

A Survey on Near Data Processing Accelerators for Graph Analytics

Makesh Tarun Chandran

*Department of Electrical Engineering and Computer Science
Pennsylvania State University
Email: mzc88@psu.edu*

1. Introduction

Graph algorithms are an interesting class of application for which a single mac mini desktop equipped with SSD outperforms a medium sized cluster! [1] The primary reason for such a counter-intuitive behavior can be attributed to the irregular data-access pattern inherent in those algorithms. They also have a low computation to communication ratio which is exacerbated in the modern memory-bandwidth limited systems. However, graph algorithms are highly iterative and there is abundance of parallelism that can be exploited. These idiosyncrasies of graph algorithm render commercial HPC systems inefficient or unsuitable for accelerating such applications. Extensive research on graph processing in the past decade have revealed new insights on the algorithmic, programming model and hardware aspects of such systems. I have tried to comprehend the relevant and significant contributions of previous work in the following section.

2. Graph Processing

Current graph analytics frameworks are based on two fundamentally different computation models - Pregel and Gather-Apply-Scatter (GAS) abstractions.

2.1. Vertex-Centrix Processing

Discuss the Vertex-based [2] and Gather-Apply-Scatter [3] based programming models here. Highlight the differences and bring out the suitability of each of them for PIM operations. Also show the pseudo code for SSSP application.

2.1.1. Pregel. Pregel is a vertex-centric programming model following the Bulk Synchronous Parallel style of execution. The computation on vertices are represented as a sequence of supersteps with a global synchronisation among the vertices participating in a superstep. A superstep only contains the set of active nodes which was enabled by the vertices of previous superstep. The vertices in a superstep can be operated on in parallel and they generate messages to vertices of the next superstep. The program halts when there are no active vertices in the next superstep. Algorithm 1 is a pseudo code of a typical Pregel program. This programming model provides an intuitive framework to construct a large

scale graph program based on the computation on a single vertex. The Synchronicity of the model makes it easy to reason about and avoids any possibility deadlock or data races. One limitation could be the presence of the global barrier synchronisation, whose performance effects can be mitigated by having adequate parallelism slack. Some of the frameworks that are based on the Pregel programming are Giraph, Giraph++, Mizan, GPS, Pregelix and Pregel+.

Algorithm 1 Pregel Abstraction

```
1: super_step  $\leftarrow$  0
2: while active_vertices  $\neq \emptyset$  do
3:   for  $v \in$  active_vertices do
4:     v.value  $\leftarrow$  v.compute(msgs)
5:     v.send_message(neighbourhood(v))
6:     v.halt()
7: super_step  $\leftarrow$  super_step + 1
```

2.1.2. Gather-Apply-Scatter. TODO: Rewrite GAS programming model splits graph processing into three steps - Gather, Apply and Scatter. During the Gather phase a vertex accumulates information about its adjacent edges and processes it. The processed value is used to update itself and its neighbors in Apply and Scatter phase respectively. The updated neighbors are activated and will go through the GAS steps. The program halts when there are no vertices to be processed. The GAS model can be realised as both synchronous and asynchronous execution depending on the application. One fundamental difference from pregel is that vertices can directly access the neighborhood without having to be explicitly communicated to it. But this requires shared memory abstraction for the GAS model to be realised, which increases the latency in a distributed systems. Techniques such as data replication, prefetching, etc. can be employed to overcome the bottleneck. A basic sequence of steps involved in GAS model is presented in Algorithm 2. GraphLab, GraphChi and PowerGraph are some of the frameworks developed on GAS Abstraction.

The paper [21] does a performance evaluation of the various GAS and Pregel based large scale graph processing framework. GAS models enables data partitioning based on edges whereas pregel model is limited to vertex based partitioning. As many natural graphs follow power law

like vertex-degree distribution, edge-based graph partitioning leads to better load balancing for such graphs.

Algorithm 2 GAS Abstraction

```

1: super_step  $\leftarrow$  0
2: while active_vertices  $\neq \emptyset$  do
3:   for  $v \in$  active_vertices do
4:     /* Pull data from replicas */
5:     accumulator  $\leftarrow$  sum( $v$ .gather(neighbourhood( $v$ ))
6:     /* Push updates to replicas */
7:      $v$ .value  $\leftarrow$   $v$ .apply(accumulator, val( $v$ ))
8:      $v$ .scatter(neighbourhood( $v$ ))
9: super_step  $\leftarrow$  super_step + 1

```

2.2. Matrix Algebra-Based Systems

Before the advent of graph style processing frameworks such as Pregel, graph computations were modelled as matrix-vector multiplication. Every iteration executes the user defined function on the adjacency list(vector) of a vertex and update the adjacency matrix. Matrix based computation enables the use of optimisations such as efficient data layouting, prefetching and use of standard linear algebra sub-routines. PEGASUS, GBASE and SystemML are three frameworks which incorporate matrix based graph computation in their framework. One drawback of matrix based graph computation is the amount of redundant work that is performed every step. Each iteration triggers a full matrix computation even though only a few of the node update their value. In contrast, in pregel like systems, computation is performed only on the active nodes which are bound to update their values. To overcome the limitation, some frameworks present a vertex-centric programming framework and translate it to matrix computation underneath to utilise the best of both the worlds.

3. In/Near-Memory Processing

- Reduce the cost of data movement - PIM Analog characteristics

3.1. Near-Data Processing

Discuss [16] like architectures where the computations are moved closer to the memory elements but not inside them. End the section by indicating that more bandwidth can be exploited when the processing is integrated with memory.

3.2. In-memory Processing

3.2.1. DRAM-PIM. Present a basic idea of the [11] paper which proposes a architecture for processing-in-memory architecture built on commercially demonstrated Hybrid Memory Cube. Also highlight the loopholes/assumptions with such design. Also list other architectures, like GraphH [17], which builds on the technology.

3.2.2. NVM-PIM. Briefly list the problems with DRAM in terms of energy and performance, and bring out NVMs as an alternative to DRAM. [12] introduces a efficient NOR operation which can be performed on the bitlines of ReRAM based memory. Show example architecture of NVM PIN and explain it working with relevant figures.

4. Accelerators

An Ideal graph processing system is where the performance increases proportional to the size of the graphs that can be stored in the system. Unfortunately, in conventional systems, memory bandwidth remains almost constant irrespective of the memory capacity due to pin count limitation per chip. Traditional caching mechanisms fail to sustain the memory throughput requirement due to the application's poor locality.

Following is the summary of various PIM-based techniques developed to mitigate the performance bottleneck in graph processing systems.

4.1. Graphicionado

[5] replaces the conventional on-chip memory hierarchy (caches) with explicitly managed ScratchPad memory. The accelerator achieves upto 6x speedup while consuming less than 2% of energy of conventional systems. However, larger graph applications that doesn't fit in the on-chip memory will incur off-chip communication which is detrimental to the performance.

4.2. Tesseract

[6] exploits the internal bandwidth of the HMC-RAM, which is an order of magnitude higher than the off-chip bandwidth, by integrating a logic layer within the memory die. The accelerator out performs conventional systems by a factor of 10 and achieves memory-capacity-proportional bandwidth to provide scalable performance improvements.

4.3. GraphR

[7] is a RRAM based in-memory graph processing accelerator. They perform analog computations using RRAM crossbar structure and achieve a speedup of 4x and energy efficiency of 11x compared to Tesseract. They exploit the fact that the graph algorithms can inherently tolerate imprecision and are resilient to errors.

4.4. GraphP

[8] takes a data-organisation centric hardware-software design approach to minimise communication in a PIM based accelerator, to achieve a speedup of up to 1.7x and energy efficiency on of 89% compared to Tesseract.

	Graphicianado	GraphH	Tesseract	GraphR	GraphP
Programing Model	Vertex-centric	Gather-Apply-Broadcast	Vertex-centric	Vertex-centric	Two phase vertex
Execution Model	Matrix computation	Message passing and shared memory	Message passing and iterative	Sparse matrix computation	Message passing and iterative
Computing Mode	Sync		Sync		
Generality	Vertex program		Vertex program	Vertex program in SpMV	
Data access	Explicitly managed in Scratchpad	In-memory processing	DRAM-PIM	ReRAM-PIM	DRAM-PIM

4.5. GraphH

[17] GraphH takes a fresh look at the graph analytics systems and proposes a DRAM-PIM based architecture to tackle challenges in graph processing. The proposed architecture outperforms Graphicianado by a factor of 5.12x.

All these systems have customised their programming and execution model to derive maximum performance from their hardware.

5. Completed and Future Work

The architectures listed above were proposed in the last three years, with the advent of new memory technologies which made in-memory processing feasible. An in-depth study of the hardware and software stack of these PIM-based accelerators and their design choices, may expose untapped design points to improve performance. The plan was to begin by understanding various design points listed in Section ??, evaluate the accelerators for strengths and limitations, and identify potential improvements in their design, leading to a survey paper or an original research paper.

I have read papers GraphR [7], GraphP [8] and other papers based on NVM-PIM [18]. Apart from that I also ramped up with Pregel and ZSim Simulation framework. The next step would be to prototype one of the architectures (possibly GraphR or GraphH) and identify bottlenecks to improve the performance. A survey paper on this topic makes little sense as the in-memory processing of graphs is a relatively narrow and emerging field. Moreover, there are recent resources [19], [20] that contain an exhaustive analysis of the topic.

References

- [1] Aapo Kyrola, Guy Blelloch, and Carlos Guestrin. 2012. GraphChi: large-scale graph computation on just a PC. In Proceedings of the 10th USENIX conference on Operating Systems Design and Implementation (OSDI'12). USENIX Association, Berkeley, CA, USA, 31-46.
- [2] Malewicz, Matthew H. Austern, Aart J.C Bik, James C. Dehnert, Ilan Horn, Naty Leiser, and Grzegorz Czajkowski. 2010. Pregel: a system for large-scale graph processing. In Proceedings of the 2010 ACM SIGMOD International Conference on Management of data (SIGMOD '10). ACM, New York, NY, USA, 135-146. DOI: <https://doi.org/10.1145/1807167.1807184>
- [3] Joseph E. Gonzalez, Yucheng Low, Haijie Gu, Danny Bickson, and Carlos Guestrin. 2012. PowerGraph: distributed graph-parallel computation on natural graphs. In Proceedings of the 10th USENIX conference on Operating Systems Design and Implementation (OSDI'12). USENIX Association, Berkeley, CA, USA, 17-30.
- [4] Amber Hassaan, Martin Burtscher, and Keshav Pingali. 2011. Ordered vs. unordered: a comparison of parallelism and work-efficiency in irregular algorithms. In Proceedings of the 16th ACM symposium on Principles and practice of parallel programming (PPoPP '11). ACM, New York, NY, USA, 3-12. DOI: <http://dx.doi.org/10.1145/1941553.1941557>
- [5] T. J. Ham, L. Wu, N. Sundaram, N. Satish and M. Martonosi. "Graphicianado: A high-performance and energy-efficient accelerator for graph analytics," 2016 49th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO), Taipei, 2016, pp. 1-13.
- [6] J. Ahn, S. Hong, S. Yoo, O. Mutlu and K. Choi, "A scalable processing-in-memory accelerator for parallel graph processing," 2015 ACM/IEEE 42nd Annual International Symposium on Computer Architecture (ISCA), Portland, OR, 2015, pp. 105-117.
- [7] L. Song, Y. Zhuo, X. Qian, H. Li and Y. Chen, "GraphR: Accelerating Graph Processing Using ReRAM," 2018 IEEE International Symposium on High Performance Computer Architecture (HPCA), Vienna, 2018, pp. 531-543.
- [8] M. Zhang et al., "GraphP: Reducing Communication for PIM-Based Graph Processing with Efficient Data Partition," 2018 IEEE International Symposium on High Performance Computer Architecture (HPCA), Vienna, 2018, pp. 544-557.
- [9] Muhammet Mustafa Ozdal, Serif Yesil, Taemin Kim, Andrey Ayupov, John Greth, Steven Burns, and Ozcan Ozturk. 2016. Energy efficient architecture for graph analytics accelerators. In Proceedings of the 43rd International Symposium on Computer Architecture (ISCA '16). IEEE Press, Piscataway, NJ, USA, 166-177. DOI: <https://doi.org/10.1109/ISCA.2016.24>
- [10] C. Zhang, T. Meng and G. Sun, "PM3: Power Modeling and Power Management for Processing-in-Memory," 2018 IEEE International Symposium on High Performance Computer Architecture (HPCA), Vienna, 2018, pp. 558-570.
- [11] R. Nair et al., "Active Memory Cube: A processing-in-memory architecture for exascale systems," in IBM Journal of Research and Development, vol. 59, no. 2/3, pp. 17:1-17:14, March-May 2015. doi: 10.1147/JRD.2015.2409732
- [12] S. Kvatinisky et al., "MAGICMemristor-Aided Logic," in IEEE Transactions on Circuits and Systems II: Express Briefs, vol. 61, no. 11, pp. 895-899, Nov. 2014. doi: 10.1109/TCSII.2014.2357292
- [13] Muhammad Imran, Carlos Castillo, Fernando Diaz, and Sarah Vieweg. 2015. Processing Social Media Messages in Mass Emergency: A Survey. ACM Comput. Surv. 47, 4, Article 67 (June 2015), 38 pages. DOI: <https://doi.org/10.1145/2771588>
- [14] C. Unsalan and B. Sirmacek, "Road Network Detection Using Probabilistic and Graph Theoretical Methods," in IEEE Transactions on Geoscience and Remote Sensing, vol. 50, no. 11, pp. 4441-4453, Nov. 2012.
- [15] Bullmore, Ed Sporns, Olaf Complex brain networks: graph theoretical analysis of structural and functional systems DOI: <http://dx.doi.org/10.1038/nrn2575>
- [16] Boncheol Gu, Andre S. Yoon, Duck-Ho Bae, Insoon Jo, Jinyoung Lee, Jonghyun Yoon, Jeong-Uk Kang, Moonsang Kwon, Chanhoo Yoon, Sangyeun Cho, Jaehoon Jeong, and Duckhyun Chang. 2016. Biscuit: a framework for near-data processing of big data workloads. SIGARCH Comput. Archit. News 44, 3 (June 2016), 153-165. DOI: <https://doi.org/10.1145/3007787.3001154>

- [17] G. Dai et al., "GraphH: A Processing-in-Memory Architecture for Large-scale Graph Processing," in IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems. doi: 10.1109/TCAD.2018.2821565
- [18] Y. K. Rupesh, P. Behnam, G. R. Pandla, M. Miryala and M. Nazm Bojnordi, "Accelerating k-Medians Clustering Using a Novel 4T-4R RRAM Cell," in IEEE Transactions on Very Large Scale Integration (VLSI) Systems. doi: 10.1109/TVLSI.2018.2808468
- [19] Systems for Big Graph Analytics by Yan, Da, Tian, Yuanyuan, Cheng, James
- [20] Gupta, S. (2018). Processing in Memory using Emerging Memory Technologies. UC San Diego. ProQuest ID: Merritt ID: ark:/13030/m5bg7kwn. Retrieved from <https://escholarship.org/uc/item/95z3z84c>
- [21] Omar Batarfi, Radwa El Shawi, Ayman G. Fayoumi, Reza Nouri, Seyed-Mehdi-Reza Beheshti, Ahmed Barnawi, and Sherif Sakr. 2015. Large scale graph processing systems: survey and an experimental evaluation. Cluster Computing 18, 3