. Performance comparison between Xeon E5-4655v4, GeForce GTX1070, Tesla P40 and KNL (normalized to Xeon E5-4655v4).

memory system, while P40 and GTX1070 are also hosted in the Xeon E5 machine with 512GB DDR4 DRAM. The thermal design power (TDP) values of Xeon E5, GTX1040, P40 and KNL are 135W, 150W, 250W and 230W. TDP represents the maximum amount of heat generated by a hardware component that the cooling system can dissipate under any workload.

Profiling tools. We adopted Intel VTune Analyzer to collect architectural statistics and likwid-powermeter [26] to profile power consumption on CPU and KNL. The likwid-powermeter is a tool for accessing RAPL registers on Intel processors, so it is able to measure the power consumption of both CPU package and DRAM memories. The register reading interval we set in all experiments is $1ms$. For GPUs, we used Nvidia visual profiler (*nvprof*) to collect both performance and power characterization results. The performance results we reported in Section V do not include the overhead of profiling tools on CPU, GPU and KNL.

Metrics. We measured the performance of graph applications as *Traversed Edges Per Second* (TEPS) [19]. The TEPS for non-traversing-based benchmarks, e.g., PR, is computed by dividing the number of processed edges by the mean time per iteration, since each vertex processes all its edges in one iteration [19]. We aim to understand the performance and power consumption of multi-threaded graph application kernels, and thus all results ignore disk-to-memory and host-to-device data transfer time during application initializations. But we measured page faults on CPU and data transfers between CPU and GPU inside application kernels.

V. EVALUATION

We cannot answer which type of hardware is the best for shared memory parallel graph processing, since the platforms we selected have completely different power budgets and hardware configurations such as operating frequencies, core numbers, cache sizes and micro-architecture details. By allocating more power budgets and hardware resources, theoretically speaking, any type of platform can possibly outperform the others. Our purpose of the comparison between KNL and other types of hardware is to demonstrate the KNL MIC processor achieves encouraging performance and power efficiency, and thus it can become one of the most promising alternatives to traditional CPUs and GPUs when processing graphs.

To investigate the KNL performance of parallel graph processing, we empirically analyze the performance comparison between four hardware platforms. And then, we characterize performance details of six graph benchmarks on KNL. We explain the impact of KNL architectural innovations such as many OoO cores, SMT, VPU and MCDRAM on multi-threaded graph applications.

A. Performance Comparison

The performance comparison between four hardware platforms is shown in Figure 4, where all results are normalized to the performance (TEPS) of Xeon E5 CPU. We configured KNL cache clustering in the Quadrant mode and MCDRAM in the cache mode. The KNL performance reported in this section is achieved by its optimal thread configuration including thread number and thread affinity. We also performed an exhaustive search on all possible configurations on CPU and GPU to attain and report the best performance.

The primary weakness of GPU is the load imbalance problem introduced by its sensitivity to graph topologies in traversal-based graph applications. For BF, GPUs achieve better performance than KNL on graphs with small diameter and large average vertex degree, e.g., social and kron. Huge traversal parallelism exists in these graphs, and hence, the more cores one platform has, the better performance it can achieve. In this case, KNL is slower than GTX1070, although they have similar peak SFP computing throughput. This is because graph traversal operations can barely take advantage of VPUs on KNL. However, compared to KNL, GPUs process less edges when traversing graphs with large diameter and small average vertex degree, e.g., road and delaunay. There is little traversal parallelism in these graphs, so the OoO cores of KNL prevail due to their powerful sequential execution capability. When processing large graphs, GPUs suffer from frequent data migrations between CPU and GPU memories, due to their limited GDDR5 memory capacity. Therefore, KNL shines on these large graph datasets. P40 outperforms GTX1070 on all graph inputs, because it has more CUDA cores and a larger capacity GDDR5 memory system. As another traversal-based graph benchmark, SS shares the same performance trend as that of BF.

KC and GC are two applications operating on graph structures. Their computations involve a huge volume of

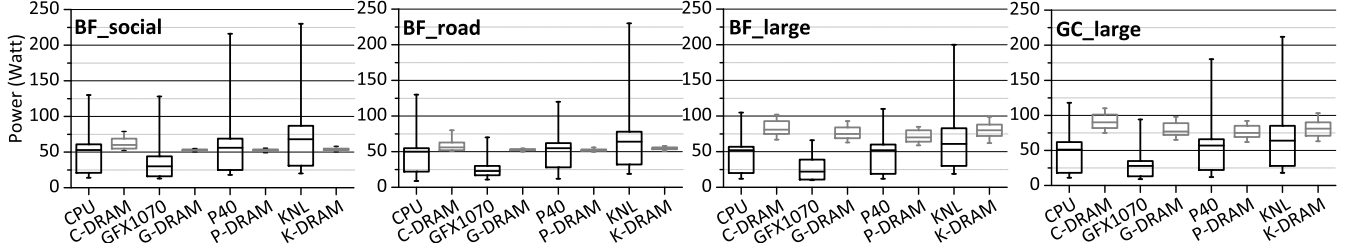


Figure 5. Power comparison between Xeon E5-4655v4, GeForce GTX1070, Tesla P40 and KNL.

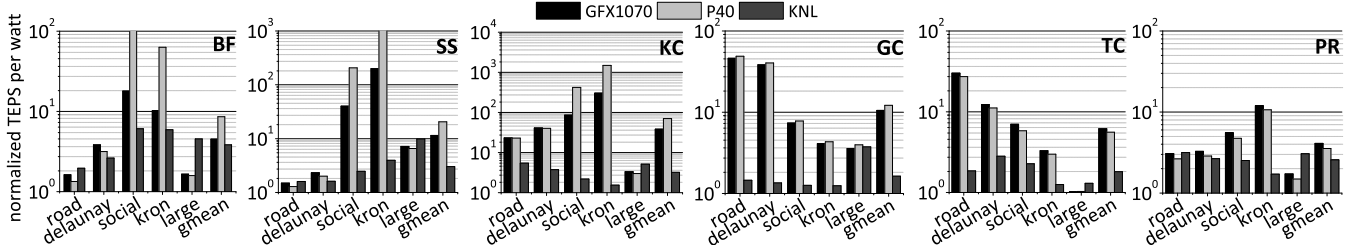


Figure 6. Performance per watt comparison between Xeon E5-4655v4, GeForce GTX1070, Tesla P40 and KNL (normalized to Xeon E5-4655v4).

basic arithmetic operations, e.g., integer in/decrementing, on vertices or edges. A large number of GPU cores support these massive arithmetic operations better than a small number of KNL cores. GPUs perform fine-grained scheduling on warps, while each KNL core can switch between only four threads, each of which cannot be fully vectorized. Therefore, the GPU performance is generally $10\times$ better than that of KNL. For GC, the improvement brought by thousands of CUDA cores on P40 even mitigates the data transfer penalty due to its limited GDDR5 memory capacity, i.e., P40 improves the processing performance over KNL by 16.7% on large graph datasets.

As graph benchmarks computing on vertex properties, TC and PR exhibit hybrid workload behaviors, since they require not only graph traversals but also relatively heavy arithmetic operations on vertex properties, e.g., multiply-accumulate operations. In TC, arithmetic operations dominate the performance of graph processing, so compared to KNL, GPUs boost the application performance by 112%~145% averagely. In PR, graph traversals dominate on graphs with large diameter and small average vertex degree, e.g., road and delaunay, so KNL obtains better edge processing speed on these graphs. For almost all combinations of graph benchmarks and datasets, Xeon E5 achieves the worst performance, i.e., only 1%~33.9% of CPU performance or 8%~50% of KNL performance. Only for TC, GPUs are slightly slower than Xeon E5 on large graph datasets, due to large penalties of data transfers on PCIe links between CPU and GPU memories.

B. Power and Performance Per Watt Comparison

The details of power characterization is shown in Figure 5, where we list power consumptions of four types of hardware and their 512GB DRAM main memories represented by X -DRAM (X can be Xeon E5 {C}PU, {G}FX1070, {P}40

and {K}NL.). The KNL power includes the power of MC-DRAM, while the GPU power contains the power of its GDDR5 memory.

For social graph datasets, all platforms are able to fully occupy almost 100% hardware resources and approach their peak power in a short time, because large traverse parallelism exists in these datasets. Only CPU DRAM is heavily accessed during BFs, while other platforms barely access their 512GB main memories. This is because both GPUs has their own GPU memories and KNL has a 16GB MCDRAM-based cache. As a result, only CPU DRAM achieves $> 60W$ average power consumption. For road graph datasets, the largest power spent by P40 (GFX1070) during searches is only around 48% (46%) of its TDP, due to the lack of traverse parallelism in these datasets. The average power consumption of two GPUs when processing road datasets decreases by 16%~24% over that dissipated during traversing on social datasets. In contrast, compared to social graphs, CPU and KNL slightly decrease their average power consumption by only 6%~9%, and the power consumption of CPU DRAM decreases by 10% averagely when processing road graphs. Compared to GPUs, CPU and KNL are insensitive to topologies of graph datasets due to their limited number of cores. Therefore, they have more consistent power results throughout various graph inputs. For large graph datasets, GPUs barely approach their TDPs due to frequent data transfers through PCIe links. CPU and KNL also suffers from frequent page faults, and thus their largest power values are only around 70%~80% of their TDPs. However, compared to other datasets, DRAMs in all platforms significantly boost the power consumption by 39%~62% when dealing with large graph inputs. SS shares a similar power consumption trend with BF. The power consumption of TC, KC, GC and PR with small datasets is similar to BF with social datasets, since all platforms

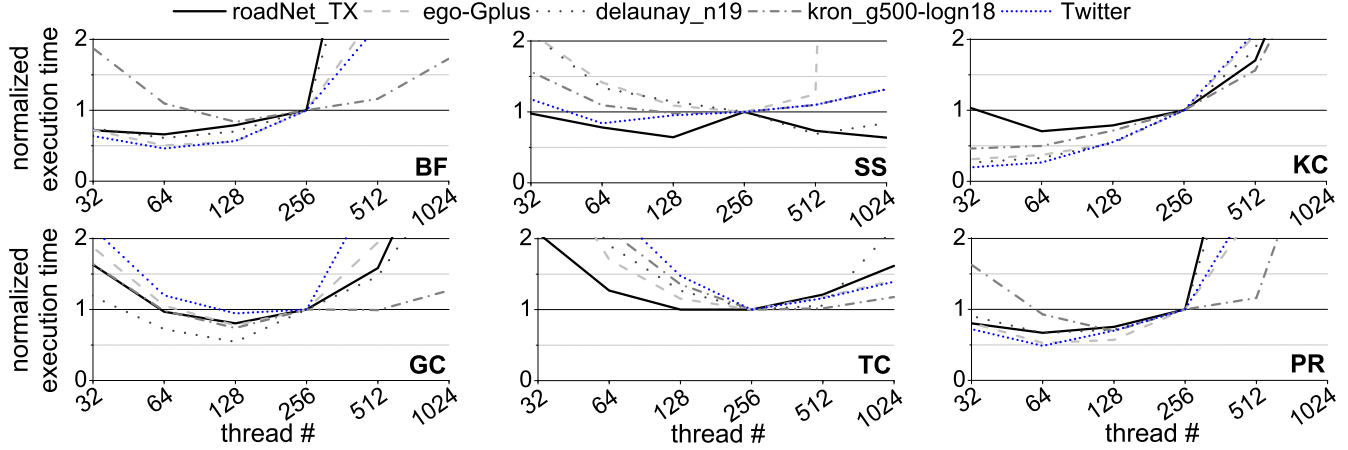


Figure 7. The KNL performance with varying thread numbers (normalized to the performance achieved by 256 threads).

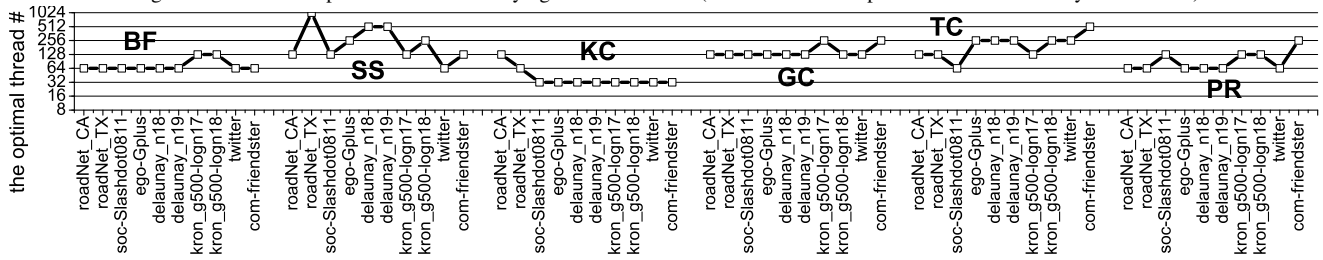


Figure 8. The optimal thread number for all benchmarks with all datasets

enjoy similarly large graph processing parallelism and small datasets do not need to frequently visit DRAMs. But when running TC, KC, GC and PR with large datasets, DRAMs spend 4%~11% more power than that when processing large datasets in BF. This is because graph applications operating on both structures of large graphs and a huge volume of vertex properties send more intensive memory accesses to DRAMs than traversal-based benchmarks.

The performance per watt comparison among all hardware platforms is exhibited in Figure 6, where all results are normalized to the performance per watt result of Xeon E5 CPU. Compared to CPU, P40 (GFX1070) achieves 2.5~69× (3~37×) better performance per watt averagely. On average, KNL obtains only 60%~285% better performance per watt over CPU. But for large graph datasets, compared to GPUs, KNL improves the performance per watt by 74%~186% when running BF and PR, and has similar results on performance per watt when running other benchmarks. Frequent data transfers between GPU and CPU memories caused by the limited capacity of GDDR5 memory waste a large amount of power and slow down applications on GPUs when processing large graphs. When running GC, TC and PR, although P40 is faster than GFX1070, the power consumption of P40 is also larger. Therefore, the performance per watt of P40 is not significantly better than that of GFX1070 for these benchmarks.

C. Threading

1) *Thread Scaling*: We show the graph application performance with varying OpenMP thread numbers on KNL in

Figure 7, where all performance results are normalized to the performance achieved by 256 threads. 256 (4 threads × 64 cores) is the KNL logic core number. Due to the limited figure space, we selected only the larger dataset in each graph input category. With a small number of threads (< 32), there are not enough working threads to do the task, and thus no benchmark obtains its best performance. When having more threads, the performance of almost all benchmarks improves more or less. However, when the thread number goes beyond 256, the execution time significantly rises, since the thread synchronization overhead dominates and degrades the performance of most benchmarks. The peak benchmark performance is typically achieved by 32 ~ 512 threads.

SMT and oversubscription allow multiple independent threads of execution on a core to better utilize hardware resources on that core. To implement SMT, some hardware sections of the core (but not the main execution pipeline) are duplicated to store architectural states. When one thread is stalled by long latency memory accesses, SMT stores its state to backup hardware sections and switches the core to execute another thread. SMT transforms one physical KNL core to four logic cores, each of which supports one thread. Some applications with certain datasets, e.g., TC with ego-Gplus and SS with kron_g500-logn18, fulfill their best performance by SMT (256 threads). In contrast, oversubscription requires the assistance from software such as OS or OpenMP library to switch threads, when the running thread is stalled. For applications suffering from massive concurrent cache misses, like SS with roadNet_TX and delaunay_n19,

oversubscription supports more simultaneous threads and outruns SMT. When running these applications, compared to the penalty of long latency memory accesses, the OS context switching overhead is not significant.

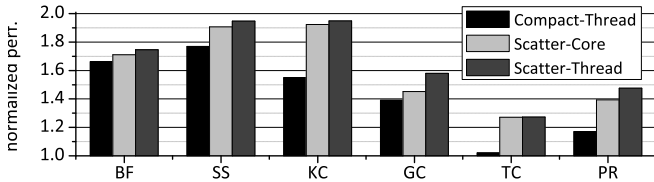


Figure 9. Performance comparison between various configurations of thread placement and affinity (normalized to Compact-Core).

2) *The Optimal Thread Number*: We define the optimal thread number as the thread number achieving the best performance for each benchmark. Figure 8 describes the optimal thread number for all applications with all datasets. There is no universal optimal thread number, e.g., the physical core number or the logic core number, that can always have the best performance for all applications with all datasets. Different applications require distinctive optimal thread numbers for their own best performance. Moreover, even for the same application, the optimal thread numbers for various graph datasets are different. As Figure 7 shows, naively setting the thread number to 256 decelerates KC with twitter7 by $4.1\times$.

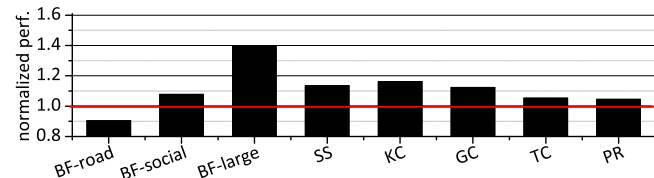


Figure 10. Performance comparison between DDR4 and MCDRAM (normalized to DDR4 DRAM main memory without MCDRAM).

3) *Thread Placement and Affinity*: Because graphBIG depends on the OpenMP library, we can configure the thread placement and affinity by $KMP_AFFINITY=X$, $granularity=Y$. Here, X indicates thread placement and has two options: assigning thread $n+1$ to an available thread context as close as possible to that of thread n (Compact) or distributing threads as evenly as possible across the entire system (Scatter). And Y denotes granularity and includes two choices: allowing all threads bound to a physical core to float between different thread contexts (Core) or causing each thread to be bound to a single thread context (Thread). The performance comparison between all configurations of thread placement and affinity is shown in Figure 9, where each bar represents one X - Y combination and all bars are normalized to Compact-Core. We see that Compact configurations with Core and Thread have worse performance, since Scatter configurations better utilize all physical cores and distribute memory requests evenly among all memory controllers. In two Scatter configurations, granularity Thread wins slightly better performance, because it scatters consecu-

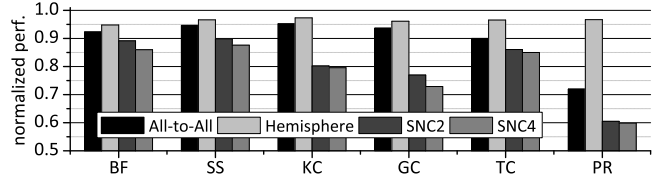


Figure 11. Performance comparison between different cache clustering modes (normalized to *Quadrant*).

tive threads sharing similar application behaviors to different physical cores.

D. MCDRAM

We configured MCDRAMs as a hardware-managed L3 cache for KNL as default. However, we can also disable MCDRAMs to use only DDR4 DRAM as main memory. The performance improvement of MCDRAM is exhibited in Figure 10, where all results are normalized to the performance of KNL with only DDR4 DRAM main memory. Compared to DDR4, MCDRAM can supply $4\times$ bandwidth. The MCDRAM-based cache boosts the performance of most benchmarks by its 16GB capacity and larger bandwidth. Particularly, large graph datasets benefit more from the MCDRAM-based cache, since they enlarge the working set size for most applications. On the contrary, the access latency of MCDRAM is longer than that of DDR4 DRAM [18]. Some benchmarks processing graphs with large diameter and small average degree, e.g. BF with road datasets, are more sensitive to the prolonged memory access latency and actually decelerated by the MCDRAM-based cache, because they have a limited number of concurrently pending memory accesses and small memory footprints.

E. Vectorization

We compiled graph programs with `icc -O3` with various vectorization choices: no vectorization, *AVX2* and *AVX512*. *AVX2* improves the graph processing performance by $47\%\sim 324\%$ over *NOVEC*. However, compared to *AVX2*, *AVX512* does not significantly further boost the performance of graph applications. This is because the implementation of System G does not explicitly represent vertices or edges by floating point or integer arrays. Instead, vertices and edges are encapsulated into lists or maps in graph frameworks, and thus it is difficult to vectorize these data structure on KNL. Moreover, most graph datasets we used are sparse, so they can barely be improved by wide *AVX512* SIMD instructions.

F. Cache Clustering Mode

The performance comparison between various cache clustering modes is shown in Figure 11, where all results are normalized to *Quadrant*. As explained in Section II-D3, among all hardware-managed modes including *All-to-All*, *Hemisphere* and *Quadrant*, *Quadrant* achieves the best performance, since it can keep both L2 accesses and memory accesses served within the local quadrant on KNL. Two

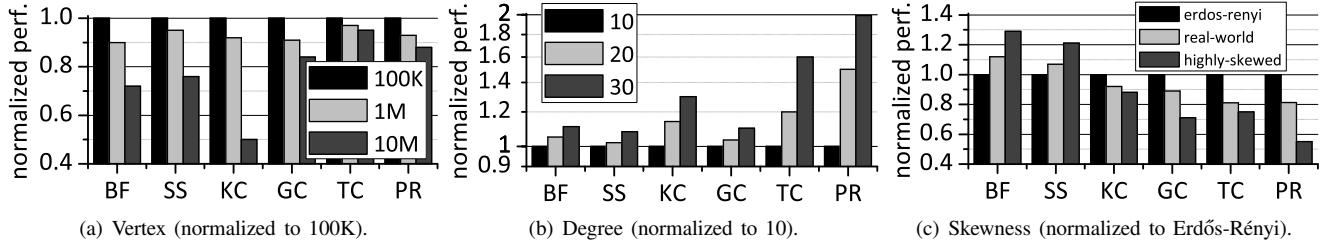


Figure 12. The KNL performance with various R-MAT datasets.

software-managed modes (*SNC-2* and *SNC-4*) have even worse graph processing performance, since benchmarks in graphBIG do not have NUMA-awareness and have to pay huge penalty for frequent communications between different NUMA regions. *SNC-4* offers four NUMA regions, therefore, compared to *SNC-2*, it degrades the graph application performance more by expensive communications between NUMA regions.

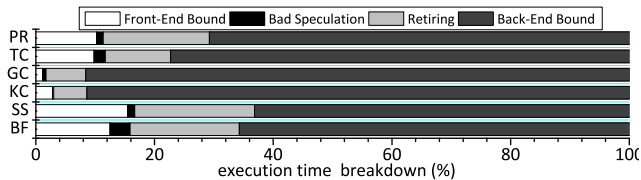


Figure 13. KNL execution time breakdown (with twitter7).

G. Execution Time Breakdown

To understand bottlenecks of applications, we show KNL execution time breakdown of all benchmarks with twitter7 graph input in Figure 13. *Bad Speculation* is the time stall due to branch mis-prediction. *Retiring* denotes the time occupied by the execution of useful instructions. *Front-End Bound* indicates the time spent by fetching and decoding instructions, while *Back-End Bound* means the waiting time due to a lack of required resources for accepting more instructions in the back-end of the pipeline, e.g., data cache misses and main memory accesses. It is well known that graph applications are extremely memory intensive and have irregular data accesses. The breakdown of execution time on KNL also supports this observation. For *KC* and *GC*, $> 90\%$ execution time is used to wait for back-end stalls. Although in Figure 10, the 16GB MCDRAM-based cache improves the performance of *KC* and *GC* by $> 10\%$ on average, we anticipate a larger capacity MCDRAM-based cache can further boost their performance, especially when processing large graph datasets.

H. Dataset Sensitivity Analysis

We generated scale-free graphs (R-MAT [5]) with varying numbers of vertices, average vertex degrees and vertex distributions by PaRMAT generator [10]. The dataset sensitivity studies on KNL are shown in Figure 12. Among four R-MAT parameters (a , b , c and d) [5], we always enforced $d = 0.11$ and $b = c$ for symmetry. To produce different skewnesses, we set the ratio (*skewness*) between a and b to 1

(Erdős-Rényi model), 3 (real-world) and 8 (highly-skewed). In Figure 12(a), we fixed the average vertex degree to 20 and the *skewness* to 3. With an increasing number of vertices, the performance of all benchmarks degrades more or less, mainly because D-TLB and L2 miss rates are increased by larger graph sizes. We fixed the vertex number to 1M and the *skewness* to 3 in Figure 12(b). All applications demonstrate better performance with an enlarging average vertex degree, since more edges per vertex lead to increased L2 hit rate. Particularly, the performance of *TC* and *PR* increases more obviously, since they operate on neighbor sets in vertex properties and higher vertex degree brings more accesses within vertices. In Figure 12(c), we set the number of vertices and the average vertex degree to 1K and 20 respectively. And we explored the *skewness* among 1, 3 and 8. For traversal-based applications, *BF* and *SS*, as graph skewness increases and graph eccentricity decreases, the application performance increases. Higher-skewed graphs have smaller diameter resulting in faster traversals. On the contrary, the other graph benchmarks suffer from severer load imbalance, when their datasets are more skewed.

VI. CONCLUSION AND FUTURE WORK

In this paper, we present a performance and power characterization study to show the potential of KNL MIC processors on parallel graph processing. To fully utilize KNLs, in future, we need to overcome challenges from both hardware angle and software perspective.

First, from hardware angle, KNL supplies many architectural features that can be configured by knobs. Different graph applications may favor different configurations. For instance, different graph benchmarks require distinctive numbers of threads to achieve the best performance. Furthermore, some applications benefit from high bandwidth MCDRAM, others may be improved by low latency DDR4 DRAM. Therefore, it is vital to have auto-tuning tools to search the optimal configuration of these knobs to achieve the best performance on KNLs. Previous works propose exhaustive iteration-based optimizations [11] and machine-learning-based tuning techniques [25]. For KNLs, we believe that future auto-tuning schemes have to consider both the MIC architecture and the heterogeneous main memories.

Second, from software perspective, though a state-of-the-art multi-threaded graph framework fully optimized for traditional multi-core CPUs can run on KNLs, we observe

that hardware resources such as VPU are underutilized and advanced software-managed architectural features, e.g., the SNC-4 cache clustering mode, may even hurt the performance of graph applications. In future, instead of optimizing a single benchmark, we need to create a fully vectorized graph framework offering AVX512 friendly primitives to support a wide variety of graph applications on KNLs. Moreover, we should incorporate a OS-based [15] or library-based [14] NUMA-aware memory management technique into future graph frameworks, so that the graph applications can benefit from the lowest local memory access latency provided by SNC-4 without incurring large communication overhead between NUMA regions. The future graph frameworks on KNLs also need to be rewritten with heterogeneous memory supporting libraries such as MEMKIND to allocate latency sensitive pages to DDR4 DRAMs and bandwidth sensitive pages to MCDRAMs.

ACKNOWLEDGMENT

We gratefully acknowledge the support from Intel Parallel Computing Center (IPCC) grant, NSF OCI-114932 and CIFDIBBS-143054. We appreciate the support from IU PHI, FutureSystems team, and ISE Modelling & Simulation Lab.

REFERENCES

- [1] M. Ahmad, *et al.*, “CRONO: A Benchmark Suite for Multithreaded Graph Algorithms Executing on Futuristic Multi-cores,” in *IEEE International Symposium on Workload Characterization*, pages 44–55, 2015.
- [2] T. Barnes, *et al.*, “Evaluating and Optimizing the NERSC Workload on Knights Landing,” in *International Workshop on Performance Modeling, Benchmarking and Simulation of High Performance Computer Systems*, pages 43–53, 2016.
- [3] S. Beamer, *et al.*, “Locality Exists in Graph Processing: Workload Characterization on an Ivy Bridge Server,” in *IEEE International Symposium on Workload Characterization*, pages 56–65, 2015.
- [4] M. Burtscher, *et al.*, “A quantitative study of irregular programs on GPUs,” in *IEEE International Symposium on Workload Characterization*, pages 141–151, 2012.
- [5] D. Chakrabarti, *et al.*, “R-MAT: A Recursive Model for Graph Mining,” in *SIAM International Conference on Data Mining*, 2004.
- [6] L. Chen, *et al.*, “Efficient and Simplified Parallel Graph Processing over CPU and MIC,” in *IEEE International Parallel and Distributed Processing Symposium*, pages 819–828, 2015.
- [7] L. Chen, *et al.*, “Exploiting Recent SIMD Architectural Advances for Irregular Applications,” in *International Symposium on Code Generation and Optimization*, pages 47–58, 2016.
- [8] P. Jiang, *et al.*, “Reusing Data Reorganization for Efficient SIMD Parallelization of Adaptive Irregular Applications,” in *International Conference on Supercomputing*, pages 16:1–16:10, 2016.
- [9] F. Khorasani, *et al.*, “CuSha: Vertex-centric Graph Processing on GPUs,” in *International Symposium on High-performance Parallel and Distributed Computing*, pages 239–252, 2014.
- [10] F. Khorasani, *et al.*, “Scalable SIMD-Efficient Graph Processing on GPUs,” in *International Conference on Parallel Architectures and Compilation Techniques*, pages 39–50, 2015.
- [11] P. A. Kulkarni, *et al.*, “Practical Exhaustive Optimization Phase Order Exploration and Evaluation,” *ACM Transactions on Architecture and Code Optimization*, 6(1):1:1–1:36, April 2009.
- [12] A. Kyröla, *et al.*, “GraphChi: Large-scale Graph Computation on Just a PC,” in *USENIX Conference on Operating Systems Design and Implementation*, pages 31–46, 2012.
- [13] J. Leskovec and A. Krevl, “SNAP Datasets: Stanford Large Network Dataset Collection,” <http://snap.stanford.edu/data>, June 2014.
- [14] S. Li, *et al.*, “NUMA-aware Shared-memory Collective Communication for MPI,” in *International Symposium on High-performance Parallel and Distributed Computing*, pages 85–96, 2013.
- [15] T. Li, *et al.*, “Efficient Operating System Scheduling for Performance-asymmetric Multi-core Architectures,” in *ACM/IEEE Conference on Supercomputing*, pages 53:1–53:11, 2007.
- [16] Y. Low, *et al.*, “Distributed GraphLab: A Framework for Machine Learning and Data Mining in the Cloud,” *Proceedings of the VLDB Endowment*, 5(8):716–727, APR 2012.
- [17] G. Malewicz, *et al.*, “Pregel: A System for Large-scale Graph Processing,” in *ACM SIGMOD International Conference on Management of Data*, pages 135–146, 2010.
- [18] J. D. McCalpin, “Memory Latency on the Intel Xeon Phi x200 KNL processor,” <http://sites.utexas.edu/jdm4372/tag/memory-latency/>, December 2016.
- [19] R. C. Murphy, *et al.*, “Introducing the Graph 500,” in *Cray User’s Group (CUG)*, 2010.
- [20] L. Nai, *et al.*, “GraphBIG: Understanding Graph Computing in the Context of Industrial Solutions,” in *International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 69:1–69:12, 2015.
- [21] A. Sodani, “Knights landing (KNL): 2nd Generation Intel® Xeon Phi processor,” in *IEEE Hot Chips Symposium*, pages 1–24, 2015.
- [22] T. Suzumura, *et al.*, “Performance characteristics of Graph500 on large-scale distributed environment,” in *IEEE International Symposium on Workload Characterization*, pages 149–158, 2011.
- [23] I. Tanase, *et al.*, “A Highly Efficient Runtime and Graph Library for Large Scale Graph Analytics,” in *Workshop on GRaph Data Management Experiences and Systems*, 2014.
- [24] Y. Tian, *et al.*, “From “Think Like a Vertex” to “Think Like a Graph,”” *Proceedings of the VLDB Endowment*, 7(3):193–204, November 2013.
- [25] A. Tiwari, *et al.*, “A Scalable Auto-tuning Framework for Compiler Optimization,” in *IEEE International Symposium on Parallel&Distributed Processing*, pages 1–12, 2009.
- [26] J. Treibig, *et al.*, “LIKWID: A lightweight performance-oriented tool suite for x86 multicore environments,” in *International Workshop on Parallel Software Tools and Tool Infrastructures*, 2010.
- [27] Y. Wang, *et al.*, “Gunrock: A High-performance Graph Processing Library on the GPU,” in *ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, pages 265–266, 2015.
- [28] J. Zhong and B. He, “Medusa: A Parallel Graph Processing System on Graphics Processors,” *ACM SIGMOD Record*, 43(2):35–40, December 2014.