

# Big data for all!

Clement Fung

cfung1@cs.ubc.ca

Stewart Grant

sgrant09@cs.ubc.ca

## Abstract

The modern internet generates petabytes of data per day. Processing vast amounts of data is an increasingly common task for both scientists and modestly experienced programmers. Often this data is naturally represented as a graph, such as social media friends, and requires clusters of machines to process. Concurrent trends in data centre architecture suggest that the rack is the new server, and shared memory is now a feasible interface between rack collocated servers. These trends made us wonder *How simple can fast graph processing be on a rack of servers?*. We investigated the tradeoffs of the conventional pregal style *"think like a vertex"* programming model, and found its performance unacceptable. In contrast we find that a *"Think like a sub-graph"* model respects locality in common graphs, provides a more holistic programming interface, and runs fast!

## 1 Introduction

Big data processing is complicated. Typically scientists, and experienced programmers alike struggle with managing and configuring clusters of machines to processes large amounts of data. Frameworks like Hadoop and Pregal have significantly eased the difficulty of big data processing, but they remain intimidating for the lay man. We noticed a trend in non-systems scientists, and researchers, that they wanted to *"Just write python code that ran on a bunch of computers"*. For the benefit of science we investigated how to make this dream a reality.

Making all Python code run on clusters of machines is impractical, and due to network overhead would lead to sluggish un-optimal code. Instead we concentrated our efforts on a common, but difficult big data processing task, graph processing. Many frameworks exist for distributed graph processing [14, 5, 12, 13, 19, 6], and many general frameworks exist which are used for graph processing [18, 20, 10, 16]. These frameworks vary in their complexity, but none are *"accessible"* for programmers with no systems experience.

With the exception of [12] the aforementioned systems suffer a common pitfall to accessibility; they expose the complexity of a distributed message passing system to the user. Extensive work has been done to hide this complexity in the abstraction of distributed shared memory (DSM) [11, 17, 15, 7, 9]. The benefits of DSM have been ignored in recent years due its flaws, mainly fate sharing and sub optimal performance. Dismissing DSM may well have been a shortsighted mistake. Ultra dense memory, and the approach of terabit bandwidth within a rack give modern racks the appearance of single machines, and has lead towards disaggregated architectures [1, 2, 3, 4, 8]. Such futuristic systems lend themselves naturally to DSM which motivates our proposal for a corresponding computation framework.

The largest disadvantage of DSM is performance. Programmers can write terribly performant programs by failing to reason about the location of memory, leading to memory thrashing. In computation where a high degree of consistency between shared resources DSM is the wrong tool for the job. In contrast when large amounts of computation can be performed between memory synchronizations DSM provides a simple and efficient programming model. Graph processing suffers from a lack of locality. In computations such as PageRank a single iteration may require edge updates which require the synchronization of every machine in a cluster. This problem can be largely avoided in practice by carefully pre-processing graphs into partitions where the minimum number of graph edges cross machines. The cost of pre-processing a graph can be large, in some cases the complexity of finding a good partition is greater than solving the initial problem! Here we demonstrate that the cost of graph partitioning is worth it for the benefits that DSM provides.

In this paper we attack the problem of developing a simple and efficient graph processing interface for DSM. Specifically we make the following contributions.

- A simple graph processing api
- A graph partitioning scheme optimal for DSM
- An evaluation of processing performance between

partitioned DSM processing, and pregal style graph processing

The remainder of this paper is organized as follows. In Section 2 we over view related work. In Section 3 we describe our graph processing API. Section 4 we describe our approach to graph partitioning. In Section 5 we evaluate our framework against a pregal style *think like a vertex model*. Section 6 describes our experiences with our system, and Section 7 concludes the paper.

## 2 Related work

## 3 API

## 4 Partitioning

## 5 Evaluation

## 6 Discussion

## 7 Conclusion

## References

- [1] Facebook disaggregated rack, <http://goo.gl/6h2ut>.
- [2] Hp the machine, <http://www.hpl.hp.com/research/systems-research/themachine/>.
- [3] Intel rsa. <https://software.intel.com/en-us/articles/intel-performance-counter-monitoring>.
- [4] Seamicro technology overview, <http://seamicro.com/>.
- [5] A. Ching, S. Edunov, M. Kabiljo, D. Logothetis, and S. Muthukrishnan. One trillion edges: Graph processing at facebook-scale. *Proc. VLDB Endow.*, 8(12):1804–1815, Aug. 2015.
- [6] J. E. Gonzalez, Y. Low, H. Gu, D. Bickson, and C. Guestrin. Powergraph: Distributed graph-parallel computation on natural graphs. In *Proceedings of the 10th USENIX Conference on Operating Systems Design and Implementation, OSDI’12*, pages 17–30, Berkeley, CA, USA, 2012. USENIX Association.
- [7] I. F. Haddad and E. Paquin. Mosix: A cluster load-balancing solution for linux. *Linux J.*, 2001(85es), May 2001.
- [8] S. Han, N. Egi, A. Panda, S. Ratnasamy, G. Shi, and S. Shenker. Network support for resource disaggregation in next-generation datacenters. In *Proceedings of the Twelfth ACM Workshop on Hot Topics in Networks, HotNets-XII*, pages 10:1–10:7, New York, NY, USA, 2013. ACM.
- [9] Z. Huang, W. Chen, and et al. Vodca: View-oriented, distributed, cluster-based approach to parallel computing. In *DSM WORKSHOP 2006, IN: PROC. OF THE IEEE/ACM SYMPOSIUM ON CLUSTER COMPUTING AND GRID 2006 (CC-GRID06)*, IEEE COMPUTER SOCIETY, 2006.
- [10] M. Isard, M. Budiu, Y. Yu, A. Birrell, and D. Fetterly. Dryad: Distributed data-parallel programs from sequential building blocks. In *Proceedings of the 2Nd ACM SIGOPS/EuroSys European Conference on Computer Systems 2007, EuroSys ’07*, pages 59–72, New York, NY, USA, 2007. ACM.
- [11] P. Keleher, A. L. Cox, S. Dwarkadas, and W. Zwaenepoel. Treadmarks: Distributed shared memory on standard workstations and operating systems. In *Proceedings of the USENIX Winter 1994 Technical Conference on USENIX Winter*

- 1994 Technical Conference, WTEC'94, pages 10–10, Berkeley, CA, USA, 1994. USENIX Association.
- [12] A. Kyrola, G. Blelloch, and C. Guestrin. Graphchi: Large-scale graph computation on just a pc. In *Proceedings of the 10th USENIX Conference on Operating Systems Design and Implementation*, OSDI'12, pages 31–46, Berkeley, CA, USA, 2012. USENIX Association.
  - [13] Y. Low, D. Bickson, J. Gonzalez, C. Guestrin, A. Kyrola, and J. M. Hellerstein. Distributed graphlab: A framework for machine learning and data mining in the cloud. *Proc. VLDB Endow.*, 5(8):716–727, Apr. 2012.
  - [14] G. Malewicz, M. H. Austern, A. J. Bik, J. C. Dehnert, I. Horn, N. Leiser, and G. Czajkowski. Pregel: A system for large-scale graph processing. In *Proceedings of the 2010 ACM SIGMOD International Conference on Management of Data*, SIGMOD '10, pages 135–146, New York, NY, USA, 2010. ACM.
  - [15] C. Morin, R. Lottiaux, G. Vallee, P. Gallard, D. Margery, J.-Y. Berthou, and I. D. Scherson. Kerighed and data parallelism: Cluster computing on single system image operating systems. In *Proceedings of the 2004 IEEE International Conference on Cluster Computing*, CLUSTER '04, pages 277–286, Washington, DC, USA, 2004. IEEE Computer Society.
  - [16] D. G. Murray, F. McSherry, R. Isaacs, M. Isard, P. Barham, and M. Abadi. Naiad: A timely dataflow system. In *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles*, SOSP '13, pages 439–455, New York, NY, USA, 2013. ACM.
  - [17] R. Power and J. Li. Piccolo: Building fast, distributed programs with partitioned tables. In *Proceedings of the 9th USENIX Conference on Operating Systems Design and Implementation*, OSDI'10, pages 293–306, Berkeley, CA, USA, 2010. USENIX Association.
  - [18] V. K. Vavilapalli, A. C. Murthy, C. Douglas, S. Agarwal, M. Konar, R. Evans, T. Graves, J. Lowe, H. Shah, S. Seth, B. Saha, C. Curino, O. O'Malley, S. Radia, B. Reed, and E. Baldeschwieler. Apache hadoop yarn: Yet another resource negotiator. In *Proceedings of the 4th Annual Symposium on Cloud Computing*, SOCC '13, pages 5:1–5:16, New York, NY, USA, 2013. ACM.
  - [19] R. S. Xin, J. E. Gonzalez, M. J. Franklin, and I. Stoica. Graphx: A resilient distributed graph system on spark. In *First International Workshop on Graph Data Management Experiences and Systems*, GRADES '13, pages 2:1–2:6, New York, NY, USA, 2013. ACM.
  - [20] M. Zaharia, M. Chowdhury, T. Das, A. Dave, J. Ma, M. McCauley, M. J. Franklin, S. Shenker, and I. Stoica. Resilient distributed datasets: A fault-tolerant abstraction for in-memory cluster computing. In *Proceedings of the 9th USENIX Conference on Networked Systems Design and Implementation*.