



## Scalable Graph Processing Frameworks: A Taxonomy and Open Challenges

Journal:	<i>Computing Surveys</i>
Manuscript ID	CSUR-2016-0561
Paper:	Regular Paper
Date Submitted by the Author:	28-Oct-2016
Complete List of Authors:	<p>Heidari, Safiollah; The University of Melbourne, Cloud Computing and Distributed Systems (CLOUDS) Laboratory, Computing and Information Systems Department</p> <p>Simmhan, Yogesh; Indian Institute of Science, Department of Computational and Data Sciences (CDS)</p> <p>Calheiros, Rodrigo; University of Melbourne, Cloud Computing and Distributed Systems (CLOUDS) Laboratory Department of Computing and Information Systems</p> <p>Buyya, Rajkumar; University of Melbourne, Cloud Computing and Distributed Systems (CLOUDS) Laboratory, Department of Computing and Information Systems</p>
Computing Classification Systems:	Graph processing, Parallel processing, Big Data, Distributed processing and Architectures, Graph algorithms, Computing methodologies

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

Scalable Graph Processing Frameworks: A Taxonomy and Open Challenges

SAFIOLLAH HEIDARI, The University of Melbourne, Australia  
YOGESH SIMMHAN, Indian Institute of Science, India  
RODRIGO N. CALHEIROS, The University of Melbourne, Australia  
RAJKUMAR BUYYA, The University of Melbourne, Australia

The world is becoming a more conjunct place and the number of data sources such as social networks, online transactions, web search engines and mobile devices are increasing even more than had been predicted. A large percentage of this growing dataset exists in the form of linked data, more generally, graphs, and of unprecedented sizes. While today’s data from social networks contain 100’s of millions of nodes connected by billions of edges, inter-connected data from globally-distributed sensors that forms the Internet of Things (IoT) can cause this to grow exponentially larger. Although analyzing these large graphs is critical for the companies and governments that own them, big data tools designed for text and tuple analysis such as MapReduce cannot process them efficiently. So, graph distributed processing abstractions and systems are developed to design iterative graph algorithms and process large graphs with better performance and scalability. These graph frameworks propose novel methods or extend previous methods for processing graph data. In this article, we propose a taxonomy of graph processing systems and map existing systems to this classification. This captures the diversity in programming and computation models, runtime aspects of partitioning and communication, both for in-memory and distributed frameworks. Our effort helps to highlight key distinctions in architectural approaches, and identifies gaps for future research in scalable graph systems.

• Information systems→Database management system engines • Computing methodologies→Massively parallel and high-performance simulations. This is just an example, please use the correct category and subject descriptors for your submission. The ACM Computing Classification Scheme:

<http://www.acm.org/about/class/class/2012>. Please read the [HOW TO CLASSIFY WORKS USING ACM'S COMPUTING CLASSIFICATION SYSTEM](#) for instructions on how to classify your document using the 2012 ACM Computing Classification System and insert the index terms into your Microsoft Word source file.

Additional Key Words and Phrases: Big data, graph processing, large graphs, parallel processing, distributed systems

ACM Reference Format:

Safiollah Heidari, Yogesh Simmhan, Rodrigo N. Calheiros, Rajkumar Buyya, , 2016. Scalable graph processing frameworks: A taxonomy and open. *ACM Comput. Surv.*

1. INTRODUCTION

The growing popularity of technologies such as Internet of Things (IoT), mobile devices, smart phones and social networks has led towards the emergence of “Big Data”. Such applications produce not just gigabytes or terabytes of data, but soon petabytes of data that need to be actively processed. Such large volumes of data gathered from billions of connected people and devices around the world is causing unprecedented challenges in terms of how data can be stored, retrieved and managed; how data security, integrity, availability and sharing can be ensured; how massive datasets can be mined; and how they can benefit from new computing paradigms such as cloud computing for data analysis [Petey 2011, Dias de Assuncao et al 2015].

According to the National Research Council of the US National Academies [Committee on the Analysis on Massive Data 2013], graph processing is among the seven major computational methods of huge data analysis. Graph computations are used in business analytics, social network analytics, image processing, hardware design and deep learning to an increasing extent. Wide-spread techniques for processing large graphs had, till recently, been limited to shared memory [Hong et al 2012, Bader and Madduri 2008] and High Performance Computing systems [Graph500, Karypis and Kumar, 1995, Harshvardhan et al 2013]. Although distributed approaches have been proposed for processing big graphs since 2001 [Huberman, 2001], graph processing systems for commodity clusters and Clouds have become particularly popular after Google introduced its Pregel [Malewicz, et al., 2010] vertex-centric graph processing system in 2010. Since then, several distributed graph processing frameworks with diverse programming models and features have been proposed to facilitate operations on large graphs. Each of these frameworks has specific characteristics with its own strengths and weaknesses.

The aim of this paper is to provide taxonomy of scalable graph processing systems and frameworks. It identifies strengths and weaknesses in the field and proposes future directions. First, it proposes a comprehensive taxonomy of programming abstractions and runtime features offered by graph processing systems, and maps the existing systems to this taxonomy. Second, it utilizes a top-down approach for investigating graph processing frameworks and their components along with examples to support them. Third, the paper identifies gaps in existing systems which need further investigation, and discuss these open problems and future research directions in detail. In summary, this survey gives readers an overarching picture about what graph processing is, what improvements have been gained through recent frameworks, different programming and runtime techniques that have been used, and the applications that benefit from

35:2

S. Heidari et al.

them. It emphasizes on scalable graph processing platforms for shared-memory and distributed processing, that fall within the ambit of Big Data processing platforms. It also contrasts them against graph frameworks for supercomputing systems, as evidenced through the Graph500 benchmark<sup>1</sup>. On the other hand, the existing works such as [McCune et al 2015] and [Doekemeijer and Varbanescu 2014] only focus on survey and they have limited focus on key elements of graph processing.

The rest of the paper is organized as follows: Section 2 includes a definition of graphs and graph processing systems, contrasts graph processing from other big data processing methods, outlines the lifecycle of a typical graph processing system, and gives examples of real graph-based applications and algorithms. Contemporary graph processing frameworks and architectures are explained in Section 3, along with distributed coordination and computational models. Section 4 categorizes partitioning, communication models, in-memory execution, fault tolerance and scheduling. Graph databases are reviewed in section 5. A taxonomy and discussion on challenges is also presented for each section. A gap analysis and open challenges, with a perspective on future directions are discussed in Section 6 and 7 respectively. Finally, we conclude the paper in Section 8.

## 2. BACKGROUND

A Graph  $G = (V, E)$ , consists of a set of vertices,  $V = \{v_1, v_2, \dots, v_n\}$  and a set of edges,  $E = \{e_1, e_2, \dots, e_m\}$  that indicate pairwise relationships,  $E \subseteq V \times V$ . If  $(v_i, v_j) \in E$ , then  $v_i$  and  $v_j$  are neighbors [Sedgewick and Wayne 2011]. The edges may be directed or undirected. So  $V$  and  $E$  are the two defining characteristics of a graph which most of graph processing frameworks implement. Frameworks typically support a single attribute value associated with the vertex and edge (e.g., label, weight). In addition, some of the platforms also support a set of named and/or typed attributes for their vertices and edges as part of their data model.

Pairwise relationships between entities play an important role in various types of computational applications. These relationships that are implied by different connections (edges) between items (vertices) give rise to domain questions to draw value from the data, such as: Is it possible to identify transitive relationships between items by following the connections? How many items are connected to a typical item? What is the shortest distance between these items? Which groups of items are similar to each other? How important is an item relative to others? Various graph-like applications and environments are mentioned in Table 1 [Sedgewick and Wayne 2011]. As can be seen in Table 1, many applications process data that naturally fits into a graph data model. Several of these applications from social networks, eCommerce and telecom domains handle large graph datasets which need to be processed and mined to draw disparate business intelligence, ranging from the interests of people about products for targeted advertising, to tracing call logs for cyber-security. Processing large graphs poses some intrinsic challenges due to the nature of graphs themselves. These characteristics make graph processing ill-suited to existing data processing approaches, and usually inhibit efficient parallelism [Pellegrini 2011]. According to [Lumsdaine et al 2007], these properties include:

Table 1. Graph-like application and environments

Application	Item (Vertices of the Graph)	Connection (Edges of the Graph)
Social network	Members	Friendships
Computer network	Computers	Network connectivity
Web content	Web Pages	Hyperlinks
Transportation	Cities	Roads
Electrical circuit	Devices	Wires
Commerce	Customers, Goods	Purchase transactions
Factory	Machines	Production lines
Supply chain	Providers	Distances
Telecommunication	Mobile Phones	Phone calls

(1) *Data-driven computations*: Graph computations are usually entirely data-driven. Graphs are made up of sets of vertices and edges that dictate the computations performed by every graph algorithm.

(2) *Irregular problems*: Graph problems are highly irregular due to the non-uniform edge degree distribution and topological asymmetry rather than being uniformly predictable problems which can be optimally partitioning for concurrent computation.

(3) *Poor locality*: The inherent irregular characteristics of graphs leads to poor locality during computation, which is in conflict with locality-based optimizations supported of many existing processors, making it difficult to achieve high performance for graph algorithms.

(4) *High data access to computation ratio*: A large portion of graph processing is usually dedicated to data access in graph algorithms. Therefore, waiting for memory or disk fetches is the most time-consuming phase relative to the actual computation on a vertex or edge itself.

<sup>1</sup> Graph 500 benchmark, <http://www.graph500.org/>

To streamline the processing of big data, MapReduce, a distributed programming framework for processing large datasets with parallel algorithms, was introduced by Google in 2004 [Dean and Ghemawat 2004]. MapReduce has two significant advantages: 1) The programmer has a simple and familiar interface using Map and Reduce functions, inspired by functional programming concepts [Hudak 1989], and 2) the application is automatically parallelized when defined using Map and Reduce methods, without the programmer needing to know how data will be distributed, grouped and replicated, and how the tasks are scheduled.

Although MapReduce addresses many deficiencies in traditional parallel and distributed computing approaches, it has several limitations that make it less efficient for processing large graphs [Cohen, 2009] [Afrati et al 2012] [Grabowski et al 2013]: 1) MapReduce is limited to a two-phased computational model that is not naturally suited for graph algorithms that run over many iterations, 2) In common MapReduce implementations, the input graph and its state are not retained in main memory across even these two phases, let alone across iterations, and consequently requires repetitive disk I/O, 3) MapReduce’s tuple-based approach that is unaware of the linked nature of graph datasets is poorly suited to design many graph applications, and 4) Graph operations using MapReduce have poor I/O efficiency - because of frequent checkpoints on completed tasks and data replication - which is a bottleneck for many graph algorithms [Lee et al 2011].

Apache Hadoop [Apache™ Hadoop] is a popular open-source implementation of the MapReduce programming model. Besides flexible batch processing applications that can be built using Hadoop, it is also the basis for NoSQL querying platforms such as Pig [Apache Pig] and Hive [Apache Hive TM] to work with large datasets. In addition, various high level languages such as SCOPE [Zhou, et al. 2012], Sphere [Gu et al 2010] and Swazal [Pike et al 2005] are available for MapReduce-like systems. Platforms like Apache Spark [Apache Spark] have extended the programming model of MapReduce, and offer incremental batch and in-memory computation with better performance. Further, Hadoop’s distributed file system storage mechanism (HDFS) as well as a Map-only model of Hadoop is used as the storage and distributed scheduling mechanism in many graph processing frameworks we discuss.

While a number of systems such as PEGASUS [Kang et al 2009] have brought innovative approaches for processing and mining peta-scale graphs, those systems are based on the MapReduce model and suffer from the above limitations. As a consequence, iterative graph processing systems started to emerge in 2010 with Google’s Pregel [Malewicz, et al 2010], a graph processing framework that uses Valiant’s Bulk Synchronous Parallel (BSP) processing model [Valiant 1990] for its computation. Pregel was the first system that promoted a “Think Like A Vertex” notion for processing large graphs, similar to MapReduce that operates on <key,value> pairs to process large data volumes. These and other contemporary graph processing systems are discussed further in this survey.

2.1 Overall Scheme of Graph Processing

In general, typical graph processing systems execute a graph algorithm over a graph dataset across different logical phases, as shown in Figure 1.

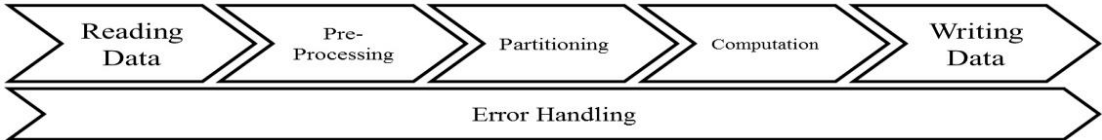


Fig. 1. Graph processing phases.

- (1) *Read/Write input/output datasets:* The first step is reading the graph data from a source dataset which can be either on disk or in memory. In the last phase, the processed data should be written back again, either to disk or memory. Graph processing systems typically do not have a custom persistence layer optimized for reading and writing graph datasets, and tend to use the standard file system, HDFS or some out of bound mechanisms. Hence, they can present a bottleneck when reading and writing large graph datasets. Many studies even ignore read/write time when they measure the execution time for evaluation. Instead, they try to improve other aspects of the systems like efficient in-memory data structures for computation.
- (2) *Pre-processing:* In some approaches, the graph data will be partitioned before being passed to the graph processing system to decrease the overall burden and runtime of the system. The main advantage of this approach is that the programmer does not need to worry about the complexity of the partitioning and it is a one-time cost that is paid upfront. Also, partitioning mechanism and computation mechanism can be two different modules which work independently and thus can be designed and implemented separately. The main drawback for this approach is that it works well only for static partitioning strategies, not dynamic partitioning or repartitioning.
- (3) *Partitioning:* In this phase, partitioning will be done dynamically within the graph processing system and not as a separate module. Both partitioning and computation phases can collaborate to choose the best partitioning method at each step, so, dynamic partitioning and repartitioning can be implemented in the processing system. Although programming

35:4

S. Heidari et al.

such a system is more complicated than implementing two independent modules, it provides more runtime flexibility and can be well suited to support diverse graph algorithms.

(4) *Computation*: Different graph processing systems have different computation approaches. This programming model and runtime is at the heart of the whole framework and there have been many proposals for efficient computation methods to decrease the graph application's runtime. More details about this phase, with a taxonomy on various computational models, is presented in Section 3.4.

(5) *Error Handling*: This fault-tolerant and failure recovery phase will be applied to the system either during the computation phase or after the computation phase is completed. There are various techniques that can be used here, such as check-pointing or restarting applications. Typically, the time taken for error handling is not considered in experimental results due to the overheads it causes and some frameworks even avoid considering this capability. However, given the use of large commodity clusters that are prone to failures and long running big data applications, fault tolerance is essential for graph processing frameworks used in an operational setting.

## 2.2 Large Graph-Oriented Applications

As noted in Table 1, there are many fields and applications that generate and provide big data in the form of graphs. With improvements in computer hardware and processing models such as cloud computing and emerging concepts like IoT, an even greater growth in datasets is expected. Here, we present some typical applications and environments where large graph data are generated and used.

(1) *Social Networks*: Social Networks and applications have grown exceedingly popular during the past decade and are constantly adding features to make effective use of the data they collect and to grow their customer network. Social networks are an important source of big graph data, and even big data in general, with large amounts of data created every day [Commission 2010]. People are sharing their personal activities with their friends and the whole world, talking about their beliefs, sharing photos and videos, and posting their interests and health information [Cha et al 2007] [Stelzner 2015]. In the year 2014, in each minute, 2,460,000 content posts are shared on Facebook, 3,472 photos are pinned on Pinterest, 72 hours of new videos are uploaded to YouTube, 278,000 tweets are shared on Twitter, 20 million photos are viewed on Flickr. These rates continue to grow, and form just a part of the whole big data social network landscape [Gunelius 2014].

Social networks are native generators and consumers of graph datasets, with an additional temporal dimension added to them. "Users" form the vertices of a huge social graph while "friendship" connections between them form the edges of the graph. Connections can be probabilistic and node's states change over time. Each node or edge can contain different values and information about a member's personal details, his/her interests, friends, groups and people, his/her followers, the pages are visited, locations, business information and so on with many other meta-data about his/her history of activities. All these form a digital trail for every user that needs to be processed and analyzed by social network providers.

From the users point of view, the network provider needs to suggest relevant pages or communities for them to follow based on their interests, and offer meaningful service offerings [Akbari et al 2013] [Bagci and Karagoz 2015]. The providers themselves benefit from leading users through targeted advertising to paid services. Although popular, social network sites are still in their infancy as they figure out how to monetize this massive dataset they have access to and make their business model sustainable. New methods and mechanisms are emerging in the area of analyzing social network data on distributed systems, clusters and clouds [Leimbach et al 2014].

(2) *Computer Networks and the Internet*: Every machine in a computer network, including clients, servers, routers and switches, is a node of a network graph and physical or network connections between these machines form the edges of the network graph. When various networks from all over the world are connect together to provide different services, it forms the "Internet" which is an extremely large graph [Chen et al 2005]. Computer networks need to be analyzed to discover whether there might be intruders, resource wasters, low efficiency, dead paths, and also to gain statistical reports about the states of the network [Ammann et al 2002]. This is particularly the case as a bulk of the network traffic moves toward rich content such as streaming video and multi-player gaming. These types of graphs should be processed in real-time as their state changes, and need a fast response, say to configure switches to allocate bandwidth to traffic, or detect malware and denial of service attacks. Network delays lead to customer dissatisfaction or worse, outages can cripple the functioning of modern society.

(3) *Smart Utilities*: Many large graph datasets are owned by public utility and service providers such as city and rail roads, and power and water grids. Take city road datasets as an example. Logistics companies need to find the shortest path between cities and streets to decrease their fuel consumption and ensure timely delivery of their goods, Governments need to plan maintenance and provision emergency services in case of power disruptions or natural disasters, and people need to find the most convenient means for travel between different locations [Lochert et al 2005].

Further, with a wider deployment of IoT and city services getting smarter [Paul 2013], the ability to monitor and collect real-time information about these physical infrastructure networks will grow, and graph analytics will be essential



for ensuring the smartness of these utilities. In fact, IoT will be a natural extension and an exponential expansion of the Internet. Graph applications can be used to drive real-time management of power grid operations with back-to-grid intermittent renewables like solar and wind, pumping operations for water networks, signaling of traffic lights based on current flow patterns, and even scheduling of public transit on-demand.

There are many other examples such as telecommunication [Marburger and Westfechtel 2010], web search engines [Page et al 1998], environmental analysis [Committee on the Analysis on Massive Data 2013], astronomy [Szalay 2011], mobile computing [Tian et al 2002], machine learning [Zha et al 2009] and so on where large graph data is required to be processed and, as we mentioned before, traditional approaches are not suitable.

2.3 Algorithms and Challenges

We discuss algorithms that are commonly used in most large graph processing studies and experiments. Table II shows the taxonomy of algorithms according to [Dominguez-Sal et al 2010], which are used in a number of works.

Table II. Graph algorithms categorization

Traversal	Breadth first search (BFS) Single source shortest path (SSSP)
Graph Analysis	Diameter Density Degree distribution
Components	Connected Components Bridges Triangle Counting
Communities	Max-flow min-cut K-means, Semi Clustering
Centrality Measures	PageRank Degree centrality Betweenness centrality
Pattern Matching	Path/subgraph matching
Graph Anonymization	K-degree anonym. K-neighborhood anonym.
Other Operations	Structural equivalence, Similarity, ranking, etc.

- (1) *Graph Traversal Algorithms*: These algorithms travel through all the vertices in a graph according to a specific procedure to check or update the vertices' values [Sedgewick and Wayne 2011]. Amongst the most common algorithms in this type are Breadth-First-Search (BFS) and Depth-First-Search (DFS) [Kozen 1992]. Both of these algorithms traverse the graph tree to find a particular node, or visit every node in a specific order. Single Source Shortest Path (SSSP) is used to find the shortest path between a particular node and any arbitrary node of the graph that might be based on the minimum cost or weight [Roy 2014]. Dijkstra's algorithm and Bellman-Ford algorithm are popular algorithms in this category [Bannister and Eppstein 2012].
- (2) *Graph Analysis Algorithms*: These algorithms peruse the topology of the graph to specify graph objects and analyze its complexity. These graph statistics and topological measures are extensively used in protein interplay analysis and social network analysis.
- (3) *Components*: Connected components algorithms find subgraphs in which a path exists between any two nodes in the subgraph and none are connected to nodes in other subgraphs [Hirschberg et al 1979]. So, each vertex only belongs to one connected component of the graph. Weakly connected components work on undirected graphs, while strongly connected components are relevant to directed graphs. Another component identification problem is counting triangles.
- (4) *Communities*: A community is a set of vertices in which each vertex in the community is closer to other vertices of the same community than any other vertices of the graph. Various topological and attribute measures can be used to define the closeness and quality of communities, and K-Means clustering and semi-clustering are popular algorithms in this category.
- (5) *Centrality Measures*: The aim of these algorithms is to give an approximate indication of the importance of a vertex in its community according to how well it is connected to the network. The most used algorithm of this type is PageRank (Page, Brin, Motwani, & Winograd, 1998), an algorithm that used by Google search engine to rank websites. Betweenness centrality is another common metric.
- (6) *Pattern Matching*: These algorithms are used to recognize the presence of input patterns in the graph, which can be an exact or approximate recognition.
- (7) *Graph Anonymization*: These algorithms are used to create a new graph based on an original graph where the latter emulates specific topological or attribute properties of the original one. This prevents any possible intruders to re-identify the network.
- (8) *Other Operations*: There are also other algorithms such as random walk algorithms where we choose a vertex randomly from neighbors of a vertex to start or continue process from there and try to converge in a probabilistic point [Fouss et al 2007].

35:6

S. Heidari et al.

In another categorization used by [Han et al 2013] and [Kang et al 2011], algorithms have been categorized based on the types of graph queries which result in two classes of algorithms (other types of query classification can be found in [Sarwat et al 2013] [Jamadagni and Simmhan 2016]):

(1) *Global Queries*: These queries need to traverse the whole graph. So, algorithms such as diameter estimation, PageRank, connected components, random walk with restart (RWR), degree distribution, etc. are in this group.

(2) *Targeted Queries*: These queries only need to access part of the graph, not all the graph. [Kang et al 2011] has formulated seven types of queries including neighborhood (1-step and n-step), induced subgraph, egonet (1-step and n-step), k-core and cross-edges.

Although there are many algorithms that can be implemented on a graph processing system, there are some challenges that these algorithms faced. First, according to [Lumsdaine et al 2007], many graph systems have limited memory that can be exclusively allocated to the processing algorithm, in addition to other processes and threads that simultaneously use and access the memory. Graph algorithms, in particular those that operate in a shared-memory system, can exceed available physical memory for large graphs processed on single machines [Murphy and Kogge 2007]. Recent graph processing systems address this by using a distributed computing paradigm. Second, the level of granularity in an algorithm can influence the level of parallelism it can exploit, especially those with linear runtime. So, a more fine-grained level of parallelism results in better scaling of such algorithms [Hendrickson and Berry 2008]. Third, algorithms should deal with diverse workloads and need to reassign tasks to processors when the visited nodes in a graph algorithm have spatial locality in the global memory. Finally, graph processing systems and algorithms should deal with the additional degree of parallelism exposed by submitting multiple concurrent queries when working on a large graph, or algorithms that operate over dynamic graphs. However, most of the systems that we review operate one graph algorithm or query over a single (large) graph.

### 3. GRAPH PROGRAMMING MODELS

We present different dimensions of graph programming and computation models, and classify and analyze prominent literature on graph frameworks based on these categories. A comprehensive list of graph processing systems based on this taxonomy is tabulated in Table 3.

#### 3.1 Graph Processing System Architectures

Graph processing systems can be categorized into three types of architecture models as depicted in Figure 2.

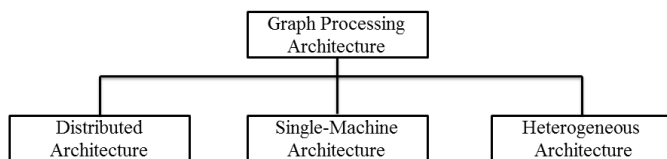


Fig. 2. Graph processing architectures.

##### 3.1.1. Distributed Architecture

A distributed system includes several processing units (host) and each host has access to only its own private memory. Each partition of the graph is typically assigned to one host to be processed while the hosts interact with each other by explicit or implicit message passing [Strandmark and Kahl 2011]. Such systems are meant to weakly scale by supporting larger graphs as more hosts are added to the system. From a cloud computing point of view, these map to an infrastructure as a service (IaaS) [Buyya et al 2009] architecture, where the hosts are Virtual Machines (VMs). Distributed graph processing systems utilize master-slaves (workers) architecture, as shown in Figure 3, where there is one master that is responsible for managing the whole system, assigning partitions to workers, managing fault-tolerance, coordinating the operations of the workers and so on; and there are multiple workers that are responsible for performing computation on the partitions.

Although the programmer has to adapt their algorithms and applications to suit the abstractions provided by the distributed graph processing systems, such systems ease the scaling of the applications on distributed environments, without the challenges of races and deadlocks that are associated with distributed computing [Coulouris et al 2012]. In contrast, shared-memory frameworks that have been developed for single machines are easier to program but are limited by their ability to hold only parts of large graphs in memory [Shun and Blelloch 2013].

Google's Pregel is a distributed vertex-centric framework that uses a master-worker architecture on multiple hosts of a cluster. GraphLab, developed in Carnegie Mellon University and later supported by GraphLab Inc., was developed for single machine processing [Low, et al 2010], but evolved into a distributed one [Low, et al 2012]. There are other Pregel-like systems such as GPS [Salihoglu and Widom 2013], Mizan [Khayyat et al 2013] and GoFFish [Simmhan et al 2014],

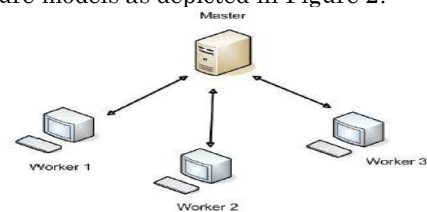


Fig. 3. Master-workers architecture.

and non-Pregel-like systems such as Presto [Venkataraman et al 2013], Trinity [Shao et al 2013], and Surfer [Chen et al 2010], which have been developed as distributed graph processing systems. Even frameworks such as GraphX [Xin et al 2013] are built on top of Spark distributed dataflow system. All of these systems use multi-node clusters or cloud VMs for their execution environment. However, as yet none of these exploit the elasticity property of Clouds, and rather treat captive VMs as a commodity cluster.

Beside the aforementioned graph frameworks, there are several graph processing libraries developed for high performance computing (HPC) clusters. Boost graph library (BGL) [Siek et al 2002] is a generic graph processing library that provides generic interfaces to the graph's structure and common operations, but hides the details of its implementation. This allows graph algorithms using BGL to have interoperable implementations on shared-memory and parallel computing platforms. Graph500 [Graph500] is a graph processing benchmark by which various metrics of supercomputers such as communication performance, memory size for graph storage and the performance of random access to memory are measured. It contrasts with Top500 [Top500] which is designed for compute-intensive applications. Although there have been many other attempts for providing parallel graph frameworks for high performance computing including several libraries such as MPI [El-Rewini and Abd-El-Barr 2005], PVM [Geist et al 1994], BLAS [Lawson et al 1979], JUNG [Java Universal Network/Graph Framework] and LEAD [Mehlhorn and Näher 1995], none of them provide the required flexibility for a general-purpose graph processing platform [Gregor and Lumsdaine 2005].

3.1.2. Shared-memory Architecture

Prior to the recent growth in distributed graph processing systems, there have been several works on processing large scale graphs on a single machine. A single machine consists of one processing unit (host) which can have one or more CPU cores, and physical memory that ranges from a few to hundreds of gigabytes that is shared across all the cores.

In 2012, Microsoft researchers conducted a study [Rowstron et al 2012] on whether using Hadoop on a cluster for analyzing big data is the right approach for data analytics. They concluded that for many data processing tasks, a single machine with large memory is more efficient than using clusters. They also investigated the cost aspect of using a single machine in big data processing and mentioned that "... for workers that are processing multi gigabytes rather than terabytes+ scale, a big memory server may well provide better performance per dollar than a cluster." [Lorica 2013].

Shared memory frameworks are inherently limited in the amount of memory and CPU cores present in that single machine [Doekemeijer and Varbanescu 2014]. The main challenge is that single hosts often have limited physical memory whereas processing large, real-world graphs can require a significant amount of memory to retain them fully in memory for many graph applications, or keeping and managing a subset of the graph out-of-memory. Novel techniques to address this limitation have been proposed.

In Strata Startup Showcase 2013, SiSense, which is a business intelligence solutions provider company, won the audience award with a software system called "Prism" that can exploit a terabyte of data on a single machine with only 8 GB of RAM [Lorica 2014]. It relies on disk for storage, transfers data to memory when needed and benefits from L1/L2/L3 caches of the CPU. It utilizes a column store and an interface which allows scalability to a hundred terabytes. Among major IT companies, for instance, Twitter uses Cassovary [Twitter 2012], an open-source graph processing system that has been developed to handle graphs that fit in the memory of a single machine. It has been claimed that Cassovary is a viable system for "most practical graphs" because of using a space efficient data structure. WTF (who to follow) [Gupta et al 2013] is a recommendation algorithm which is used by Twitter to suggest users with common interests and connections that implemented on Cassovary.

GraphChi [Kyrola et al 2012] is a vertex-centric graph processing framework that proposes a parallel sliding window (PSW) method for leveraging external memory (disk) and is suited for sparse graphs. PSW needs a small number of sequential disk-block transmissions, letting it to perform well on both SSD (solid state drive) and HDD (hard disk drive). Besides, GraphChi can process an ongoing in-flow of graph updates while performing advanced graph mining algorithms simultaneously, like Kineograph [Cheng et al 2012]. GraphChi uses space-efficient data structures such as a degree file that is created at the end of processing to save in-degree or out-degree for each vertex as a flat array. It also uses dynamic selective scheduling that lets *update function* and *graph amendments* to enlist vertices to be updated. It was extended later as a graph management system called GraphChi-DB [Kyrola and Guestrin 2014] and tried to address some of these challenges.

Many other graph processing systems have been developed based on single machines. Signal/Collect [Stutz et al 2010], for example, is a vertex-centric framework made to improve the semantic web computational performance. In this model signals will be sent along edges where they will be collected in vertices. The advantage of this model is that it provides flexibility for synchronous, asynchronous and prioritized execution. Other systems such as RASP [Yoneki et al 2014], X-stream [Roy et al 2013] which provides an edge-centric framework, FlashGraph [Zheng et al 2015], Galois [Nguyen et al 2013], TOTEM [Gharaibeh et al 2013], BPP [Najeebullah et al 2014], etc. also make processing graphs possible on single machines using various computational models and processing systems.



35:8

S. Heidari et al.

Graph processing on a single machine would be much cheaper and easier to program than processing on distributed systems, due to efficient communication and simpler debugging and fault-tolerance, if the entire graph fits within the local resources on that machine. But this limits their scalability beyond a certain graph size. New approaches for processing graphs on single machines are targeting SSDs [Yamato 2015, Koo et al 2015] whose speeds are matching main memory and offer advantages for graph processing [Zheng et al 2015, Nilakant et al 2014].

### 3.1.3. Heterogeneous Architecture

In a heterogeneous environment, not every processing unit is equally powerful [Guo et al 2015]. This may be a single machine and additional on-board accelerators and specialized devices, or it can also consist of distributed, non-homogeneous systems. Because of this, we considered them as a separate group in this taxonomy. For example, processing systems such as RASP [Yoneki et al 2014] and FlashGraph [Zheng et al 2015] have tried to optimize the storage part of the system by using SSD which is much faster and more reliable than traditional hard drives [Geier 2015]. Many graph processing systems have proposed utilizing graphic processing units (GPU) alongside CPU for computation [Zhang et al 2015]. Medusa [Zhong and He 2013], for instance, was developed to make processing graphs using GPUs easier. Medusa is a programming framework that enables users to write C/C++ APIs to promote the capabilities of GPUs to execute the APIs in parallel. Its extended version also can be run on multiple GPUs within a single machine. Gharaibeh et al [2013] developed a system called TOTEM that assigns the low-degree vertices to the GPU and operates high-degree vertices processing on the CPU. On the other hand, systems like CuSha [Khorasani et al 2014] compute the entire graph on GPU. Another possibility is to exploit non-uniform VM sizes on Clouds for a distributed, heterogeneous architecture, which has been less explored.

## 3.2 Programming Models for Graph Processing Frameworks

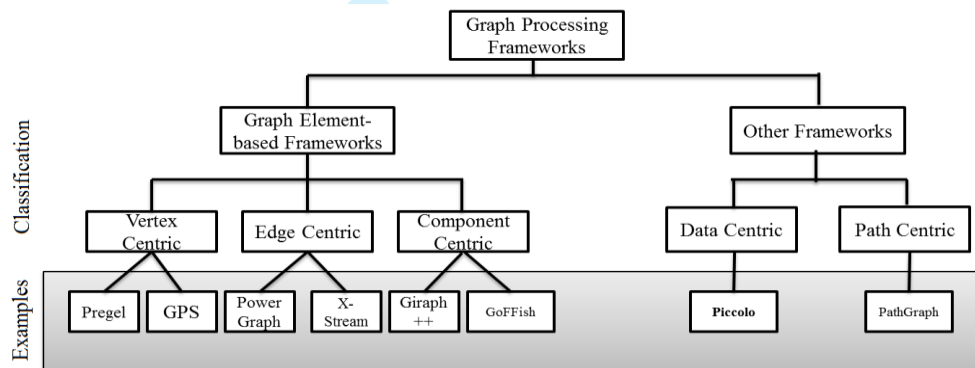


Fig. 4. Taxonomy of programming models used by graph processing frameworks

Graph processing frameworks enable graphs to be processed on different infrastructures such as clusters and clouds. Here, we restrict ourselves to distributed memory systems that are designed for commodity, rather than High Performance Computing or Supercomputing clusters. The programming abstraction for each framework is designed either based on a graph topology element, such as vertices and edges, or other alternative approaches. Figure 4 depicts the taxonomy of graph processing frameworks according to main characteristics of the graph and other alternatives. We discuss these further below.

### 3.2.1. Vertex-Centric Frameworks

Vertex-centric programming is the most mature distributed graph processing abstraction and several frameworks have been implemented using this concept [Doekemeijer and Varbanescu 2014]. A vertex-centric system partitions the graph based on its vertices, and distributes the vertices across different partitions, either by hashing them without regard to their connectivity [Malewicz et al 2010, Apache Giraph] or by trying to reduce the edge cuts across partitions [Salihoglu and Widom 2013]. Edges that connect vertices lying in two different partitions either form remote edges that are shared by both partitions or owned by the partition with the source vertex.

In the vertex-centric programming abstraction introduced by Google's Pregel [Malewicz et al 2010], computation centers around a single vertex – its state and its outgoing edges – and interactions between vertices are through explicit message passing between them. This gives a fine-grained degree of vertex-level data parallelism that can be exploited for concurrent execution. Pregel's execution follows a bulk synchronous parallel (BSP) model, where vertex computation and inter-vertex messaging are interleaved, and the application iteratively progresses along barrier-synchronized supersteps. The Pregel API allows developers to focus on the vertex-centric graph algorithms while abstracting away communication and coordination details to the runtime. In Pregel, the domain of a vertex's user-defined *compute* function is restricted to

the vertex and its outgoing edges, while LFGraph [Hoque and Gupta 2013] considers incoming edges to be restricted. Figure 5.b shows vertex-centric processing approach for a sample graph shown in Figure 5.a.

A vertex-centric model makes programming of graph processing *intuitive and easy*, similar to the advantages of Map-Reduce for tuple-centric programming. Parallelization is done *automatically*, and *race conditions* on distributed execution are avoided. Primitives like *combiners* and *aggregators* are available for application-level message optimizations and global state exchange. The model also allows for *graph mutations*, where the structure of the graph can be changed as part of the execution (useful, for e.g., when iteratively coarsening the graph for partitioning, clustering or coloring).

However, Pregel does have a few shortcomings: 1) While the vertex-centric model exposes parallelism at the level of individual vertices, which can be computed in negligible time, massive graphs can impose coordination overheads on this degree of parallelization that may out-weigh the benefits [Tian et al, 2013], 2) The number of barrier-synchronized supersteps taken for traversal algorithm can be proportional to the diameter of the graph with the number of message exchanges required between partitions also being high, proportional to the number of edges [Simmhan et al 2014], 3) Mapping shared memory graph algorithms to this model is not trivial and requires new vertex-centric algorithms to be developed [Simmhan et al 2014] and 4) Using a vertex-centric programming model without regard to the graph partitioning and data layout on disk can lead to punitive I/O initialization and runtime performance [Simmhan et al 2014].

Apache Giraph [Apache Giraph], is a popular open-source implementation of Pregel. Giraph uses Map-only Hadoop jobs to schedule and coordinate the vertex-centric workers and uses Hadoop distributed file system (HDFS) for storing and accessing graph datasets. It is developed in Java and has a large community of developers and users such as Facebook [Jackson 2013, Salihoglu et al 2015]. Giraph has a faster input loading time compared to Pregel because of using byte array for graph storage. On the other hand, this method is not efficient for graph mutations which lead to decentralized edges when removing an edge. Giraph inherits the benefits and deficiencies of the Pregel vertex-centric programming model. Its performance and scalability is algorithm and graph dependent, and works very fast, for e.g., on stationary algorithms like PageRank but not as fast on traversal algorithms like single source shortest path (SSSP) [Roy 2014] and weakly connected components (WCC) [Salihoglu and Widom 2014], particularly for graphs with a large diameter. However, the ease of use of this framework and the community support has made it a popular platform over which to develop other Pregel-like systems with feature enhancements to the vertex-centric concept.

Other distributed platforms like Apache Hama and GraphX also offer a vertex-centric programming model, with features comparable to Giraph. GraphX, developed on top of Apache Spark, determines *transformation on graphs* where every operating action produces a new graph. This framework uses a programming abstraction called Resilient Distributed Graph (RDG) interface, which builds upon Spark’s in-memory storage abstraction – Resilient Distributed Datasets (RDD). The graph in GraphX includes the *directed adjacency structure* along with *user defined attributes* connected to each node and edge, and both are encoded as RDGs. Using RDG, the implementation of frameworks such as Pregel and PowerGraph on Spark needs less efforts.

Pregelix [Bu et al 2014] is vertex-centric framework that tries to model Pregel as an iterative dataflow on top of Hyracks [Borkar et al 2011] parallel dataflow engine. Pregelix has been developed to address three main challenges in distributed Pregel-like systems: 1) Many Pregel-like systems have limitations to support out-of-core vertex-storage, 2) Existing Pregel-like systems have specific strategies and implementations for communication, node storage, message delivery, and so on. Therefore, a user cannot choose between different implementation strategies based on what is better for a particular algorithm, dataset or cluster. Pregelix improves physical flexibility and scalability of the processing system to address this challenge, and finally, 3) Pregelix tries to leverage current data-parallel platforms to streamline the implementation of Pregel-like systems.

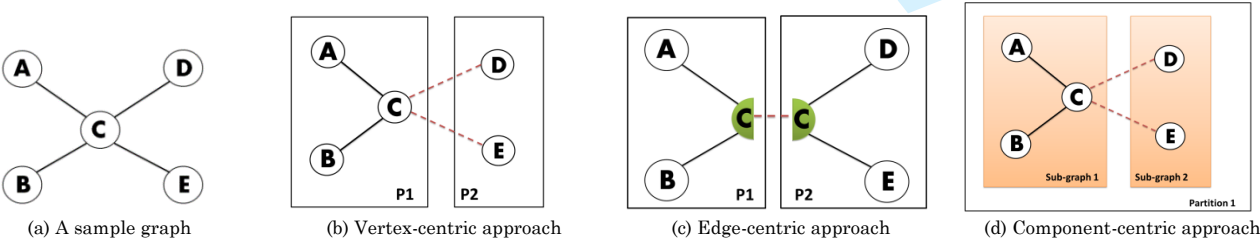


Fig. 5. Graph element-based approaches for graph processing frameworks

3.2.2. Edge-Centric Frameworks

In edge-centric frameworks, edges are the primary unit of computation and partitioning, and vertices that are attached to edges lying in different partitions are replicated and shared between those partitions. It means that each edge of the graph will be assigned to one partition, but each vertex might exist in more than one partition. Figure 5.c depicts this approach. While edge-based partitioning is more costly, this model shows better graph processing performance compared

35:10

S. Heidari et al.

to vertex centric approaches [Rahimian et al 2014]. However, programming an edge-centric system is more difficult than vertex-centric systems [Yuan et al 2014]. It is also important to create edge-balanced partitions in this method to load balance the computation across workers, just as vertex balancing is important for vertex-centric frameworks. Decreasing the vertex cuts have been investigated in some research [Feige et al 2005, Beseri Sevim et al 2012, Liu et al 2006].

[Catalyurek and Aykanat 1996] and [Devine et al 2006] have suggested a vertex-cut method for distributed graph placement in hyper graph partitioning, where the edge-centric problem can be solved by converting each edge into a vertex and vice versa. PowerGraph [Gonzalez et al 2012] is the first system to use an edge-centric model of computing that was developed in Carnegie Mellon University. The motivation for developing PowerGraph is that systems such as Pregel and GraphLab [Low et al 2010] are effective for flat graphs but have shortcomings with graphs that follow a power law edge degree distribution due to low quality partitioning and vertices with high edge degrees. Real-world graphs such as social networks are such power-law graphs where a small set of vertices have high edge degrees that connect to a large part of the graph, e.g. celebrities in social networks, or hubs in airline networks. Partitioning and representing power-law graphs in a distributed environment is also difficult [Leskovec et al 2008, Abou-Rjeili and Karypis 2006].

PowerGraph [Gonzalez et al 2012] was originally implemented as a vertex-centric model, with a Gather-Apply-Scatter (GAS) decomposition (Section 3.4) to parallelize high-degree vertices and effectively distribute large power-law graphs. But later, this approach changed to an edge-centric model due to strictly lesser communication and storage. PowerGraph also borrows features from Pregel and GraphLab. It offers developers a shared-memory view of computation based on GraphLab, which is lacking in a vertex-centric model, to eliminate an explicit synchronized message passing model of information flow by the user; and it uses the “gather” notion from Pregel. Each vertex in PowerGraph has multiple copies. One of these copies at random is chosen as the master, which keeps the authoritative data for a vertex, and other replicas of that vertex are mirrors that keep a read-only version of the vertex data. PowerGraph has some shortcomings such as inefficient out-of-core storage and significant fault-tolerant overhead for large graphs. These are discussed in subsequent sections.

X-Stream [Roy et al 2013] is another well-known edge-centric system that processes out-of-core and in-memory graphs using a gather-scatter approach (Section 3.4). It bases this approach on the intuition that storage media such as solid state drives (SSD), main memory and magnetic disk perform significantly better with a sequential access to data than random access. The authors have implemented different algorithms on their system and observe that many of them can work on edge-centric mode. It can even return results from unsorted edge lists. However, it causes overheads when new edges are added to the graph. X-Stream is not suitable for very large graphs that do not fit onto the SSD, it wastes remarkable consecutive bandwidth for certain algorithms and finally, X-Stream is not suitable for graphs that their structure needs so many iterations [Yuan, et al., 2014].

### 3.2.3. Component-Centric Frameworks

Component-centric approaches have been recently introduced, where components are collections of vertices and or edges that are coarser than a single vertex or edge. Tian et al [Tian et al 2013] from IBM introduced a “Think Like A Graph” instead of “Think Like A Vertex” [Tian et al 2013] abstraction after observing shortcomings in vertex-centric and edge-centric methods of graph processing. In their partition-centric view, they divide the whole graph into partitions and assign those partitions to machines for being processed. A partition, which is a collection of vertices and edges in the graph, forms the unit of computing. Figure 5.d shows this approach. Giraph++, based on Apache Giraph [Apache Giraph], implements this model and uses this coarse-grained parallelism. In contrast to vertex-centric model that hides partitioning and component connectivity details from users, Giraph++ gives exposes the partition’s structure to the users to allow optimizations. So, the performance of the system depends on the partitioning strategy that is used and how effectively users exploit the access to the coarse components in their execution. On the other hand, communication within a partition is by direct memory access, which is faster than passing messages between each single vertices in vertex-centric model. This results in fewer network message passing and lower time of execution per iteration (superstep), with a reduction in the number of iterations needed for convergence. It also benefits from local asynchrony in the computation which means that vertices in the same partition can exchange their state and perform consequent computations to the extent possible in the same iteration.

Simmhan et al developed GoFFish [Simmhan et al 2014] which has a subgraph-centric computation model to merge both the scalability and flexibility of vertex-centric programming approach with the extensibility of shared-memory subgraph computation. A subgraph (weakly connected component) is the unit of computation. A partition may contain one or more subgraphs, whereas each subgraph only belongs to one partition of the graph. Vertices in the same subgraph have a local path between each other, so existent shared memory graph algorithms can directly be exerted to each subgraph. This gives a programming and algorithm design advantage over partition-centric frameworks like Giraph++ that offer no guarantee on connectivity between vertices in a partition, while retaining the advantages of fewer iterations and shared-memory access of those frameworks. Subgraphs, or vertices that span subgraphs, communicate by passing messages, similar to a vertex-centric model.

GoFFish consists of two major components: 1) A distributed graph oriented file system, *GoFS*, which partitions, stores and provides accesses to graph datasets in a cluster across hosts, and 2) A subgraph-centric programming framework, *Gopher*, which executes applications designed using the subgraph-centric abstraction using the Floe [Simmhan and Kumbhare 2013] dataflow engine on top of GoFS. However, subgraph-centric programming algorithms are vulnerable to imbalances in the number of subgraphs per iteration as well as non-uniformity in their sizes. The time complexity per iteration also can be larger since it often runs the single machine graph algorithm on each subgraph, even as it often takes much fewer iterations. The benefits are also more pronounced for graphs with large diameter, where algorithms tend to be several times faster than a vertex-centric equivalent, rather than small-diameter power-law graphs. GoFFish supports applications that operate on single property-graphs as well as on time-series graphs [Simmhan et al 2014]. Frameworks like Blogel [Yan et al 2014] also adopt connected components (blocks) as units of computation.

3.2.4. Other Graph Frameworks

In addition to the aforementioned programming abstractions, other alternatives have been developed as well. Several data-centric models offer a declarative dataflow interface to users to access and process data without needing to explicitly define communication mechanisms. For example, MapReduce provides a dataflow programming model that is popular for processing bulk on-disk data, but not for in-memory computations across multiple iterations, and applications do not have online access to the intermediate states. Piccolo [Power and Li 2010], was developed in New York University as a data-centric programming method for writing parallel in-memory applications in several machines. It uses a key-value interface with a user-defined accumulator function that automatically combines concurrent updates on the same key. Like many other data-flow models such as Pig, Hive, Dryad [Isard et al 2007, Yu et al 2008], Flume Java [Chambers et al 2010] and Swazal, developers in Piccolo operate at a higher level of dataflow programming abstraction but need to know the framework and system behavior well to leverage its scalability for different applications. For e.g., the programmer has to *a priori* specify the number of partitions while creating a table. Further, these are tuple-oriented data flow models rather than graph specific ones.

Yuan et al [Yuan et al 2014] introduces PathGraph which aims to leverage memory and disk locality on both out-of-core and in-memory graphs using a path-centric approach. Their path-centric abstraction utilizes a set of tree-based partitions to model the graph and benefits from a path-centric computation instead of a vertex or edge centric computation. It means that the graph will be partitioned into paths including two forward and reverse edge traversal trees for each partition. It applies iterative computation per traversal tree partition in parallel, and then merges partitions by examining border vertices. Two functions, *gather* and *scatter* (Section 3.4), are used to traverse each tree by a user-defined algorithm. In addition to the computation tier, PathGraph has a path-centric storage tier to better the local accessibility for the computation. The storage structure is based on a tree partition and uses vertex-based indexing for tree-based edge chunks. The system outperforms the vertex-centric GraphChi [Kyrola et al 2012] and edge-centric X-Stream frameworks.

3.3 Distributed Coordination

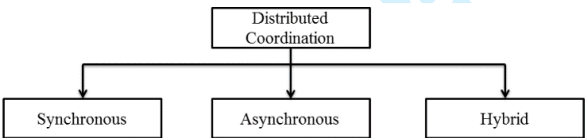


Fig. 6. Distributed coordination

3.3.1. Synchronous

When a graph algorithm executes synchronously, it means that concurrent workers process their share of the work iteratively, over a sequence of *globally coordinated* and well-defined iterations. Synchronization may be applied to vertex-centric, edge-centric and component-centric models, and both on distributed and single machine systems. For example, Pregel-like systems call their barrier synchronized iterations as superstep, and workers coordinate their computation and communication phases in each superstep – everyone completes a superstep before starting the next. Initially, the master assigns partitions to the workers in the first iteration; the workers update their set of vertices based on the assigned partitions and wait for a global barrier, which tells them all workers are ready with their partition [Malewicz et al 2010]. Subsequent supersteps indicate actual computation based on the application logic. Updated vertices in each partition send messages to (typically) vertices in neighboring partitions between iterations. Within an iteration, vertices can only access information about their local vertex’s state and messages received from the previous iteration. Such a synchronized execution is possible even in a shared-memory system, across workers (threads, processes) on a single server.

These regular intervening periods make the system appropriate for algorithms where sizeable computation and communication can take place within each iteration since there is an overhead associated with the coordination. The bulk messaging at iteration boundaries can utilize the bandwidth efficiently if there is heavy communication between partitions [Ediger and Bader 2013, Xie et al 2015]. It is easy to program, debug and deploy such systems, without concerns of distributed race conditions and deadlocks. Another advantage of synchronous processing is that the outcome of each



35:12

S. Heidari et al.

superstep is known immediately and provides real-time response of incremental application progress and easier error recovery in case at superstep boundaries. Synchronous execution is also suitable for balanced workloads that are computed symmetrically, with all workers having adequate work, so that the overhead of the global barrier and idle time for faster workers waiting for slower workers to synchronize is reduced [Xie et al 2015]. These advantages make synchronous execution very popular such that several graph processing systems like Pregel, GPS, Kineograph, Mizan, GasCL [Che 2014] and Medusa [Zhong and He 2013] use this model. Some like GoFFish [Simmhan et al 2015] have two levels of such synchronized supersteps, an outer loop over different graphs in the context of time-series graphs, and an inner loop as supersteps over a single graph from the outer loop.

Synchronous execution model has some disadvantages as well that should be considered while designing or choosing a processing system for graphs. First, this model is not suitable for unbalanced workloads in which computation converges asymmetrically [Suri and Vassilvitskii 2011]. Likewise, if the distributed machines are not homogenous, the performance of the hardware may also cause some partitions to operate slowly. In such cases, it is possible to have stragglers where when all partitions have been computed on workers except one slower worker which has not finished and hence delays all workers in the superstep. So, the runtime in this model is completely dependent on the slowest machine in each iteration [Salihoglu and Widom 2013]. Some of these shortcomings have been identified and addressed through elastic load balancing of partitions across workers [Dindokar and Simmhan 2016]. Another drawback is that the intermediate processing updates between supersteps, in the form of messages or state, has to be retained in memory and this causes additional memory pressure [Redekopp et al 2013]. A third disadvantage is that a synchronous execution model is well suited for applications and algorithms that need coordination between adjacent vertices [Doekemeijer and Varbanescu 2014]. For example, in a graph coloring algorithm in which vertices try to choose a different color from their neighbors, two adjacent vertices might pick conflicting colors frequently [Gonzalez et al, 2012, Tasci and Demirbas 2013] and the algorithm will converge slowly. Lastly, based on the drawbacks mentioned above, the cost for the systems that use synchronous model of execution is higher because the throughput must always remain high and running time would be longer [Zhang et al 2014].

### 3.3.2. Asynchronous

An asynchronous execution model does not have any global barrier and a subsequent phase of execution will be started on a worker immediately after its current iteration finishes its computation [Xie et al 2015]. So, some of the challenges of load balancing and long tail computation in the synchronous model is addressed by asynchronous computation, where workers do not have to wait for the slowest worker to start their subsequent iteration. This approach is useful where the workload is imbalanced and convergence can occur faster than synchronous approach. Therefore, we can say that asynchronous model is the preferred model when computation across workers is heavily skewed and there is little communication that can benefit from bulk operations [Xie et al 2015]. In other words, this model is more preferable for CPU-based algorithms while synchronous model would perform much better on I/O-bound algorithms. Another advantage for this model is that it can use dynamic scheduling to implement prioritized computation to execute more units of computation before others, to obtain better performance [Zhang et al 2012]. Normally, asynchronous execution provides more flexibility than synchronous execution by utilizing dynamic workloads which makes it outperform synchronous methods in many cases; however, the exact comparison between these two models depends on various properties of the input graph, platforms that the system has been deployed on, execution stages and applications [Xie et al 2015]. Finally, using asynchronous approach provides a non-blocking process because resources could be free and become available for the next iteration, whereas in a synchronous approach, they are blocked until the global barrier declares the end of superstep which leads to a competition for resources at the beginning of next superstep.

As before, there are disadvantages to this model as well. The key disadvantage is that programming asynchronous processing systems is more difficult than synchronous systems. The programmer should deal with irregular communication intervals, unpredictable response time, complex error handling and more complicated scheduling issues. For example, for error recovery in such a system, many factors have to be considered: which machine has faced a fault, in which iteration of a particular worker the error happened, which resources caused the errors, should new resources be allocated to the computation or it should only be rearranged, and so on. This also results in more complex debugging and deployment, and careful programming to avoid deadlocks. In case of pull-based communication model (Section 4.2), which is usually implemented in an asynchronous manner, many redundant communication may happen because there are several intertwined reads and writes while adjacent vertices values do not change [Han et al 2014] [Zhang et al 2012]. On the other hand, regardless of these drawbacks, many single machine systems have preferred an asynchronous execution approach since the shared memory makes it easier to asynchronously use the latest data without waiting for a barrier.

### 3.3.3. Hybrid

There are many systems that use only synchronous execution mode; for example Pregel, GoldenOrb [Cao 2011], GBASE [Kang et al 2011], Chronos [Han et al 2014] and GraphX [Xin et al 2013], while many other systems utilize asynchronous mode like GiraphX [Tasci and Demirbas 2013], GraphHP [Chen et al 2014], Ligra [Shun and Blelloch 2013], RASP

[Yoneki et al 2014] and GraphChi. But recently a new approach called hybrid execution model has been implemented in a few systems that tries to take advantages of both asynchronous and synchronous approaches or incorporate them with new additional solutions. Such graph processing systems have developed to improve system performance by overcoming the shortcomings of existing methods, and use both synchronous and asynchronous models of coordination to benefit from their relative strengths.

GRACE [Wang et al 2013], for instance, is a single machine framework that combines synchronous programming with asynchronous execution features. It actually separates execution policies from application logic. In an asynchronous execution, a processing sequence of vertices can be intelligently ordered by dynamic scheduling to remarkably speed up the convergence of computation. GRACE uses the BSP computational model and message passing communication model as two primary paradigms of synchronous model. It helps GRACE to improve its automatic scalability by applying prioritized execution of vertices and receiving messages selectively outside of the last iteration. Various workloads like topic-sensitive PageRank, social community detection and image restoration have been used in GRACE and it shows comparable running time to other asynchronous systems such as GraphLab with even better scalability.

Another hybrid approach distinguishes between local vertices that are within a partition and remote vertices connecting across partitions. These type of systems exploit both local asynchronous computation when they still need global barrier for synchronization of remote vertices values [Chen et al 2014]. In the local computation phase, messages will be passed very fast across local vertices using shared memory. In the next phase (synchronous phase), remote nodes (boundary nodes) that are connected by edges across the partitions, will be updated by exchanging messages. In fact, component-centric frameworks such as Giraph++ and GoFFish allow users to exploit this. Others like Giraph Unchained [Han and Daudjee 2015] also allow incremental forward progress within a future superstep based on partial messages that are received, even before the previous superstep completes. These straddle synchronous and asynchronous models.

Apart from these methods, a novel approach has been introduced by Xie et al. [Xie et al 2015] in the PowerSwitch system, which sequentially switches between synchronous and asynchronous execution mode according to a heuristic prediction. This is because some properties of the processing might change as it progresses. For example, processing single source shortest path (SSSP) algorithm begins with just a few vertices active, which means that few messages are passed; this is suitable for an asynchronous model. But as the process goes further, the number of vertices involved in the computation will increase which means that the number of messages passing among them increase as well, and this is suitable for synchronous execution model. PowerSwitch can effectively predict the proper heuristic for each step and it can switches between the two modes if required.

3.4 Computational Models

Performing computation is at the heart of a graph processing system where data (vertex or edge) will be processed and updated. Computational models that are used in existing graph processing systems can be divided into two major groups: 1) two-phase models, and 2) three-phase models. Figure 7 shows the classification of these two models with examples from each group. Computational model of some systems is a combination of these methods with other approaches.

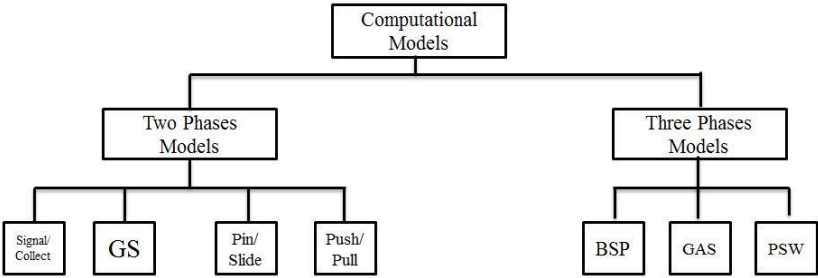


Fig. 7. Classification of computational models in graph processing systems

3.4.1. Two-Phase Computational Models

There are usually two functions that are applied on data (vertices or edges) in a two-phase computation model. Signal-collect approach is the first two-phase programming model for large scale graph processing on the semantic web within a system with the same name (signal/collect) [Stutz et al 2010]. Computation in the vertices are completed by collecting the signals that are coming from edges and performing some processing on them using the vertex's state and then sending (signaling) their adjacent vertices in the compute graph. Signal/collect has been implemented for working with both synchronous and asynchronous execution models. In both models, some parameters should be set to determine when the computation should be stopped: *signal\_threshold* and *collect\_threshold* parameters in which a minimum level of "importance" will be set for the execution, and a *num\_iterations* parameter which is the number of iterations in synchronous mode, and *num\_ops* that is the number of executed operations in asynchronous mode. All these parameters should be set by the user.

35:14

S. Heidari et al.

Scatter-Gather (GS) is another two-phase model that includes two major functions like signal/collect. *Scatter* function broadcasts a vertex value to its adjacent vertices similar to “signal” while *gather* function performs like “collect” to collect the values from neighbors and use them to update itself. X-Stream [Roy et al 2013] uses scatter-gather method as its programming model on its edge-centric framework. The edge-centric scatter grabs an edge as input, performs the computation according to its source node data field when the updated value should be sent to its destination vertex. The edge-centric gather, on the other hand, takes a vertex value as input and uses it to recalculate or compute its destination data field. A new concept called Scatter-Combine has been introduced in [Yan et al 2013] that performs like scatter-gather, but on the graphs that contain messages carrying data and also the operation that should be executed on it.

Pin-and-slide was first introduced by TurboGraph [Han et al 2013] for parallel execution of large datasets on a single machine. A pin-and-slide mechanism consists of a graph dataset, a buffer pool and two different threads, as explained in Section 3-1, callback threads and execution threads. When the buffer manager is being sent an asynchronous input/output by a callback thread, it sends the demand to the FlashSSD via the operating system after when the control of execution goes back to the calling thread immediately. The main goal in this system is to reach all related adjacency lists efficiently. To achieve this goal, first, the pages that contain these adjacency lists should be identified. The most important challenge here is pinning large adjacency pages (LA pages) which means that a number of smaller adjacency pages must be unpinned first, then the LA page can be pinned. To overcome this challenge, LA pages will be pinned when all of related LA pages for a big vicinity list are completely loaded to maximize the buffer exploiting. When execution threads or callback threads terminated the processing of a page, this page will be unpinned and an asynchronous input/output demand will be sent to the FlashSSD by the execution thread. As soon as the processing has been completed, the execution window can be slid by the size of the pinned pages in the buffer [Han et al 2013].

Nguyen *et al* [Nguyen et al 2013] have used another two-phase approach called push/pull styles. The value of an active vertex will be pushed (flowed) from that vertex to its neighbors, which is more like scatter phase. The pull style occurs when the data from an active vertex's neighbors flow into that vertex which is more like a gather phase. In an algorithm like SSSP, the push-style is applied to the active vertex neighbors by updating the destination label of the siding nodes of the active vertex by doing a relaxation with them; and the pull-style function updates the destination label of the active vertex by doing relaxation with all neighbors of that node. The pull mode also needs less synchronization because there is just one writer for each active vertex. KineoGraph [Cheng et al 2012] is another system that uses this model for computation.

### 3.4.2. Three-Phase Computational Models

Bulk synchronous parallel (BSP) is a parallel programming and also the most representative model in this category that has been used in several graph processing systems [Malewicz et al 2010, Salihoglu and Widom 2013, Vaquero et al 2013, Khayyat et al 2013]. To deal with the scalability challenges of parallelizing tasks across a number of workers, BSP, which utilizes a message passing interface, was developed. In BSP, as a vertex-centric computational model, each node is able to have two modes of “active” or “inactive”. The computation comprises a series of superstep that come with synchronization hurdle between them. So, in each superstep: 1) The node that is involving in computation, obtain its adjacent nodes updated values from last superstep (except the first superstep), 2) Then, the node will be updated based on the obtained values, and 3) The node forwards its updated value to its neighbors that will be available for them in the next iteration. In each iteration, a vertex may choose to vote to halt, in case it does not receive any messages from its neighbors or it has reached a locally stable state. That means it will not participate in the processing anymore until it receives new messages that convert its state from inactive to active. So, in each iteration only active vertices will be computed. If there is no active vertex in the graph, then the computation is finished.

Some research have modified the BSP model and introduced new models. For example, temporally iterative BSP (TI-BSP) [Simmhan et al 2015] is a computational model for time-series graphs on a subgraph centric model such as GoFFish. It has used BSP as a building block to support the design pattern. TI-BSP is a series of BSP loops (nested supersteps) in which each *outer loop*, as a *timestep*, runs on one graph instance in time. Supersteps using a subgraph programming model form the *inner loop* that operate over a single instance. As a result, the design pattern will be decided based on the order of timesteps execution and the messages between them.

There is another BSP model that stands for “BiShard Parallel” and has been introduced by the single machine based system, BPP [Najeebullah et al 2014], to empower full CPU parallelism for graph processing. This model has also three phases that include: 1) Loading a sub-graph of the large graph from disk, 2) Performing compute operation on the sub-graph and update the values of edges and vertices, and 3) Writing back the modified values on disk. BiShard Parallel performs under an asynchronous execution model and needs more disk space compared to one shard mechanism that was used in GraphChi, because two copies of each edge is managed in this model.

GAS (Gather-Apply-Scatter) is another three-phase computational model that was introduced by PowerGraph. The data about the adjacent nodes and edges is obtained and collected using a general summation over all adjacent vertices and edges of a vertex in the *gather* phase. The *apply* operation should be defined by user and must be both associative and

commutative, and can vary from a numeric summation to the aggregation of data across all adjacent edges and vertices. The results from gather phase is used to update the central vertices values in the apply phase. Finally, the recent data of the central vertex is used to renew the values on neighboring edges in the *scatter* phase. The critical challenge in this model is that graph parallel abstractions should be able to perform computations with high fan-in and high fan-out where both of them are specified by gather/scatter phases. GAS model is used to develop a runtime system mapping in parallel on GPUs as a graph application called GasCL [Che 2014]. This model is like the model have been used in systems such as Pregel and GraphLab, but in a different way.

GraphChi uses a different computation model called parallel sliding window (PSW). PSW is an asynchronous computation model that can efficiently process the graph with changeable edge value from disk, with a few number of non-consecutive disk access. PSW performs three phases for processing a graph as follow: it loads a subset of the graph from disk, then, applies update operation on vertices and edges, and eventually, the new updated values will be written on disk. The number of “reads” from disk is exactly equal to the number of “writes” to the disk in this model.

4. RUNTIME ASPECTS OF GRAPH FRAMEWORKS

4.1 Partitioning

Graph partitioning is a method in which graph data is divided into smaller parts with specific properties [Buluç et al 2013]. For example, in a *k*-way partitioning, the graph is partitioned into *K* smaller parts of equal size while minimizing the edge cuts between partitions. This is an NP-complete problem [Andreev and Racke 2004]. Graph partitioning is a fundamental research problem and several reviews have been done on different methods and perspectives of graph partitioning [Buluç et al 2013, Bader et al 2013]. In a graph processing system, partitioning is applied on the large graph in order to assign each smaller partition to a worker to be computed. The most important challenge in this context is “how do we partition the graph to achieve better cuts while taking load balancing and simplicity of computation into consideration?”.

Many novel heuristics have been proposed for partitioning large graphs. METIS [Karypis and Kumar 1995], for instance, is a popular tool that uses multi-level partitioning. It is able to perform high quality partitioning that decreases the overall communication (edge cuts) and reduce imbalances across partitions. However, METIS is computational costly and high random access needs make it unsuitable for large graphs. ParMETIS is a parallel MPI-based version of METIS that mitigates some of these performance limitations.

There are distributed partitioning algorithms, some of which have been implemented on top of graph processing frameworks as well. Spinner [Martella et al 2015], for instance, runs on top of Giraph and utilizes an iterative node migration approach using Label Propagation Algorithm (LPA) to deal with scalability and changing partitions. It allows Spinner to scale to billion-vertex graphs by avoiding costly synchronization among vertices. Blogel implements a *Graph Voronoi Diagram (GVD)* partitioner using a vertex-centric computing method by operating as a multi-source breadth-first search (BFS). It partitions the vertices into blocks using multi-source BFS over linear workloads.

Some of graph processing systems create additional topological constructs on top of the partitioned graph. In GoFFish, which is a subgraph centric framework, each partition may have more than one subgraph (weakly connected component), and these subgraphs by definition cannot have an edge between them. So GoFFish has a post-processing stage once the graph is partitioned, in which it identifies all subgraphs within a partition that form the unit of processing during the programming model’s execution.

In general, two partition creation strategies can help to improve the runtime performance during distributed graph processing: 1) Creating more partitions than workers and allocate more than one partition to each worker, and 2) allocating one partition per worker, yet using multiple workers on each processing host [Salihoglu and Widom 2013]. We next discuss alternative perspectives towards partitioning to support graph processing systems, also shown in Figure 8.

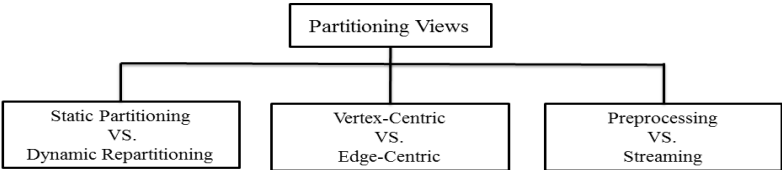


Fig. 8. Partitioning views in graph processing systems

4.1.1. Static Partitioning vs. Dynamic Partitioning

Several graph processing frameworks utilize static partitioning which means that they consider the graph and the processing environment to stay unchanged [Schloegel et al 2001]. These systems assume that the I/O bandwidth, latency, processing units and the graph itself are constant and predictable. So, this method of one-time *a priori* graph partitioning is easy to program and load balancing can be easily achieved, if the assumptions hold and the problem domain does not change [Elsner 2002].



35:16

S. Heidari et al.

On the other hand, dynamic repartitioning assumes that runtime behavior of an algorithm, the processing environment and even the graph itself can be variable. They try to repartition the current state of the graph according to the system and algorithm behavior at a given point in time, and assign them to the available workers. This approach has been used for graph databases and a number of graph processing systems [Salihoglu and Widom 2013, Nicoara et al 2015]. Dynamic repartitioning can be applied in-flight when, for instance, workers are waiting for a straggler worker to finish. In this situation, the vertices that have been assigned to the slowest worker can be repartitioned and reassigned to other idle workers to be computed faster and also balances the workload to reduce overall runtime. This does need support for dynamic migration by the graph framework [Salihoglu and Widom 2013]. Another reason to use dynamic repartitioning is when the number of active vertices in the graph change due to mutations or when the algorithm is non-stationary, causing vertices to become inactive, and it is suitable for iterative programming models such as Pregel [Xu et al 2014].

According to GPS [Salihoglu and Widom 2013], three major questions should be answered in a dynamic repartitioning process: 1) Which nodes should be reallocated?, 2) When and how to transfer the reallocated nodes to their new workers? and 3) How to place the reallocated nodes? These decisions can impose a heavy cost and affect the overall run-time. Some researchers have also shown that dynamic repartitioning does not offer significant performance improvements except under particular conditions. For example, the vertices in a PageRank algorithm are always active which makes dynamic repartitioning moot due to predictable and stationary load through the entire application's lifetime [Lu et al 2014]. But, despite these concerns, systems such as GPS, xDGP [Vaquero et al 2013], Mizan and XPreGel [Thien Bao and Suzumura 2013] have incorporate dynamic repartitioning and migrate the vertices synchronously across the workers, along with their incoming messages.

#### 4.1.2. Vertex-centric Partitioning vs. Edge-centric Partitioning

Vertex-centric frameworks partition the graph by assigning vertices to partitions and cutting some edges across partitions, in the process, while minimizing the number of such crossing edges. This is a common and well-supported partitioning approach. On the other hand, edge-centric frameworks partition the graph by cutting vertices and assigning edges to each partition. This approach minimizes the number of crossing vertices which is useful for many real-world graphs that have a power-law degree distribution to balance edges across the partitions well [Gonzalez et al 2012, Abou-Rjeili and Karypis 2006]. As was shown in Figure 10, vertices shared by edges belonging to different partitions would be cut and replicated across all the partition. One copy of the vertex is considered as the master and other copies are ghosts or mirrors. When updated, each ghost vertex sends its locally updated value to the master, and the master vertex applies updates from all ghost vertices to itself and sends the globally updated value back to the ghost vertices. We can see that many messages need to be passed across the network to maintain the cut vertices up to date.

To avoid this, PowerGraph does partitioning based on high-degree vertices of the graph and many systems have adopted such edge-centric partitioning [Kim and Candan 2012, Rahimian et al 2014, Xin et al 2013, Jain et al 2013]. There are also a number of additional optimizations that have been done [Rahimian et al 2014, Kim and Candan 2012]. Authors in [Feige et al 2008] and [Bourse et al 2014] have investigated several edge-centric (vertex-cut) approaches with vertex-centric (edge-cut) approaches and found that in many cases the former outperforms the latter.

#### 4.1.3. Pre-processing vs. Streaming

As seen in Section 2-1, there might be a pre-processing phase before the computation or the main processing starts. In the pre-processing approach, the large graph which is present on disk, will be partitioned before entering the system. Single-machine frameworks such as GraphChi [Kyrola et al 2012], TurboGraph [Han et al 2013], BPP [Najeebullah et al 2014] and CuSha [Khorasani et al 2014] use this method because they do not have enough memory to keep all the processing states in the single system. So, they partition the graph before starting the processing so help cope with large graphs. It also limits the partitioning operation, which can be costly, to a single time. Distributed frameworks like GoFFish do partitioning and subgraph identification in such a pre-processing phase for the same reason.

In streaming partitioning, the graph is partitioned once or as it is loaded into the graph processing system. The graph data enters the system sequentially, say a vertex and its adjacency list at a time, and the vertex is mapped to a partition on the fly. In this model, the order in which the vertices enter the system is important as each placement depends on the previous placements [Stanton and Kliot 2012]. Streaming can also benefit from a pre-processing model of partitioning, where specific vertex or edge ordering has been performed. Random partitioning, round-robin and range algorithms are three most common algorithms for steaming whereas linear deterministic greedy (LDG) [Stanton and Kliot 2012] and FENNEL [Tsourakakis et al 2014] are two greedy heuristics that improve the performance and quality of such online partitioning.

## 4.2 Communication Models

Graph processing systems use different approaches to communicate among their vertices, edges and partitions. In this Section we discuss these methods as shown in Figure 9 and enumerate the advantages and disadvantages of each method.

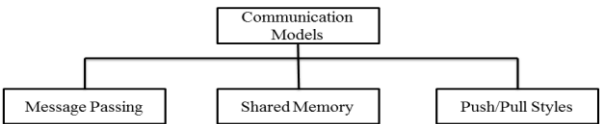


Fig. 9. Communication models in graph processing systems

4.2.1. Message Passing

Many distributed graph programming models offer explicit or implicit communication between their entities. In the message passing technique, communication is carried out by sending messages explicitly from one entity to another in the graph. The entity can be a vertex, edge or a component in a local or remote partition [Malewicz et al 2010]. Message passing is performed in many graph processing systems using communication libraries. For example, Pregel allows developers to pass messages from a vertex in the graph to another by calling an API. As part of the BSP execution model, messages sent in one iteration are received by the destination vertex in the subsequent iteration using bulk messaging. The receiving vertex updates its state based on incoming messages and sends its modified state to one or more of its neighbors by sending additional messages. Each source vertex maintains a list of its adjacent vertex IDs or outgoing edges IDs. Further, vertices also have queues where incoming messages from its neighbors and outgoing messages to its neighbors are stored between superstep boundaries.

A message passing model of communication is used by many graph processing frameworks and architectures, including vertex-centric, edge-centric and component-centric frameworks. Programming using synchronous message passing is also intuitive and the complexity is limited to an API to send a message to a destination entity, which is a familiar model for many programmers [Hudak 1989]. Although message passing is common in the frameworks with synchronous model of execution, it can be used in asynchronous execution models as well. Asynchronous message passing method is used extensively between vertices in the same partition or subgraph where they do not need to wait for other vertices to send their message in-bulk. Vertices and subgraphs can communicate asynchronously while communication between partitions can be synchronous. However, the asynchronous model of message passing brings more complexity to the programming paradigm because it requires more resources for storing and rebroadcasting data in a system where its components do not execute concurrently [Gehani 1990]. On the other hand, buffer management is an important issue that should be considered by message passing implementation. Issues like: How many buffers should each worker have? What should be the size of each buffer? When should a buffer block a sender? And what if the buffer is full but there are new messages coming? This model also has overheads because of the number of message replicas that exists in the network.

The Message Passing Interface (MPI) [El-Rewini and Abd-El-Barr 2005] is a common protocol used in graph processing systems, and is used by systems such as Pregel, GraphLab, Piccolo and Mizan. Portability and velocity are two significant advantages of using MPI where creating overhead is its most noticeable disadvantage. Communication can be done by passing actual messages between servers or by serialization. For Java based systems like Giraph and Hama, protocols such as Thrift [Apache Thrift] and ActiveMQ [Apache ActiveMQ] can be used for message passing. They utilize remote procedure call (RPC) to communicate seamlessly without the need to change the messages structures. Also protocols such as Avro [Apache Avro] and Protocol Buffer [Protocol Buffers] can be used for serialization by which the data will be serialized to be able to be sent between different platforms.

Systems also propose optimizations on top of these standard messaging libraries to reduce the communication overhead and minimize the runtime of the algorithm by reducing the number and size of messages that are passed. For example, GasCL has considered two different message buffering strategies: first, messages from the same source are stored together, second, messages for the same destination are stored together. The second approach is more common and when a message is dispatched from the origin, it is instantly put down to the right target node. It also uses a reverse edge index to store the message which utilizes array offset to facilitate message combining. Giraph++ [Tian et al 2013] introduces two types of messages in its hybrid model. “Internal messages” that are messages sent from a vertex within a partition to another vertex within the same partition, and “External messages” that are messages sent from the vertices of one partition to another partition. It provides two incoming message buffers for each vertex inside a partition as well: *inbox<sub>in</sub>* for internal messages and *inbox<sub>out</sub>* for external messages. In GPS [Salihoglu and Widom 2013], instead of sending multiple copies of the same message to multiple vertices in another partition, the system only sends one message to the remote partition and then, in the remote partition, the message will be copied to the vertices that need to receive it. This can dramatically reduce the network traffic.

4.2.2. Shared Memory

Using shared memory for communication is well suited for frameworks running on a single server, but can also be used for distributed systems in place of explicit message passing. In this model, the memory location can be simultaneously accessed by multiple processing modules, including both read and write to that location. In contrary to message passing, the shared memory method avoids extra memory copying and intermediate buffering. As single machine frameworks have limited memory and CPU resources, this shared-memory model that is often natively supported by the operating systems

35:18

S. Heidari et al.

is preferred [Nitzberg and Lo 1991]. Locks or semaphores are usually used in this model to prevent race conditions because concurrent tasks can read and write to the same memory [Low et al 2012, Chen et al 2008]. To maintain memory consistency, shared memory machines provide mechanisms for invoking the appropriate job (sequential consistency) or reordering a collection of jobs to be executed consecutively (relaxed consistency).

In distributed systems, it is also possible to have a distributed shared memory, where changes to memory locations are internally transfer using messages between different machines. From the programmer's points of view, they only perform memory accesses and the development is easier as explicit messages need not be passed. The concept of data "ownership" is lacking as well since anyone can write to that location. On the other hand, data locality cannot be controlled easily and if many remote workers access a particular memory location, it puts pressure on the processor and memory holding that location and can also lead to higher bus traffic and cache misses.

Virtual shared memory can be realized by using **ghost vertices or mirrors which are the copies of distant adjacent vertices** [Low et al 2010]. One machine keeps the main vertex and another machines work on copies of this vertex. Main vertex and ghost copies are shown in Figure 10. By keeping the mirror copies immutable during the computation with distributed write locking or an accumulator, the consistency is maintained [Low et al 2012, Power and Li 2010]. Both GraphLab and PowerGraph use this approach for communication. In particular, this is suitable for edge-centric frameworks because vertices should be cut in these frameworks and the partitioning is done based on edge-grouping. So, vertices can be easily cut (Figure 10). **Trinity [Shao et al 2013] is another memory-based graph processing system (Section 4-3)** that uses ghost vertices for communication. The Trinity specification language (TSL) maps the data storage and graph model together. The parallel boost graph library (PBGL), is a parallel graph processing library also uses ghost nodes but with message passing mechanism [Siek et al 2002].

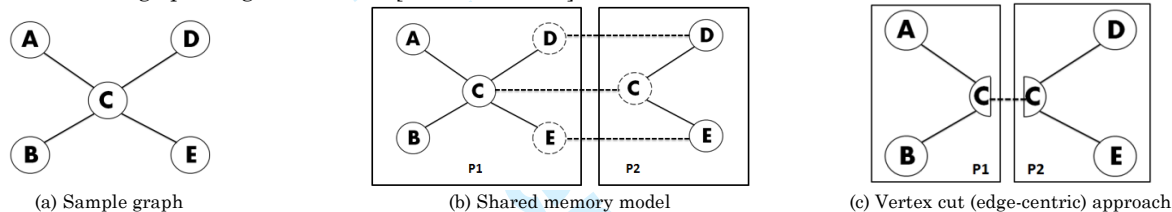


Fig. 10. Shared memory model and vertex cut (edge-centric) approach

#### 4.2.3. Push/Pull Styles

Pull/Push model is used with active messages. Active messages are those that carry both data and the operator that should be applied on them [Han et al 2014, Zhang et al 2012]. In the push style, the information flows (is pushed) from an active vertex to its adjacent vertices and in the pull style, the information flows (is pulled) from the neighbors of an active vertex to that active vertex. This kind of communication is using in the push/pull computation model that was discussed in Section 3.4.1. In terms of consistency, the pull style is naturally consistent because the active vertex would be updated in this phase, but the push style needs to use locks for every neighbor's update. Active messages are sent asynchronously in this model and they will be executed when they are received by the destination vertex. The sending and receiving messages are even combined in a framework such as GRE [Yan et al 2013] and it does not need to save intermediate states anymore. It should be noticed that some algorithms such as betweenness centrality cannot be expressed using pull style communication, but it is just a benefit for some algorithms such as PageRank to use it as an optimization [Gharaibeh et al 2013].

#### 4.3 In-Memory Execution

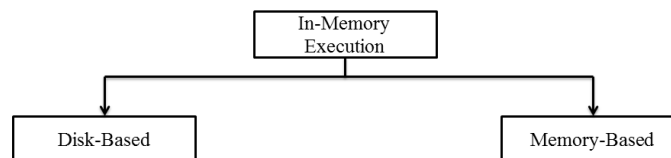


Fig. 11. In-Memory execution

##### 4.3.1. Disk-based

According to a storage view of execution models, two common approaches can be used: 1) **disk-based approach**, 2) **memory-based approach**. A disk-based execution fetches the graph data from physical disks, not just when loading the graph initially by also actively writes and reads parts of the graph state to/from disk during the execution. The advantages of using a disk-based approach is that it is cheaper to add disk capacity rather than memory, some large graphs do not fit in distributed memory either, and one can persist the partial state of execution in the middle of the processing to enable

recovery from faults [Chockler et al 2009]. Disk management is also easy, so many graph processing systems use this approach (Table 3).

On the other hand, the computation that is performed by graph algorithms is data-driven [Lumsdaine et al 2007], and they need many random data accesses and hard disks are still slow and inefficient compared to main memory. One of the challenging issues for disk-based graph processing is how to make disk access more efficient. **BPP (BiShard Parallel Processor)** [Najeebullah et al 2014], for example, provides a disk-based engine for processing large graphs on a single server. A novel storage structure called BiShard (BS) has been introduced which divides the graph into subgraphs containing equal number of edges and stores the in-edges and out-edges independently. This technique decreases the number of non-sequential I/O considerably and has two advantages compared to single shard storage mechanism that is presented by GraphChi. First, by storing in-edges and out-edges separately, access to each of them becomes independent and the system does not need to read the whole shard for every subset of vertices. Second, each edge has two copies in BS (one in each direction) which eliminates race condition among vertices to access their edges. Furthermore, BPP uses a novel asynchronous vertex-centric parallel processing model that leverages BS to provide full CPU parallelism for graph processing.

Other frameworks such as Giraph also support out-of-core execution using disk. When a graph is too big to fit into main memory (like small clusters) or a certain algorithm creates very large message sets (many messages or large ones) these frameworks can spill the excess messages or partitions to disk, later to be incrementally loaded and computed from disk. In addition, some frameworks such as FlashGraph [Zheng et al 2015] and PrefEdge [Nilakant et al 2014] use SSD instead of HDD to make data transmission and computation faster for out-of-core computation.

4.3.2. Memory-based

In the memory-based approach, the graph and its state are exclusively stored in memory during runtime for storing and processing the big data. For example, Giraph runs the whole computation in memory and reaches the disk for checkpointing and I/O; and Blogel keeps all neighbors of a typical high-degree vertex in the same block to be processed by in-memory algorithms and avoid message passing. The most important benefit of this approach is that using RAM or cache for processing is much faster than disk-based approach since the **CPU can access memory much quicker than disk** [Mittal and Vetter 2015]. However, memory is much more expensive than spinning disks and this becomes challenging when we consider larger graphs. So, memory-based systems must be efficient in retaining only relevant data in memory and in a compact form. Memory-based models also have less scalability than disk-based models, especially in a single machine system.

In GoFFish, the framework only loads a subset of properties for a given property graph from disk into distributed memory based on those attributes that are used by the algorithm, in addition to the complete topology of the graph that is always loaded [Jamadagni and Simmhan 2016]. This limits the memory footprint of the graph application during runtime. Many memory-based systems also use columnar representation since this offers a compact representation of data. Zhong and He [Zhong and He 2013] have indicated that GPU acceleration cannot reach considerable speedup if the data has to be loaded from disk because of the I/O costs that are themselves comparable to the total query runtime.

Microsoft's Trinity [Shao et al 2013] is a distributed graph processing engine over a memory cloud. It is supporting both online query processing which requires low latency (finding a path between users in a social network), and offline query processing which requires high throughput (PageRank). Trinity uses TSL (Trinity specification language) for communication that supports both synchronous and asynchronous modes. It stores objects as blobs of bytes that is economical, compact, with no serialization and deserialization burden. As a storage infrastructure, it structures the memory of numerous hosts to a distributed memory address space which is universally addressable, for maintaining huge graphs. Trinity has three main components including: 1) *slave* that stores graph data and computes on them, 2) *client* that acts as a user interface between Trinity and the user that communicates with Trinity proxies and slaves through the APIs provided by Trinity library, and 3) *proxy* that is an optional component for handling messages as a middle tier between clients and slaves. These components, along with other features like user-defined communication protocol, graph schema and computation models through TSL, enable Trinity to process the graph efficiently on memory cloud.

4.4 Fault Tolerance

Fault tolerance enables a system to continue performing properly even if some of its components face failures [Elena 2013]. Since graph processing systems are created from distributed and commodity components, it is possible that components confront failures which in turn will affect the execution and correctness of the applications. In order to improve the reliability and robustness of these systems, several techniques have been developed to support error handling and fault tolerance of the graph framework. Figure 12 shows the techniques that are used in many graph processing systems.



35:20

S. Heidari et al.

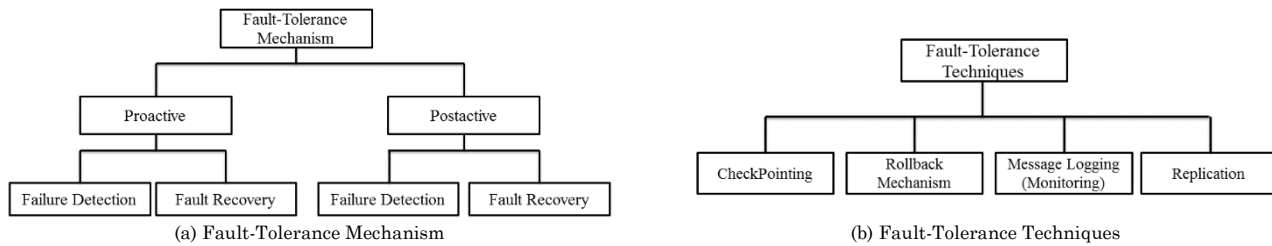


Fig. 12. Fault-tolerance in graph processing systems

Error handling in a graph processing system, as with other systems, has two main phases: 1) *failure detection* in which the system discover the error, and 2) *fault recovery* in which the system tries to resolve the problem and resume the operation. Several researches have been done on various fault-tolerance techniques on parallel and distributed systems [Kavila et al 2013, Treaster 2005, Elnozahy et al 2002]. In [Treaster 2005] for example, two types of components in an application, called central components and parallel components, are investigated where both mostly use rollback and replication methods for fault recovery. On the other hand, some graph processing systems do not support any error handling because it increases the complexity of the system, and the overheads can strongly affect the execution time.

Most graph processing systems use *checkpointing* and *rollback mechanism* [Egwutuoha et al 2013] for failure recovery, such as Pregel and Pregel-like systems like Giraph. Pregelix [Bu et al 2014], for instance, checkpoints states to HDFS at any superstep boundary that is selected by the user. The checkpointing applies to vertices and messages at the end of each superstep and ensures that the user does not need to know anything about the failure. Whenever a host or disk failure occurs, the unsuccessful machine will be added to a blacklist. For recovery, Pregelix reloads the state of the latest checkpoint to a set of “failure-free” workers that is periodically updated.

Piccolo [Power and Li 2010] uses a global checkpoint-restore method to recover from failures by providing synchronous and asynchronous checkpointing APIs. Synchronous checkpoints are suitable for iterative algorithms such as PageRank where the state in different iterations are decoupled by global barriers and it is adequate to checkpoint the state every few iterations. But asynchronous checkpoint is used to save the state of algorithms with long running algorithms, such as distributed crawler, periodically. Piccolo also utilizes Chandy-Lamport (CL) distributed snapshot algorithm [Chandy and Lamport 1985] for checkpointing. Once a failure is detected in a worker, the master will reset the status of all other machines and recover the operation from the latest finished universal checkpoint. The interior status of the master will not be checkpointed in Piccolo.

PowerGraph is another system that uses snapshots of the data-graph for fault-tolerance. The synchronous engine in PowerGraph creates the snapshot at the end of each superstep and before the start of next superstep while the asynchronous engine suspends the execution of the system to create the snapshot. Many of these systems provide task rescheduling after the recovery phase. Some systems such as Pregel, Piccolo and GraphLab benefit from rollback which allows them to continue the computation from the point that failure happened while in a number of systems, fault recovery is not completely provided and they need to restart the processing from the scratch [Han et al 2014, Krepska et al 2011, Plimpton and Devine 2011]. All these mechanisms are post-active fault tolerance approaches which means they handle the failure after it has happened.

Trinity [Shao et al 2013] uses message logging and replication for pro-active fault-tolerance, where the failure will be considered before scheduling and releasing a job for execution. Trinity utilizes heartbeat messages to proactively detect failures in machines. In addition, machines that unsuccessfully try to access the address-space in other machines also report the inaccessible machine to the master, and await for the addressing table to be updated before retrying the memory access. Meanwhile, in the recovery phase, the master reloads the data in the failed machine to another machine, updates the addressing Table and distributes it. Trinity provides checkpointing after every few iterations for synchronous BSP-based computations and provides “periodic interruption” mechanisms to generate snapshots in asynchronous computations. Buffered logging approach that is suggested in RAMCloud [Ousterhout et al 2010] has been used in Trinity to recover from failures in online queries while for read-only enquiries it only restarts the faulty machine and load the data again from the steady disk storage. GraphX [Xin et al 2013] uses lineage-based fault tolerance that assumes its RDDs cannot be updated but only created afresh. It has a very light overhead compared to the systems which use checkpointing as well as arbitrary dataset replication. So, it attains fault-tolerance without explicit checkpoint recovery while the retaining in-memory performance of Spark.

#### 4.5 Scheduling

Scheduling techniques help assign and manage jobs on the system resources [Pinedo 2012]. This is particularly useful in parallel and distributed multitasking systems in which several computations have to be done on a limited number of resources. In graph processing systems with large scale graphs having billions of vertices and edges, the vertex or edge

(depending on the programming model) will need to be scheduled for computation on a processing host. Typically, collections of vertices or edges are grouped into a coarser unit for scheduling, such as a partition or a subgraph, and it is the coarse unit that is actually scheduled on a CPU core. Within a processor, there may be multiple threads that execute individual vertices or edges in partition, leveraging, say, vertex level parallelism in a vertex-centric model.

According to [Doekemeijer and Varbanescu 2014], three different types of scheduling methods have been used for graph processing in general that is shown in Figure 13.

In batch scheduling method, the entire graph would be scheduled for processing across computing resources. This is more beneficial in bulk iteration model of computing such as Pregel. There is no priority or precedence in executing the partitions of the graph and they will be processed in any arbitrary order [Chen et al 2012]. This model has been used widely in dataflow frameworks such as Hadoop, Haloop [Bu et al 2010] and Twister [Ekanayake et al 2010]. There is always a preferred situation, like a limited number of iterations that is used as a condition for finishing the process.

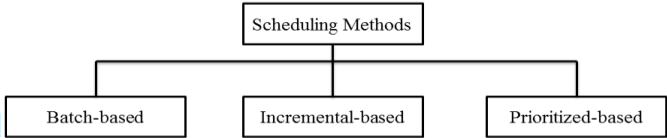


Fig. 13. Graph processing scheduling methods

Scheduling can be done once at the beginning of the application or redone at the start of each superstep. For e.g. Giraph allows partitions of the graph to be mapped to workers both at the start of the application and at each superstep, while TOTEM maps partitions to workers at the start of the application and retains that mapping. Remapping of partitions to workers as the application is executing also requires the ability to migrate both the graph and its updated state and messages to different workers. The mapping of vertices and edges to partitions may also change as a result. For example, Mizan migrates the vertices from busy workers to the one with fewer vertices to load the balance, and GPS repartitions the graph to distribute the load among idle workers during the computation.

Another aspect is on whether the partitions are mapped to a static set of compute resources or the resources themselves can be elastic over the execution of the application. For example, [135] looks at mapping partitions to an elastic set of VMs based on the expected computational complexity of the partition for stationary and non-stationary graph algorithms.

In contrast, in prioritized scheduling method, jobs will be processed according to a priority condition that is defined by the user. System such as Maiter [Zhang et al 2012] shows that using this method results in quicker convergence for many graph algorithms. For instance, a defined prioritizing function can schedule jobs based on the number of vertices in each partition. In GoFFish, the largest subgraphs in a partition are executed first so that the computing of smaller subgraphs can be interleaved with the message passing from the large subgraph. Prioritized scheduling can be helpful in processing imbalanced workloads.

Doekemeijer and Varbanescu [Doekemeijer and Varbanescu 2014] believe that incremental scheduling only processes a subdivision of data like active vertices. This model is used in a number of graph processing systems, e.g., Stratospher [Alexandrov et al 2014] and GraphLab [Low et al 2012], in which the processing continues until there are active vertices.

5. GRAPH DATABASES

A graph database is one where the data is natively stored as a graph structure that can be queried upon [Angles and Gutierrez 2008]. The data itself is typically a property graph with not just vertices and edges but also name-value properties or labels defined on vertices and edges. Graph query models support different types of traversal queries such as path, reachability and closure, in addition to filter queries over their properties [Jamadagni and Simmhan 2016, Sarwat et al 2013]. Graph databases contrast with relational databases that store graphs – the latter require multiple joins for traversal of graphs rather than having direct references from a vertex to its neighbors that allows for faster query processing in graph databases. Graph databases leverage the topological properties of graphs, including graph theory and query cost models, in answering the queries. The queries also provide a high-level declarative interface for processing and accessing the graph compared to a graph framework that requires users to write a program using their graph programming abstractions and executes it in batch [Vicknair et al 2010]. Another distinction from the graph frameworks discussed above is the need for low latency ( $O(seconds)$ ) execution of hundreds of queries rather than high throughput analysis of single programs over large graphs.

Graph datasets are getting increasing attention every day and several companies are starting to utilize graph databases to perform interactive queries to support their business needs. According to DB-Engines, which is an industry observer, “graph DBMSs are gaining popularity faster than any other database categories”, that shows remarkable growth in the last few years (Figure 14).

35:22

S. Heidari et al.

Neo4j [Neo4j 2015] is a popular graph database designed as an open-source NoSQL database. It supports ACID (Atomicity, Consistency, Isolation, Durability) properties by implementing a Property Graph Model efficiently down to the storage level. It also allows for traversal of relationships across both breadth and depth of the graph in constant time because it represents vertices and edges efficiently. It is useful for single server deployments to query over medium sized graphs due to using memory caching and compact storage for the graph. Its implementation in Java also makes it widely usable. Besides the single server model, it also provides master-worker clustering with cache sharding for enterprise deployment. However, according to some reports, the scalability of the distributed version is not as good as even relational databases and it has deadlocks problems such as not being able to handle two concurrent *upserts* if they touch the same node<sup>2</sup>.

OrientDB and Titan are two other well-used graph databases [OrientDB vs Neo4j 2015]. OrientDB can save 220,000 records per second on ordinary hardware. It supports multi-master replication and sharing which give it better scalability. It also provides a security profiling system based on roles and users in the database. Titan [TITAN Distributed Graph Database 2015] is another open-source distributed transactional graph database that provides linear elasticity and scalability for growing data, data distribution and replication for fault-tolerance and performance. It supports ACID and different storage back-ends such as Apache Hbase [Apache Hbase 2015] and Apache Cassandra [Apache Cassandra 2015]. Titan also uses the Gremlin query language [Gremlin Query Language] in which traversal operators are chained together to form path-oriented expressions to retrieve data from the graph and modify them.

Twitter has developed its own graph database called FlockDB [Twitter 2012] to store social graphs such as “who blocks whom” and “who follows whom”. FlockDB is an open-source fault-tolerant distributed graph database which aims to support online data migration, add/delete/update operations, complicated set of arithmetic queries, replication, archive/restore edges and so on. In April 2010, the FlockDB cluster had stored more than 13 billion connections (edges) and supported a peak traffic of 20K writes per second with 100K reads per second [Green 2013]. But it appears that FlockDB is not able to traverse graphs deeply and is not implementing the full stack of storage services<sup>3</sup>.

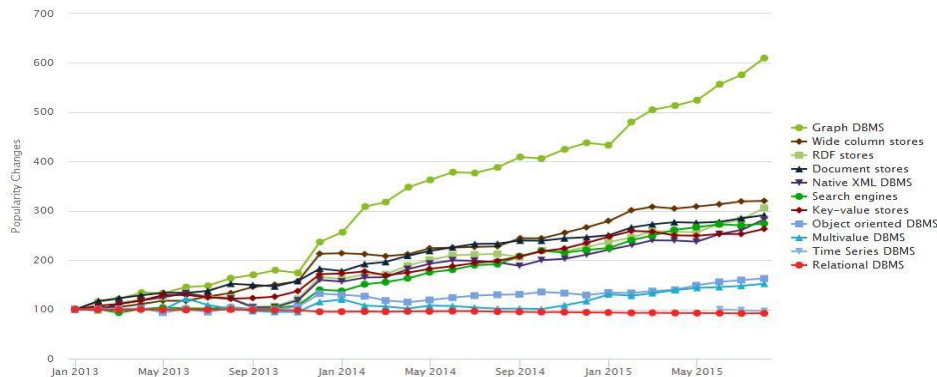


Fig. 14. Popularity changes in using databases [201]

There is also research on distributed graph databases, though this is an emerging area. Horton+ [Sarwat et al 2013] from Microsoft offers a graph query language that supports path, closure and joint queries over property graphs. It converts the query into a Deterministic Finite Automaton (DFA) that is executed over a distributed database using a vertex-centric BSP model based on Giraph. GoDB [Jamadagni and Simmhan 2016] is another research database that leverages GoFFish to offer similar query capabilities over property graphs, but with support for scalable indexes and using a subgraph-centric model of execution that offers a much faster performance relative to Titan and Horton+. GBASE [Kang et al 2011] introduces *compressed block encoding* graph storage method that utilizes adjacency matrix representation to store homogeneous regions of graphs. It also uses a grid-based selection strategy for query optimization to provide quicker answers by minimizing disk accesses. Quegel [Yan et al 2016] handles enquiries as “first-class citizens” by which the user is only required to determine the Pregel-like algorithm for a general enquiry. Then, it sets up the computing and processing of multiple inbound enquiries on demand.

There are several other open source and commercial graph databases such as HyperGraphDb [HyperGraphDb 2015], AllegroGraph [AllegroGraph 2015], InfiniteGraph [InfiniteGraph 2015], InfoGrid [InfoGrid 2015], JCoreDB Graph [JCoreDB 2015], ArangoDB [ArangoDB 2015], GraphBase [GraphBase 2015], MapGraph [MapGraph 2015, Fu et al 2014], and Weaver [Weaver 2015]. All these projects try to provide modern solutions for storing and retrieving large-scale graph

<sup>2</sup> <https://news.ycombinator.com/item?id=9699102>

<sup>3</sup> <http://stackoverflow.com/questions/2629692/how-does-flockdb-compare-with-neo4j>

data and it seems that this area is a very promising field of research and commercial investments for the future. Jouili and Vansteenbergh [2013] have compared some of these graph databases.

6. SYSTEM CLASSIFICATION

Table 3 presents the key graph processing systems with their characteristics according to the proposed taxonomy. The notations in the table for each category are as follow:

- Programming model: vertex-centric (V), edge-centric (E), component-centric (C), path-centric (P) or data-centric (Da).
- Architecture: distributed (D), single machine (S) or heterogeneous (H).
- Computational Model: different names are used by different systems
- Communication Model: message passing (MP), shared memory (SM) or dataflow (DF)
- Coordination: synchronous (Synch), asynchronous (Asynch) or both timing approach together
- Storage: disk-based (DB) or memory-based (MB) storage approach.
- N/A means that there is no specific name or method mentioned by the paper that is describing the system.

Table III. Overview of existing graph processing frameworks

Year	System	Programming Model	Architecture	Computational Model	Communication Model	Timing	Storage
2009	PEGASUS [Kang et al 2009]	N/A	D	N/A	DF	Synch	DB
2010	Pregel [Malewics et al 2010]	V	D	BSP	MP	Synch	DB
2010	Signal/Collect [Stutz et al 2010]	V	S	Signal/Collect	MP	Both	DB
2010	Surfur [Chen et al 2010]	V	D	Transfer-combine	MP	Synch	DB
2010	JPregel [Prakasam and Chandrasekhar 2010]	V	D	BSP	MP	Synch	DB
2010	GraphLab [Low et al 2010]	V	S	N/A	SM	Asynch	DB
2010	Piccolo [Power and Li 2010]	Da	D	Three phases	Dataflow	Synch	DB
2011	GoldenOrb [Cao 2011]	V	D	BSP	SM	Synch	DB
2011	GBase [Kang et al 2011]	E	D	N/A	Dataflow	Synch	DB
2011	HipG [Krepska et al 2011]	V	D	BSP	SM	Both	DB
2012	Giraph [Apache Giraph]	V	D	BSP	MP	Synch	DB
2012	Distributed GraphLab [Low et al 2012]	V	D	GAS	SM	Both	DB
2012	KineoGraph [Cheng et al 2012]	V	D	Push/Pull	MP	Synch	MB
2012	PowerGraph [Gonzalez et al 2012]	E	D	GAS	SM	Both	DB
2012	Sedge [Yang et al 2012]	V	D	BSP	MP	Synch	DB
2012	GraphChi [Kyrola et al 2012]	V	S	PSW	SM	Asynch	DB
2013	TOTEM [Gharaibeh et al 2013]	V	H	BSP	Both MP and SM	Asynch	MB
2013	Mizan [Kayyat et al 2013]	V	D	BSP	MP	Synch	DB
2013	Trinity [Shao et al 2013]	V	D	TSL	SM	Asynch	MB
2013	Grace [Wang et al 2013]	V	S	Three phases	MP	Asynch	DB
2013	GPS [Salihoglu and Widom 2013]	V	D	BSP	MP	Synch	DB
2013	Giraph++ [Tian et al 2013]	C	D	BSP	Both MP and SM	Both	DB
2013	Naiad [Murray et al 2013]	V	D	Timely Dataflow	SM	Both	MB
2013	PAGE [Shao et al 2013]	V	D	Partition-aware	MP	Synch	DB
2013	Stratospher [Ewan et al 2013]	V	D	Push/Pull	Dataflow	Synch	DB
2013	TurboGraph [Han et al 2013]	V	S	Pin-and-slide	SM	Asynch	DB
2013	xDGP [Vaquero et al 2013]	V	D	BSP	MP	Synch	DB
2013	X-Stream [Roy et al 2013]	E	S	Scatter-gather	MP	Synch	DB
2013	GiraphX [Tasci and Demirbas 2013]	V	D	BSP	SM	Asynch	DB
2013	GraphX [Xin et al 2013]	E	D	GAS	Dataflow	Synch	MB
2013	Galois [Nguyen et al 2013]	V	S	ADP	SM	Asynch	DB
2013	GRE [Yan et al 2013]	V	D	Scatter-Combine	MP	Synch	DB
2013	Ligra [Shun and Blelloch 2013]	C	S	Push-pull	SM	Asynch	MB
2013	LFGraph [Hoque and Gupta 2013]	V	D	N/A	SM	Synch	MB
2013	PowerSwitch [Xie et al 2015]	V	D	Hybrid	SM	Both	DB
2013	Presto [Venkataraman et al 2013]	V	D	N/A	Dataflow	Synch	DB
2013	Medusa [Zhong and He 2013]	V	H	EMV	MP	Synch	MB
2014	RASP [Yoneki et al 2014]	V	S	Scatter-gather	SM	Asynch	DB
2014	GoFish [Simmhan et al 2014]	C	D	Iterative BSP	Both MP and SM	Synch	MB
2014	GasCL [Che 2014]	V	H	GAS	MP	Synch	MB
2014	CuSHa [Khorasani et al 2014]	V	H	GAS	SM	Asynch	MB
2014	BPP [Najeebullah et al 2014]	V	S	BSP	SM	Asynch	DB
2014	Imitator [Wang et al 2014]	V	D	BSP	MP	Synch	DB
2014	GraphHP [Chen et al 2014]	V	D	BSP	MP	Synch	DB
2014	PathGraph [Yuan et al 2014]	P	S	Scatter-gather	SM	Asynch	DB
2014	Seraph [Xue et al 2014]	V	D	GES	MP	Synch	DB
2014	GraphGen [Nurvitadhi et al 2014]	V	H	N/A	SM	Synch	MB
2015	Pregelx [Bu et al 2014]	V	D	Join-operator based	MP	Synch	DB



35:24

S. Heidari et al.

2015	FlashGraph [Zheng et al 2015]	V	S	BSP	Both MP and SM	Asynch	DB
2015	GraSP [Battaglini et al 2015]	V	D	N/A	MP	Synch	MB
2016	iGiraph [Heidari et al 2016]	V	D	BSP	MP	Synch	DB

Although, many frameworks have been proposed for processing large-scale graphs, there are still several gaps that need to be addressed, as highlighted by this table. Among these observations are: 1) Many graph processing systems have been developed based on vertex-centric programming model because it is the simplest way of partitioning and processing large-scale graphs. Although edge-centric and component-centric systems are more difficult to be implemented, it has been empirically shown that frameworks such as PowerGraph and GoFFish can scale more efficiently than vertex-centric ones. So, those types of systems need to be investigated more. 2) Disk-based approach is the dominant mechanism that is using by most frameworks. It also includes the frameworks that support out-of-core computation. Disks are cheap but much slower than memory. On the other hand, memory is faster than disk but it is more expensive and memory management makes it more complicated to develop a system based on this approach. 3) Synchronous programming is popular on distributed systems as they avoid race conditions, but often require message passing and have longer runtimes due to the coordination. While asynchronous methods work well on single machine and heterogeneous based frameworks, its effect on distributed frameworks are necessary for scaling is less studied.

## 7. FUTURE DIRECTIONS

Although several works have looked into improving graph processing systems (see Table 3), there are still a number of issues that are open. For example not many graph processing frameworks use dynamic repartitioning which performs better than simple static partitioning in many cases. Most the frameworks use checkpointing for error handling, which can be costly, and other approaches to fault recovery are not well studied. While many researchers have studied classic graph algorithms such as PageRank and shortest paths, it is not clear whether these frameworks can still perform as well for more sophisticated and real-world applications such as machine learning algorithms.

Besides these, there are several other advances to the programming and data model of large graphs, and the runtime execution of the graph platforms that need to be examined. These are discussed below.

### 7.1 Incremental Processing Models

Regardless of the type of framework or algorithms used for processing big graphs, or how large the graph is, data can be processed in three different ways as shown in Figure 15.

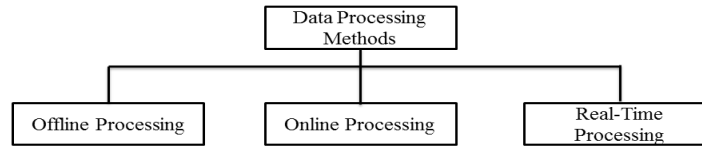


Fig. 15. Data processing approaches

According to the problem domain that a framework wants to present solutions for, each data processing approach shown above can be considered in the framework. Offline processing (batch processing) is done where a number of analogous tasks are gathered together to be processed by a computing system all at once instead of individually. In this method, which is used in many graph processing frameworks, the whole graph dataset is loaded into the system, processed for an application, and the results will be return to the user. The original graph is not changed externally, other than through modifications by the running application, and this leads to predictable partitioning and scheduling strategies which make their design easy.

In online processing, user can communicate with the system and can make changes to the graph data stored in the system. Thus, the system will be updated automatically and re-process the data with new values periodically or based on user-defined events, which is not necessarily real-time and immediate.

Real-time processing allows the graph to change over time based on incremental updates that it receives to the graph topology or properties. Processing such dynamic graphs is more like an event-driven system where sensors may generate a stream of updates about a vertex or edge that the sensor represents (e.g., road network with traffic cameras or sensors). Real-time processing requires that the computation should be done immediately after the changes happen to the data, and the updated results should be returned with very short delays. Sometimes, the operations may be performed on the delta events themselves before they are actually applied to the graph data. Such requirements are increasingly important in competitive businesses such as social networks, and in IoT domains. There is limited work in the area of real-time processing of large dynamic graphs. For example, Twitter uses a graph-based content recommendation engine called GraphJet [Sharma et al 2016] which is an in-memory framework that supplements real-time with batch processing by keeping real-time bipartite interplay graph between tweets and users.

The temporal dimension can also come through the notion of time-series graphs where different states of a graph are available, and the application has to operate over both the spatial and temporal dimension. However, this data base be collected *a priori* and available offline, and distinct from the changing states of the graph arriving in real-time. For e.g., GoFFish operates over time-series graphs for algorithms such as time-dependent shortest path and tracking meme propagation [Simmhan et al 2015].

7.2 Complex Workflows

A workflow is a dependency graph of different tasks that should be done in a specific order to complete a bigger job. Current graph processing frameworks are based on very simple workflows, typically singleton workflows with one operation executed in a data-parallel manner. They pick a dataset and an algorithm and execute the algorithm on the data. They usually try to solve very simple problems such as finding shortest path or PageRank problem. But, many real world problems are not as easy as this. For example, in a social network a typical scenario can be like this: an algorithm finds all friends and followers of somebody, then finds the common interests between them using another algorithm, draws a map of his/her communication history, combines all these information with the information from other people in that city to find the whole trends and so on. Such a complicated series of processes cannot be modeled seamlessly based on existing graph processing systems without manually creating multiple jobs and passing data explicitly between them through the file system. While there is limited work on Master-Compute model that allows a master task to change the phase of computation on the workers, this can be used to model only simple sequential operations on a single graph rather than more complex operations that may even span different graphs. So, new frameworks are needed to allow the users to perform more complicated operations.

On the other hand, such complex scenarios also require efficient resource provisioning. That is, proper scheduling is critical to minimize the monetary costs and execution time on one side, and improve resource utilization and performance of the system on the other side. Graph tasks in the workflow may have different processing needs, and may arrive at different intervals, and with different priorities and profitability metrics. Managing these graph workloads offer novel challenges as well. Some research issues on this include:

- How to schedule complex graph workflows to gain minimum cost and maximum resource utilization?
- What factors influence workflow management in graph processing systems considering graph algorithms characteristics and features of graph datasets?
- How does workflow management in graph processing frameworks—specially for complex scenarios- affect the energy consumption of resources?

7.3 Graph Databases

Relational databases have existed since the 1980s, and have grown mature. While they deal well with structured data tuples stored in tables, their use for storing and querying graph datasets is limited. As discussed in Section 5, graph database, while not a novel concept, are still in their early stages when considering large property and semantic graphs. It is because relationships in a graph database are much stronger than those hypostatized at runtime in a relational database since they are being treated as high priority entities [Robinson et al 2015]. Relational databases are much slower than graph databases for connected data, hence using graph databases is recommended for highly connected environment and applications such as social networks, IoT, business transactions, and Web searching.

Despite the usefulness of graph databases in the aforementioned environments, they are not as mature as relational databases particularly in terms of tools for data mining purposes on massive graph data on distributed systems. Therefore, future directions for research include:

- What are the canonical query models for static and dynamic graphs? What is the equivalent of a relational algebra for graph databases?
- Supporting graph queries in combination with graph kernel algorithms, e.g., find all websites hosted in Australia (property) whose PageRank (algorithm output) is greater than X.
- What are the cost models to be developed for efficient execution of graph queries on distributed environments?
- Improve the ability to sustain low-latency processing of large numbers of transactional graph queries on distributed and elastic systems like Cloud.
- How can analysis be performed across data stored in traditional relational databases and graph databases seamlessly and effectively?
- How do we manage distributed data and indexes in graph databases that have data constantly changing or streaming in?

35:26

S. Heidari et al.

#### 7.4 Cloud Features and Cost Models

The Cloud computing paradigm has modified hardware, software and datacenters implementation and design. It offers new economical and technological solutions such as utilizing distributed computing, pay-as-you-go pricing models and resource elasticity. Cloud computing offers computing as a utility in which users can have accesses to different services according to their needs without heed to where the services are hosted or how they are delivered.

Computing as a service is the infrastructure service most relevant to graph processing. While the scalability offered by VMs has been used for graph processing, these are treated as yet another distributed resource by graph processing frameworks rather than consider their ability to elastically scale, or consider their costs. There is limited work on this regard. iGiraph has started to consider cost **optimization on clouds [Heidari et al 2016]. It classifies** the graph algorithms into convergent and non-convergent algorithms and utilizes a dynamic repartitioning algorithm by which reduces the number of VMs for the graphs that are shrinking during the operation to decline the price. It also performs better on non-convergent applications compared to the famous Giraph.

[Dindokar et al 2016] has proposed an approach to model the computational behavior of non-stationary graph algorithms using a meta-graph model for subgraph-centric programming model. The meta-graph model is able to offer predictions on the subgraphs that will be active in different supersteps, and this is used to schedule subgraphs to VMs in different supersteps, **including elastically scaling the VMs in and out [Dindokar and Simmhan 2016].** Their strategies show a pricing reduction by half for large graphs like Orkut for costly graph algorithms like Betweenness centrality, with minimal increase in the runtime relative to static over-provisioning of VMs. **Elasticity has also be examined in [Pundir et al 2016] where it uses two partitioning mechanisms** called Contiguous Vertex Repartitioning (CVR) and Ring-based Vertex Repartitioning (RVR) to 1) scale in/out without interfering graph computation, 2) decrease the network overhead after scaling, 3) keep the load balanced by reducing stragglers across servers.

Cloud providers have different cost models for their VM resources, typically, on-demand VMs that you pay for based on the minutes or hours used, and spot VMs that have dynamic pricing based on demand-supply, and are pre-emptible when the demand out-strips supply. Spot VMs are much cheaper than on-demand VMs and their use can also be explored for large graph applications, while addressing the faults that can occur due to out of bid event when prices spike. While this has been examined for applications like MapReduce [Chohan et al 2010], there is no work in this regard yet for graph processing.

Service level agreement (SLA) [Patel et al 2009] is a contract between a service provider and a service user to define the service features, the time for delivering the service, the steps that should be taken in the case of service crashes, service domain, prices, etc. Using SLAs, both user and provider can ensure that the service is delivered exactly based on what had been agreed upon and penalties can be applied in case of commitment violation. Quality of a service (QoS) [Ardagna et al 2014] provides a level of performance, availability and reliability offered by software, platform and infrastructure that the service is hosted on them. If we consider graph processing as a service then the quality of this service should be in an acceptable level from both provider and customer points of views. According to the aforementioned SLA and QoS definitions, and taking graph processing characteristics into consideration, some research directions can be defined as follow:

- Which parameters have the most impact on the performance of a graph processing service and quality of that?
- What factors should be considered for selecting appropriate graph processing service among other analogous services?
- How SLA-based resource provisioning and scheduling mechanisms for managing graph processing systems and services can be?

#### 7.5 Network Optimizations

The network communication and messaging aspects are less studied **in current graph processing frameworks.** The factors such as **network latency, network bandwidth, network traffic and topology can affect the runtime performance of the system.** The problem also becomes more complicated when it comes to the Cloud environments. Most existing distributed graph frameworks have been developed for integrated clusters in which resource management and communication is more predictable. But in a Cloud-based framework, the network performance can be variable and VM placement not in the control of users. Hence it becomes essential to consider **network factors. Unlike earlier works that considered the role of the network as trivial in graph processing [Ousterhout et al 2015],** particularly for the graphs that can fit into the memory of a single machine, most recent experiments showed that the network plays a major role in the performance of a graph processing system whether the graph can fit in the memory of a single machine or it is processed on a distributed system **[McSherry and Schwarzkopf 2015].** For example, allocating larger or denser partitions to the machines with higher bandwidth on one side and reducing the network traffic by decreasing the number of messages transferred between machines on the other hand can enhance the efficiency of the system [Chen et al 2010].

7.6 Other Improvements

Since scalable graph processing is still in its infancy, there are many open issues to improve the performance and features of each component discussed in Section 2-1, and the overall performance of the system. For example, read and write from/to disk is costly in these systems and usually acts as a bottleneck. In many researches such as [Zheng et al 2015] and [Nilakant et al 2014] SSD (solid state drive) is used as a faster storage device compared to traditional HDDs (hard disk drive) and cheaper compared to main memory. Further, efficient storage models for graph datasets on disks also need to be explored. For e.g., when processing large graphs, the time to load data from disk to memory can outstrip the time to perform the analysis. Compact and compressed graph data representation on disk, loading necessary subsets of the graph on-demand, and support for efficient storage of property graph are some novel topics to explore. Literature has also examined processing of large graphs on single machines and tries to keep the whole graph and computation results in memory (Section 3.1.2). They rely on memory costs dropping and capacities increasing with new technologies like 3D stacked RAM, when single machines will become viable even for billion-vertex graphs. So there are several aspects of storage and memory management that can be explored.

In addition to these, other parts of graph processing system pipeline can be improved as well. These include:

- What initial partitioning or pre-processing techniques can improve the performance of the system and speed up the computation process? How can repartitioning improve the efficiency of the system and if it can speed up the computation process?
- How can we better model and predict the behavior of different graph algorithms for graphs with different characteristics, such as power law, small world, planar, etc.? How are these affected by the different programming models? Can we use these to determine the ideal graph processing technique or strategy to be chosen, e.g., synchronous vs. asynchronous algorithms, computation-bound algorithms vs. memory-bound algorithms, denser datasets vs. less dense datasets and so on.
- Are there any computational mechanisms that uses less memory size or can reduce network traffic by reducing the number of messages between machines?
- What fault-tolerance techniques can be used other than check-pointing to improve system reliability and performance?
- What resource provisioning and scheduling algorithms can be used to optimize the processing framework particularly in a competitive environment such as cloud spot-markets?

8. SUMMARY AND CONCLUSIONS

Huge quantities of data is being created, analyzed and used every day in the contemporary world of Internet communications and connected devices. “Big Data” is the term used to signify the challenges posed by this massive data influx. A growing majority of big data is in the form of “Graphs” which is one of the major computational methods of huge data analysis. Social network applications and web searches, Internet of things, knowledge graphs and deep learning, financial transactions, and neuroscience are some examples of large-scale graphs that need to be analyzed for various domains. Several works have investigated on creation of effective systems for processing large-scale graphs in recent years.

In this paper, we have investigated and categorized existing graph processing frameworks and systems from different perspectives. First, we explained how different parts of a graph processing system including read and write from/to disk or memory, pre-processing, partitioning, communication, computation and error handling work together to process large-scale graphs. Second, we presented a taxonomy of different abstractions and approaches that are used in existing graph processing systems within each of these phases. In addition, we described notable frameworks that have used these techniques, and analyzed their advantages and disadvantages to support our discussions. We further summarized the features of graph processing frameworks developed since 2009 in Table 3. It gives a comprehensive overview of current systems and enables making comparison between them. Finally, future research directions are discussed which shows that scalable graph processing is still at a nascent stage and there are many issues that remain unsolved.

REFERENCES

Hendrickson, B., and Berry, J. W. (2008). Graph Analysis with High-Performance Computing. *Computing in Science and Engineering*, 14-19.

Khorasani, F., Vora, K., Gupta, R., and Bhuyan, L. N. (2014). CuSha: vertex-centric graph processing on GPUs. *23rd international symposium on High-performance parallel and distributed computing (HPDC '14)*, (pp. 239-252). Vancouver, BC, Canada.

McCune, R. R., Weninger, T., and Madey, G. (2015). Thinking Like a Vertex: A Survey of Vertex-Centric Frameworks for Large-Scale Distributed Graph Processing. *ACM Computing Surveys (CSUR)*, 48(2), 1-39.

Pinedo, M. L. (2012). *Scheduling: Theory, Algorithms, and Systems* (4th Edition ed.). Springer-Verlag New York: New York, US.

Abou-Rjeili, A., and Karypis, G. (2006). Multilevel Algorithms for Partitioning Power-Law Graphs. *IEEE International Parallel and Distributed Processing Symposium (IPDPS'06)*, (pp. 124-124). Rhodes Island, Greece.



35:28

S. Heidari et al.

- Afrati, F. N., Das Sarma, A., Salihoglu, S., and Ullman, J. D. (2012). Vision Paper: Towards an Understanding of the Limits of Map-Reduce Computation. *Cloud Futures 2012 Workshop*. Berkeley, California, USA: Microsoft Research.
- Akbari, F., Tajfar, A., and Farhoodi Nejad, A. (2013). Graph-Based Friend Recommendation in Social Networks Using Artificial Bee Colony. *IEEE 11th International Conference on Dependable, Autonomic and Secure Computing (DASC)*, (pp. 464-468). Chengdu, China.
- Alexandrov, A., Bergmann, R., Ewen, S., Freytag, J.-C., Hueske, F., Heise, A., . . . Warneke, D. (2014). The Stratosphere platform for big data analytics. *The VLDB Journal — The International Journal on Very Large Data Bases*, 23(6), 939-964.
- AllegroGraph . (2015). Retrieved August 10, 2015, from <http://allegrograph.com/>
- Ammann, P., Wijesekera, D., and Kaushik, S. (2002). Scalable, graph-based network vulnerability analysis. *9th ACM conference on Computer and communications security (CCS '02)*, (pp. 217-224). Washington, DC, USA.
- Andreev, K., and Racke, H. (2004). Balanced Graph Partitioning. *16th annual ACM symposium on Parallelism in algorithms and architectures (SPAA'04)*, (pp. 120-124). Barcelona, Spain.
- Angles, R., and Gutierrez, C. (2008). Survey of Graph Database Models. *ACM Computing Surveys (CSUR)*, 40(1), 1-39.
- Apache ActiveMQ. (n.d.). (Apache) Retrieved 05 28, 2016, from <http://activemq.apache.org/>
- Apache Avro. (n.d.). (Apache) Retrieved 05 28, 2016, from <https://avro.apache.org/>
- Apache Cassandra . (2015). Retrieved August 10, 2015, from <http://cassandra.apache.org/>
- Apache Giraph. (n.d.). Retrieved from <http://giraph.apache.org/>
- Apache Hbase . (2015, July 23). Retrieved August 10, 2015, from <http://hbase.apache.org/>
- Apache Hive TM. (n.d.). Retrieved from <https://hive.apache.org/>
- Apache Pig. (n.d.). Retrieved from <https://pig.apache.org/>
- Apache Spark. (n.d.). (Apache) Retrieved 03 23, 2016, from <http://spark.apache.org/>
- Apache Thrift. (n.d.). (Apache) Retrieved 05 28, 2016, from <https://thrift.apache.org/>
- Apache™ Hadoop. (n.d.). Retrieved July 23, 2015, from <https://hadoop.apache.org/>
- ArangoDB . (2015). Retrieved August 10, 2015, from <https://www.arangodb.com/>
- Ardagna, D., Casale, G., Ciavotta, M., Pérez, J. F., and Wang, W. (2014). Quality-of-Service in Cloud Computing: Modeling Techniques and Their Applications. *Journal of Internet Services and Applications*, 5(11), 1-17.
- Bader, D. A., and Madduri, K. (2008). SNAP, Small-world Network Analysis and Partitioning: An open-source parallel graph framework for the exploration of large-scale networks. *Proceedings of IEEE International Symposium on Parallel and Distributed Processing Systems (IPDPS'08)*, (pp. 1-12). Miami, FL, USA.
- Bader, D. A., Meyerhenke, H., Sanders, P., and Wagner, D. (2013). *Graph Partitioning and Graph Clustering*. Atlanta, GA, US: American Mathematical Society.
- Bagci, H., and Karagoz, P. (2015). Context-aware location recommendation by using a random walk-based approach. *Knowledge and Information Systems*.
- Bannister, M. J., and Eppstein, D. (2012). Randomized Speedup of the Bellman-Ford Algorithm. *Analytic Algorithmics and Combinatorics (ANALCO12)*, (pp. 41-47). Kyoto, Japan.
- Barker, B. (2015, January 14). Workshop: High Performance Computing on Stampede.
- Battagliolo, C., Pienta, R., and Vuduc, R. (2015). GraSP: Distributed Streaming Graph Partitioning. *Proceedings of 1st High Performance Graph Mining workshop (HPGM '15)*. Sydney, Australia.
- Beseri Sevim, T., Kutucu, H., and Ersen Berberler, M. (2012). New Mathematical Model for Finding Minimum Vertex Cut Set. *International Conference on Problems of Cybernetics and Informatics (PCI'2012)*, (pp. 143-144). BAKU, AZERBAIJAN.
- Borkar, V., Carey, M., Grover, R., Onose, N., and Vernica, R. (2011). Hyracks: A Flexible and Extensible Foundation for Data-Intensive Computing. *IEEE 27th International Conference on Data Engineering (ICDE '11)*, (pp. 1151-1162). Hannover, Germany.
- Bourse, F., Lelarge, M., and Vojnovi, M. (2014). Balanced Graph Edge Partition. *20th ACM SIGKDD international conference on Knowledge discovery and data mining (KDD '14)*, (pp. 1456-1465). New York, NY, US.
- Bu, Y., Borkar, V., Jia, J., Carey, M. J., and Condie, T. (2014). Pregelix: Big(ger) Graph Analytics on A Dataflow Engine. *VLDB Endowment*, 8(2), 161-172.
- Bu, Y., Howe, B., Balazinska, M., and Ernst, M. D. (2010). HaLoop: efficient iterative data processing on large clusters. *The VLDB Endowment*, 3(1-2), 285-296.
- Buluç, A., Meyerhenke, H., Safo, I., Sanders, P., and Schulz, C. (2013). *Recent Advances in Graph Partitioning*. arXiv preprint arXiv:1311.3144.
- Buyya, R., Yeo, C., Venugopal, S., Broberg, J., and Brandic, I. (2009). Cloud computing and emerging IT platforms: Vision, hype, and reality for delivering computing as the 5th utility. *Future Generation Computer Systems*, 25(6), 599-616.
- Cao, L. (2011). *GoldenOrb*. Retrieved July 25, 2015, from <https://github.com/jzachr/goldenorb>
- Catalyurek, U., and Aykanat, C. (1996). Decomposing Irregularly Sparse Matrices for Parallel Matrix-Vector Multiplication. *Third International Workshop on Parallel Algorithms for Irregularly Structured Problems (IRREGULAR '96)* (pp. 75-86). Santa Barbara, CA, USA: Springer Berlin Heidelberg.
- Cha, M., Kwak, H., Rodriguez, P., Ahn, Y.-Y., and Moon, S. (2007). I tube, you tube, everybody tubes: analyzing the world's largest user generated content video system. *7th ACM SIGCOMM conference on Internet measurement (IMC '07)*, (pp. 1-14). San Diego, CA.
- Chambers, C., Raniwala, A., Perry, F., Adams, S., Henry, R. R., Bradshaw, R., and Weizenbaum, N. (2010). FlumeJava: Easy, Efficient Data-Parallel Pipelines. *31st ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '10)*, (pp. 363-375). Toronto, Ontario, Canada.
- Chandy, K., and Lamport, L. (1985). Distributed Snapshots: Determining Global States of Distributed Systems. *ACM Transactions on Computer Systems (TOCS)*, 3(1), 63-75.

Che, S. (2014). GasCL: A Vertex-Centric Graph Model for GPUs. *IEEE High Performance Extreme Computing Conference (HPEC)*, (pp. 1-6). Waltham, MA.

Chen, G., Fan, Z., and Li, X. (2005). Modelling the Complex Internet Topology. *Complex Dynamics in Communication Networks*, 213-234.

Chen, Q., Bai, S., Li, Z., Gou, Z., Suo, B., and Pan, W. (2014). *GraphHP: A Hybrid Platform for Iterative Graph Processing*. Technical Report.

Chen, R., Weng, X., He, B., and Yang, M. (2010). Large Graph Processing in the Cloud. *ACM SIGMOD International Conference on Management of data (SIGMOD '10)*, (pp. 1123-1126). Indianapolis, Indiana, USA.

Chen, R., Weng, X., He, B., Yang, M., Choi, B., and Li, X. (2012). Improving Large Graph Processing on Partitioned Graphs in the Cloud. *Third ACM Symposium on Cloud Computing (SoCC '12)*. San Jose, CA.

Chen, Y., Lee, Y.-H., Wong, W., and Guo, D. (2008). A Race Condition Graph for Concurrent Program Behavior. *3rd International Conference on Intelligent System and Knowledge Engineering (ISKE'08)* (pp. 662-667). Xiamen, China: IEEE.

Cheng, R., Hong, J., Kyrola, A., Miao, Y., Weng, X., Wu, M., Chen, E. (2012). Kineograph: Taking the Pulse of a Fast-Changing and Connected World. *7th ACM european conference on Computer Systems (EuroSys '12)*, (pp. 85-98). Bern, Switzerland.

Chockler, G., Guerraoui, R., Keidar, I., and Vukolic, M. (2009). Reliable Distributed Storage. *Computer*, 42(4), 60-67. doi:10.1109/MC.2009.126

Chohan, N., Castillo, C., Spreitzer, M., Steinder, M., Tantawi, A., and Krintz, C. (2010). See Spot Run: Using Spot Instances for MapReduce Workflows. *Proceeding of 2nd USENIX WOkshop on Hot Topics inCloud Computing (HotCloud'10)*. Boston, MA.

Cohen, J. (2009). Graph Twiddling in a MapReduce World. *Computing in Science and Engineering*, 11(4), 29-41.

Commission, T. E. (2010). *Social Networks Overview: Current Trends and Research Challenges*. Information Society and Media. Luxembourg: Publications Office of the European Union.

Coulouris, G., Dollimore, J., Kindberg, T., and Blair, G. (2012). *Distributed Systems Concepts and Design* (5th ed.). Boston, Massachusetts, US: Addison-Wesley.

Committee on the Analysis of Massive Data, Committee on Applied and Theoretical Statistics, Board on Mathematical Sciences and Their Applications, Division on Engineering and Physical Sciences and National Research Council, (2013). *Frontiers in Massive Data Analysis*. The National Academies Press.

DB-Engines Ranking Per Database Model Category. (2015, August 5). Retrieved August 15, 2015, from DB-Engines: [http://db-engines.com/en/ranking\\_categories](http://db-engines.com/en/ranking_categories)

De Mauro, A., Greco, M., and Grimaldi, M. (2014). What is Big Data? A Consensual Definition and a Review of Key Research Topics. *International Conference on Integrated Information (IC-ININFO 2014)* (pp. 97-104). Madrid: AIP Conf. Proc.

Dean, J., and Ghemawat, S. (2004). MapReduce: Simplified Data Processing on Large Clusters.: *Sixth Symposium on Operating Systems Design and Implementation(OSDI '04)*. San Francisco, California, USA : USENIX.

Devine, K. D., Boman, E. G., Heaphy, R. T., Bisseling, R. H., and Catalyurek, U. V. (2006). Parallel Hypergraph Partitioning for Scientific Computing. *20th international conference on Parallel and distributed processing (IPDPS'06)*. Rhodes Island, Greece: IEEE.

Dias de Assuncao, M., N. Calheiros, R., Bianchi, S., A.S. Netto, M., and Buyya, R. (2015). Big Data computing and clouds: Trends and future directions. *Journal of Parallel and Distributed Systems*, 79-80, 3-15.

Dijkstra, E. W. (1959). A note on two problems in connexion with graphs. *Numerische Mathematik*, 1(1), 269-271.

Dindokar, R., and Simmhan, Y. (2016). Elastic Partition Placement for Non-stationary Graph Algorithms. *Proceedings of the 16th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGrid'16)*. Cartagena, Colombia.

Dindokar, R., Choudhury, N., and Simmhan, Y. (2015). Analysis of Subgraph-centric Distributed Shortest Path Algorithm. *Proceeding of International Workshop on Parallel and Distributed Computing for Large Scale Machine Learning and Big Data Analytics (ParLearning'15)*. Hyderabad, India.

Dindokar, R., Choudhury, N., and Simmhan, Y. (2016). A Meta-graph Approach to Analyze Subgraph-centric Distributed Programming Models. *Proceeeding of IEEE International Conference on Big Data*.

Doekemeijer, N., and Varbanescu, A. (2014). *A Survey of Parallel Graph Processing Frameworks*. Delft, The Netherlands: Delft University of Technology.

Dominguez-Sal, D., Martinez-Bazan, N., Munes-Mulero, V., Baleta, P., and Lluís Larriba-Pey, O. (2010). A Discussion on the Design of Graph Database Benchmarks. *Proceedings of the Second TPC Technology Conference (TPCTC'10)* (pp. 25-40). Singapore: Springer.

Ediger, D., and Bader, D. A. (2013). Investigating Graph Algorithms in the BSP Model on the Cray XMT. *27th IEEE International Parallel and Distributed Processing Symposium (IPDPS'13)* (pp. 1638-1645). Boston, Massachusetts, USA: IEEE.

Egwutuoha, I. P., Levy, D., Selic, B., and Chen, S. (2013). A survey of fault tolerance mechanisms and checkpoint/restart implementations for high performance computing systems. *The Journal of Supercomputing*, 65(3), 1302-1326.

Ekanayake, J., Li, H., Zhang, B., Gunarathne, T., and Bae, S.-H. (2010). Twister: a runtime for iterative MapReduce. *19th ACM International Symposium on High Performance Distributed Computing (HPDC '10)*, (pp. 810-818). Chicago, IL, USA.

Elena, D. (2013). *Fault-Tolerant Design*. New York, NY, US: Springer-Verlag New York.

Elnozahy, E., Alvist, L., Wang, Y.-M., and Johnson, D. B. (2002). A survey of rollback-recovery protocols in message-passing systems. *ACM Computing Surveys (CSUR)*, 34(3), 375-408.

El-Rewini, H., and Abd-El-Barr, M. (2005). Message Passing Interface (MPI). In *Advanced Computer Architecture and Parallel Processing* (pp. 205-234).

Elsner, U. (2002). *Static and Dynamic Graph Partitioning. A Comparative Study of Existing Algorithms*. Berlin, Germany: Logos Verlag Berlin.

Ewen, S., Schelter, S., Tzoumas, K., Warneke, D., and Markl, V. (2013). Iterative Parallel Data Processing with Stratosphere: An Inside Look. *ACM SIGMOD International Conference on Management of Data (SIGMOD '13)*, (pp. 1053-1056). New York, NY, USA.

35:30

S. Heidari et al.

- Fard, A., Nisar, M., Ramaswamy, L., Miller, J. A., and Saltz, M. (2013). A Distributed Vertex-Centric Approach for Pattern Matching in Massive Graphs. *IEEE international Conference on Big Data* (pp. 403-411). Silicon Valley, CA: IEEE.
- Fegaras, L., Li, C., and Gupta, U. (2012). An optimization framework for map-reduce queries. *15th International Conference on Extending Database Technology*, (pp. 26-37). Berlin, Germany.
- Feige, U., Hajiaghayi, M., and Lee, J. R. (2005). Improved approximation algorithms for minimum-weight vertex separators. *37 annual ACM symposium on Theory of computing (STOC '05)*, (pp. 563-572). Baltimore, MD.
- Feige, U., Hajiaghayi, M., and Lee, J. R. (2008). Improved Approximation Algorithms for Minimum Weight Vertex Separators. *SIAM Journal on Computing*, 38(2), 629-657.
- Fouss, F., Pirotte, J.-M., Renders, J.-M., and Saerens, M. (2007). Random-Walk Computation of Similarities between Nodes of a Graph with Application to Collaborative Recommendation. *IEEE Transactions on Knowledge and Data Engineering*, 19(3), 355 - 369.
- Fu, Z., Personick, M., and Thompson, B. (2014). MapGraph: A High Level API for Fast Development of High Performance Graph Analytics on GPUs. *Workshop on GRAPh Data management Experiences and Systems (GRADES'14)*, (pp. 1-6). Snowbird, Utah, USA.
- Gehani, N. (1990). Message passing in Concurrent C: Synchronous versus asynchronous. *Software: Practice and Experience*, 20(6), 571-592.
- Geier, F. (2015). *The Differences Between SSD and HDD Technology Regarding Forensic Investigations*. Småland, Sweden: Linnaeus University.
- Geist, A., Beguelin, A., Dongarra, J., Jiang, W., Manchek, R., and Sunderam, V. S. (1994). *PVM: A Parallel Virtual Machine*. Cambridge, MA, USA: Scientific and Engineering Computation Series. MIT Press.
- Gharaibeh, A., Beltrão Costa, L., Santos-Neto, E., and Ripeanu, M. (2013). On Graphs, GPUs, and Blind Dating: A Workload to Processor Matchmaking Quest. *IEEE 27th International Symposium on Parallel and Distributed Processing (IPDPS '13)*, (pp. 851-862). Boston, Massachusetts USA.
- Gharaibeh, A., Reza, T., Santos-Neto, E., Beltrao Costa, L., Sallinen, S., and Ripeanu, M. (2013). *Efficient Large-Scale Graph Processing on Hybrid CPU and GPU Systems*. arXiv preprint arXiv:1312.3018.
- Gonzalez, J. E., Low, Y., Gu, H., Bickson, D., and Guestrin, C. (2012). PowerGraph: Distributed Graph-Parallel Computation on Natural Graphs. *10th USENIX conference on Operating Systems Design and Implementation (OSDI'12)*, (pp. 17-30). Hollywood, CA.
- Grabowski, M., Hidders, J., and Sroka, J. (2013). Representing MapReduce Optimisations in the Nested Relational Calculus. *29th British National Conference on Databases*. Oxford, United Kingdom.
- Graph500. (n.d.). Retrieved 03 25, 2016, from <http://www.graph500.org/>
- GraphBase. (2015). Retrieved August 10, 2015, from <http://graphbase.net/>
- Green, C. (2013, August 14). *An Introduction to Graph Databases*. Retrieved July 28, 2015, from Information Age: <http://www.information-age.com/technology/information-management/123457275/an-introduction-to-graph-databases>
- Gregor, D., and Lumsdaine, A. (2005). The Parallel BGL: A Generic Library for Distributed Graph Computations. *Parallel Object-Oriented Scientific Computing (POOSC)*.
- Gremlin Query Language. (n.d.). (Aurelius) Retrieved 08 22, 2016, from <http://s3.thinkaurelius.com/docs/titan/0.5.4/gremlin.html>
- Gu, Y., Lu, L., Grossman, R., and Yoo, A. (2010). Processing Massive Sized Graphs Using Sector/Sphere. *2010 IEEE Workshop on Many-Task Computing on Grids and Supercomputers (MTAGS)* (pp. 1-10). New Orleans, LA: IEEE.
- Gunelius, S. (2014, July 12). *The Data Explosion in 2014 Minute by Minute – Infographic*. Retrieved July 25, 2015, from ACI: <http://aci.info/2014/07/12/the-data-explosion-in-2014-minute-by-minute-infographic/>
- Guo, Y., Varbanescu, A., Iosup, A., and Epema, D. (2015). An Empirical Performance Evaluation of GPU-Enabled Graph-Processing Systems. *15th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGrid'15)* (pp. 423-432). Shenzhen, China: IEEE.
- Gupta, P., Goel, A., Lin, J., Sharma, A., Wang, D., and Zadeh, R. (2013). WTF: The Who to Follow Service at Twitter. *22nd international conference on World Wide Web (WWW '13)*, (pp. 505-514). Rio de Janeiro, Brazil.
- Han, M., and Daudjee, K. (2015). Giraph Unchained: Barrierless Asynchronous Parallel Execution in Pregel-like Graph Processing Systems. *The VLDB Endowment*, 8(9), 950-961.
- Han, W., Miao, Y., Li, K., Wu, M., Yang, F., Zhou, L., . . . Chen, E. (2014). Chronos: a graph engine for temporal graph analysis. *Ninth European Conference on Computer Systems (EuroSys '14)*. Amsterdam, Netherland.
- Han, W.-S., Lee, S., Park, K., Lee, J.-H., Kim, M.-S., Kim, J., and Yu, H. (2013). TurboGraph: A Fast Parallel Graph Engine Handling Billion-scale Graphs in a Single PC. *19th ACM SIGKDD international conference on Knowledge discovery and data mining (KDD '13)*, (pp. 77-85). Chicago, IL.
- Harshvardhan, Fidel, A., Amato, N. M., and Rauchwerger, L. (2013). The STAPL Parallel Graph Library. *Proceedings of the 25th International Workshop on Languages and Compilers for Parallel Computing (LCPC'12)* (pp. 46-60). Tokyo, Japan: Springer.
- Harshvardhan, Fidel, A., Amato, N. M., and Rauchwerger, L. (2014). KLA: A New Algorithmic Paradigm for Parallel Graph Computation. *23rd international conference on parallel architectures and compilation (PACT '14)*, (pp. 27-38). San Francisco, CA, USA.
- Heidari, S., Calheiros, R. N., and Buyya, R. (2016). iGraph: A Cost-efficient Framework for Processing Large-scale Graphs on Public Clouds. *Proceedings of 16th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGrid '16)*. Cartagena, Colombia.
- Hirschberg, D. S., Chandra, A. K., and Sarwate, D. V. (1979). Computing connected components on parallel computers. *Communications of the ACM*, 22(8), 461-464.
- Hong, S., Chaf, H., Sedlar, E., and Olukotun, K. (2012). Green-Marl: a DSL for easy and efficient graph analysis. *Proceedings of the 17th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS'12)*. London, UK.



Scalable Graph Processing Frameworks: A Taxonomy and Open Challenges 35:31

Hong, S., Chafi, H., Sedlar, E., and Olukotun, K. (2012). Green-Marl: a DSL for easy and efficient graph analysis. *Proceedings of the 17th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '12)*. London, UK.

Hoque, I., and Gupta, I. (2013). LFGGraph: Simple and Fast Distributed Graph Analytics. *First ACM SIGOPS Conference on Timely Results in Operating Systems (TRIOS '13)*. Farmington, Pennsylvania, USA.

Huberman, B. A. (2001). *The Laws of the Web: Patterns in the Ecology of Information*. MIT Press Cambridge.

Hudak, P. (1989). "Conception, evolution, and application of functional programming languages. *ACM Computing Surveys*, 21(3), 359-411.

Hudak, P. (1989). Conception, Evolution, and Application of Functional Programming Languages. *ACM Computing Surveys*, 21(3), 359-411.

HyperGraphDb . (2015). Retrieved August 10, 2015, from <http://www.hypergraphdb.org/index>

InfiniteGraph . (2015). Retrieved August 10, 2015, from <http://www.objectivity.com/products/infinitegraph/>

InfoGrid . (2015). Retrieved August 10, 2015, from <http://infogrid.org/trac/>

Isard, M., Budiu, M., Yu, Y., Birrell, A., and Fetterly, D. (2007). Dryad: Distributed Data-Parallel Programs from Sequential Building Blocks. *2nd ACM SIGOPS/EuroSys European Conference on Computer Systems*, (pp. 59-72). Lisboa, Portugal.

Jackson, J. (2013, 08 14). Facebook's Graph Search puts Apache Giraph on the map. Retrieved 07 25, 2015, from PCWorld: <http://www.pcworld.com/article/2046680/facebooks-graph-search-puts-apache-giraph-on-the-map.html>

Jain , N., Liao , G., and Willke , T. L. (2013). GraphBuilder: Scalable Graph ETL Framework. *First International Workshop on Graph Data Management Experiences and Systems (GRADES '13)*. New York, NY, US.

Jamadagni, N., and Simmhan, Y. (2016). GoDB: From Batch Processing to Distributed Querying over Property Graphs. *Proceedings of the 16th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGrid'16)*. Cartagena, Colombia.

Java Universal Network/Graph Framework. (n.d.). Retrieved 06 25, 2016, from <http://jung.sourceforge.net/>

JCoreDB . (2015). Retrieved August 10, 2015, from <https://sites.google.com/site/jcoredb/>

Jouili, S., and Vansteenbergh, V. (2013). An Empirical Comparison of Graph Databases. *International Conference on Social Computing (SocialCom'13)* (pp. 708 - 715). Alexandria, VA: IEEE.

Kang, U., Tong, H., Sun, J., Lin, C.-Y., and Faloutsos, C. (2011). GBASE: A Scalable and General Graph Management System. *17th ACM SIGKDD international conference on Knowledge discovery and data mining (KDD '11)*, (pp. 1091-1099). San Diego, California.

Kang, U., Tsourakakis, C. E., and Faloutsos, C. (2009). PEGASUS: A Peta-Scale Graph Mining System - Implementation and Observations. *2009 Ninth IEEE International Conference on Data Mining (ICDM '09)* (pp. 229-238). Miami, FL : IEEE.

Karypis , G., and Kumar, V. (1995). Multilevel Graph Partitioning Schemes. *The International Conference on Parallel Processing(ICPP'95)*, (pp. 113–122). Raleigh, NC, US.

Karypis , G., and Kumar, V. (1997). A coarse-grain parallel formulation of multilevel k-way graph partitioning algorithm. *8th SIAM Conference on Parallel Processing for Scientific Computing (SIAM PP'97)*.

Kavila, S. D., Raju, G. P., Satapathy, S. C., Machiraju, A., Kinnera, G., and Rasly, K. (2013). A Survey on Fault Management Techniques in Distributed Computing. In S. C. Satapathy, S. K. Udgata, and B. N. Biswal , *Proceedings of the International Conference on Frontiers of Intelligent Computing: Theory and Applications (FICTA)* (Vol. 199, pp. 593-602). Berlin,Germany: Springer Berlin Heidelberg.

Khayyat, Z., Awara, K., Alonazi, A., Jamjoo, H., Williams, D., and Kalnis, P. (2013). Mizan: A System for Dynamic Load Balancing in Large-scale Graph Processing. *8th ACM European Conference on Computer Systems (EuroSys '13)*, (pp. 169-182). Prague, Czech Republic.

Kim, M., and Candan, K. (2012). SBV-Cut: Vertex-Cut Based Graph Partitioning Using Structural Balance Vertices. *Data and Knowledge Engineering*, 72, 285-303.

Koo, S., Kwon, S., Kim, S., and Chung , T.-S. (2015). Dual RAID technique for ensuring high reliability and performance in SSD. *IEEE/ACIS 14th International Conference on Computer and Information Science (ICIS)*, (pp. 399 - 404). Las Vegas, NV, USA.

Kozen, D. C. (1992). Depth-First and Breadth-First Search. In D. C. Kozen, *The Design and Analysis of Algorithms* (pp. 19-24). New York, NY. US: Springer New York.

Krepska, E., Kielmann, T., Fokkink, W., and Bal, H. (2011). HipG: Parallel Processing of Large-Scale Graphs. *ACM SIGOPS Operating Systems Review*, 45(2), 3-13.

Kyrola , A., and Guestrin, C. (2014). GraphChi-DB: Simple Design for a Scalable Graph Database System - on Just a PC. CoRR abs/1403.0701.

Kyrola, A., Blelloch, G., and Guestrin, C. (2012). GraphChi: Large-Scale Graph Computation on Just a PC. *10th USENIX conference on Operating Systems Design and Implementation (OSDI'12)*, (pp. 31-46 ). Hallywood, USA.

Lawson, C. L., Hanson, R. J., Kincaid, D. R., and Krogh, F. T. (1979). Basic Linear Algebra Subprograms for Fortran. *ACM Transactions on Mathematical Software (TOMS)*, 5(3), 308-323.

Lee, J., Bagheri, B., and Kao , H. (2014). Recent Advances and Trends of Cyber Systems and Big Data Analytics in Industrial Informatics. *IEEE International Conference on Industrial Informatics (INDIN) 2014*. Porto Alegre.

Lee, K.-H., Lee, Y.-J., Choi, H., Chung, Y., and Moon, B. (2011). Parallel Data Processing with MapReduce: A Survey. *ACM SIGMOD Record*, 40(4), 11-20.

Lee, K.-H., Lee, Y.-J., Choi, H., Chung, Y., and Moon, B. (2011). Parallel Data Processing with MapReduce: A Survey. *ACM SIGMOD Record*, 40(4), 11-20.

Leimbach, T., Hallinan, D., Bachlechner, D., Weber, A., Jaglo, M., Hennen, L., . . . Hunt, G. (2014). *Potential and Impacts of Cloud Computing Services and Social Network Websites*. European Parliamentary Research Service. Brussels: Science and Technology Options Assessment (STOA).



35:32

S. Heidari et al.

- Leskovec, J., Lang, K. J., Dasgupta, A., and Mahoney, M. W. (2008). *Community Structure in Large Networks: Natural Cluster Sizes and the Absence of Large Well-Defined Clusters*. Pittsburgh, PA: arXiv:0810.1355 [cs.DS].
- Liu, X., Xiao, L., Kreling, A., and Liu, Y. (2006). Optimizing Overlay Topology by Reducing Cut Vertices. *International workshop on Network and operating systems support for digital audio and video (NOSSDAV '06)*. Newport, Rhode Island.
- Lochert, C., Mauve, M., Füßler, H., and Hartenstein, H. (2005). Geographic routing in city scenarios. *ACM SIGMOBILE- Mobile Computing and Communications Review*, 9(1), 69-72.
- Lorica, B. (2013, April 14). *Single server systems can tackle big data*. Retrieved July 25, 2015, from O'Reilly Radar: <http://radar.oreilly.com/2013/04/single-server-systems-can-tackle-big-data.html>
- Lorica, B. (2014). One year later: some single server systems that can tackle Big Data. In *Big Data Now* (pp. 29-30). Sebastopol, CA: O'Reilly Media Inc.
- Low, Y., Gonzalez, J., Kyrola, A., Bickson, D., Guestrin, C., and Hellerstein, J. (2010). GraphLab: A New Framework For Parallel Machine Learning. *The 26th Conference on Uncertainty in Artificial Intelligence (UAI 2010)*. Catalina Island, USA.
- Low, Y., Gonzalez, J., Kyrola, A., Bickson, D., Guestrin, C., and Hellerstein, J. M. (2012). Distributed GraphLab: A Framework for Machine Learning and Data Mining in the Cloud. *VLDB Endowment*, 5(8), 716-727.
- Lu, Y., Cheng, J., Yan, D., and Wu, H. (2014). Large-Scale Distributed Graph Computing Systems: An Experimental Evaluation. *The VLDB Endowment*, 8(3), 281-292.
- Lumsdaine, A., Gregor, D., Hendrickson, B., and Berry, J. (2007). Challenges in Parallel Graph Processing. *Parallel Processing Letters*, 17(1), 5-20.
- Malewicz, G., Austern, M. H., J. C. Bik, A., Dehnert, J. C., Horn, I., Leiser, N., and Czajkowski, G. (2010). Pregel: A System for Large-Scale Graph Processing. *The 2010 ACM SIGMOD International Conference on Management of data*, (pp. 135-145).
- MapGraph . (2015). Retrieved August 10, 2015, from <http://mapgraph.io/>
- Marburger, A., and Westfechtel, B. (2010). Graph-Based Structural Analysis for Telecommunication Systems. In G. Engels, C. Lewerentz, W. Schafer, A. Schurr, and B. Westfechtel, *Graph Transformations and Model-Driven Engineering* (pp. 363-392). Springer Berlin Heidelberg.
- Martella, C., Logothetis, D., Loukas, A., and Siganos, G. (2015). *Spinner: Scalable Graph Partitioning in the Cloud*. arXiv : <https://arxiv.org/pdf/1404.3861.pdf>.
- McSherry, F., and Schwarzkopf, M. (2015, July 8). *The Impact of Fast Networks on Graph Analytics*. Retrieved August 28, 2015, from <http://www.frankmcsherry.org/pagerank/distributed/performance/2015/07/08/pagerank.html#fn0>
- Mehlhorn, K., and Näher, S. (1995). The LEDA Platform of Combinatorial and Geometric Computing. *Communications of the ACM*, 38(1), 96-102.
- Mittal, S., and Vetter, J. S. (2015). A Survey of Software Techniques for Using Non-Volatile Memories for Storage and Main Memory Systems. *IEEE Transactions on Parallel and Distributed Systems*, PP(99), 1-14.
- Murphy, R. C., and Kogge, P. M. (2007). On The Memory Access Patterns of Supercomputer Applications: Benchmark Selection and Its Implications. *IEEE Transactions on Computers*, 56(7), 937 - 945.
- Murray, D. G., McSherry, F., Isaacs, R., Isard, M., Barham, P., and Abadi, M. (2013). Naiad: A Timely Dataflow System. *24th ACM Symposium on Operating Systems Principles (SOSP '13)*, (pp. 439-455). Farmington, Pennsylvania, USA.
- Najeebullah, K., Khan, K., Nawaz, W., and Lee, Y.-K. (2014). BiShard Parallel Processor: A Disk-Based Processing Engine for Billion-Scale Graphs. *International Journal of Multimedia and Ubiquitous Engineering*, 9(2), 199-212.
- Najeebullah, K., Khan, K., Waqas Nawaz, M., and Lee, Y.-K. (2014). *BPP: Large Graph Storage for Efficient Disk Based*. CoRR abs/1401.2327.
- Neo4j. (2015, August 10). Retrieved August 10, 2015, from <http://neo4j.com/>
- Nguyen, D., Lenharth, A., and Pingali, K. (2013). A Lightweight Infrastructure for Graph Analytics. *24th ACM Symposium on Operating Systems Principles (SOSP '13)*, (pp. 456-471). Farmington, Pennsylvania, USA.
- Nicoara, D., Kamali, S., Daudjee, K., and Chen, L. (2015). *Hermes: Dynamic Partitioning for Distributed Social Network Graph Databases*. Waterloo, Ontario, Canada: University of Waterloo.
- Nilakant, K., Dalibard, V., Roy, A., and Yoneki, E. (2014). PrefEdge: SSD Prefetcher for Large-Scale Graph Traversal. *International Conference on Systems and Storage (SYSTOR'14)* (pp. 1-12). ACM.
- Nitzberg, B., and Lo, V. (1991). Distributed shared memory: A survey of issues and algorithms. *Computer*, 24(8), 52-60.
- Nurvitadhi, E., Weisz, G., Wang, Y., Hurkat, S., Nguyen, M., Hoe, J. C., . . . Guestrin, C. (2014). GraphGen: An FPGA Framework for Vertex-Centric Graph Computation. *IEEE 22nd International Symposium on Field-Programmable Custom Computing Machines (FCCM '14)*, (pp. 25-28). Boston, Massachusetts.
- OrientDB vs Neo4j. (2015). Retrieved August 10, 2015, from OrientDB : <http://orientdb.com/orientdb-vs-neo4j/>
- Ousterhout, J., Agrawal, P., Erickson, D., Kozyrakis, C., Leverich, J., Mazières, D., . . . Stutsman, R. (2010). The Case for RAMClouds: Scalable High-Performance Storage Entirely in DRAM. *ACM SIGOPS Operating Systems Review*, 43(4), 92-105.
- Ousterhout, K., Rasti, R., Ratnasamy, S., Shenker, S., and Chun, B.-G. (2015). Making Sense of Performance in Data Analytics Frameworks. *12th USENIX Symposium on Networked Systems Design and Implementation (NSDI '15)*, (pp. 293-307). Oakland, CA, USA.
- Page, L., Brin, S., Motwani, R., and Winograd, T. (1998). *The PageRank Citation Ranking: Bringing Order to the Web*. Stanford InfoLab.
- Patel, P., Ranabahu, A., and Sheth, A. (2009). *Service Level Agreement in Cloud Computing*. Kno.e.sis Publications.
- Paul, A. (2013). Graph based M2M optimization in an IoT environment. *The 2013 Research in Adaptive and Convergent Systems(RACS '13)*, (pp. 45-46). Montreal, Canada.
- Pellegrini, F. (2011). Current challenges in parallel graph partitioning. *Comptes Rendus Mécanique*, 339(2-3), 90-95.

Pettey, C. (2011, June 27). *Gartner Says Solving 'Big Data' Challenge Involves More Than Just Managing Volumes of Data*. Retrieved July 21, 2015, from Gartner: <http://www.gartner.com/newsroom/id/1731916>

Pike, R., Dorward, S., Griesemer, R., and Quinlan, S. (2005). Interpreting the Data: Parallel Analysis with Sawzall. *Scientific Programming - Dynamic Grids and Worldwide Computing*, 13(4), 277-298.

Plimpton, S. J., and Devine, K. D. (2011). MapReduce in MPI for Large-scale Graph Algorithms. *Parallel Computing*, 37(9), 610-632.

Power, R., and Li, J. (2010). Piccolo: Building Fast, Distributed Programs with Partitioned Tables. *9th USENIX conference on Operating systems design and implementation (OSDI'10)*. Vancouver, Canada.

Prakasam, K., and Chandrasekhar, M. (2010). *JPregel*. Retrieved July 24, 2015, from <http://kowschik.github.io/JPregel/Protocol Buffers>. (n.d.). (Google) Retrieved 05 28, 2016, from <https://github.com/google/protobuf>

Pundir, M., Kumar, M., Leslie, L. M., Gupta, I., and Campbell, R. H. (2016). Supporting On-demand Elasticity in Distributed Graph Processing. *2016 IEEE International Conference on Cloud Engineering (IC2E'16)*. Berlin, Germany.

Rahimian, F., Payberah, A. H., Girdzijauskas, S., and Haridi, S. (2014). Distributed Vertex-Cut Partitioning. In *Distributed Applications and Interoperable Systems* (pp. 186-200). Berlin, Germany: Springer Berlin Heidelberg.

Redekopp, M., Simmhan, Y., and Prasanna, V. K. (2013). Optimizations and Analysis of BSP Graph Processing Models on Public Clouds. *Proceedings of the 27th International Symposium on Parallel and Distributed Processing (IPDPS'13)*. Boston, MA, USA.

Robinson, I., Webber, J., and Eifrem, E. (2015). *Graph Databases* (2nd Edition ed.). Sebastopol, CA: O'Reilly.

Rowstron, A., Narayanan, D., Donnelly, A., O'Shea, G., and Douglas, A. (2012). Nobody ever got fired for using Hadoop on a cluster. *1st International Workshop on Hot Topics in Cloud Data Processing (HotCDP 2012)*. Bern, Switzerland: ACM.

Roy, A., Mihailovic, I., and Zwaenepoel, W. (2013). X-Stream: Edge-centric Graph Processing using Streaming Partitions. *24th ACM Symposium on Operating Systems Principles (SOSP '13)*, (pp. 472-488). Farmington, USA.

Roy, P. (2014). A new memetic algorithm with GA crossover technique to solve Single Source Shortest Path (SSSP) problem. *2014 Annual IEEE India Conference (INDICON)*. Pune, India.

Salihoglu, S., and Widom, J. (2013). GPS: A Graph Processing System. *25th International Conference on Scientific and Statistical Database Management (SSDBM)*. Baltimore, Maryland.

Salihoglu, S., and Widom, J. (2014). Optimizing Graph Algorithms on Pregel-like Systems. *VLDB Endowment*, 7(7), 577-588.

Salihoglu, S., Shin, J., Khanna, V., Truong, B., and Widom, J. (2015). Graft: A Debugging Tool For Apache Giraph. *ACM SIGMOD International Conference on Management of Data (SIGMOD '15)*, (pp. 1403-1408). Melbourne, VIC, Australia.

Sarwat, M., Elnikety, S., He, Y., and Mokbel, M. F. (2013). Horton+: a distributed system for processing declarative reachability queries over partitioned graphs. *Proceedings of the VLDB Endowment*, 5(14), 1918-1929.

Satuluri, V., Parthasarathy, S., and Ruan, Y. (2011). Local graph sparsification for scalable clustering. *ACM SIGMOD International Conference on Management of data (SIGMOD '11)*, (pp. 721-732). Athens, Greece.

Schloegel, K., Karypis, G., and Kumar, V. (2001). Graph Partitioning for High Performance Scientific Simulations. In *CRPC Parallel Computing Handbook* (pp. 491-541). San Francisco, CA, USA: Morgan Kaufmann Publishers Inc.

Sedgewick, R., and Wayne, K. (2011). *Algorithms (4th Edition)*. Upper Saddle River, NJ: Addison-Wesley Professional.

Serrano, D., Bouchenak, S., Kouki, Y., Jr, F. A., Ledoux, T., Lejeune, J., . . . Sens, P. (2015). SLA Guarantees for Cloud Services. *FutureGenerationComputerSystems*.

Shao, B., Wang, H., and Li, Y. (2013). Trinity: A Distributed Graph Engine on a Memory Cloud. *ACM SIGMOD International Conference on Management of Data (SIGMOD '13)*, (pp. 505-516). New York, USA.

Shao, Y., Cui, B., Ma, L., and Yao, J. (2013). PAGE: A Partition Aware Engine for Parallel Graph Computation. *22nd ACM international conference on Conference on information and knowledge management (CIKM '13)*, (pp. 823-828). Burlingame, CA.

Sharma, A., Jiang, J., Bommanavar, P., Larson, B., and Lin, J. (2016). GraphJet: Real-Time Content Recommendations at Twitter. *The VLDB Endowment*, 9(13), 1281-1292.

Shun, J., and Blelloch, G. E. (2013). Ligra: A Lightweight Graph Processing Framework for Shared Memory. *ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP '13)*, (pp. 135-146). Shenzhen, China.

Siek, J., Lee, L.-Q., and Lumsdaine, A. (2002). *The Boost Graph Library: User Guide and Reference Manual*. Upper Saddle River, NJ: Addison-Wesley.

Simmhan, Y., and Kumbhare, A. (2013). *Floe: A dynamic, continuous dataflow framework for elastic clouds*. Los Angeles, CA.

Simmhan, Y., Choudhury, N., Wickramaarachchi, C., and Kumbhare, A. (2015). Distributed Programming over Time-Series Graphs. *Proceedings of IEEE International Conference on Parallel and Distributed Processing Symposium (IPDPS)*. Hyderabad.

Simmhan, Y., Kumbhare, A., Wickramaarachchi, C., Nagarkar, S., Ravi, S., Raghavendra, C., and Prasanna, V. (2014). GoFFish: A Sub-Graph Centric Framework for Large-Scale Graph Analytics. *Euro-Par 2014 Parallel Processing* (pp. 451-462). Porto, Portugal: Springer International Publishing.

Simmhan, Y., Wickramaarachchi, C., Kumbhare, A., Frincu, M., Nagarkar, S., Ravi, S., . . . Prasanna, V. (2014). Scalable analytics over distributed time-series graphs using goffish. *arXiv preprint arXiv:1406.5975*.

Snijders, C., Matzat, U., and Reips, U.-D. (2012). Big Data: Big gaps of knowledge in the field of Internet. *International Journal of Internet Science*, 7(1), 1-5.

Stanton, I., and Kliot, G. (2012). Streaming Graph Partitioning for Large Distributed Graphs. *18th ACM SIGKDD international conference on Knowledge discovery and data mining (KDD '12)*, (pp. 1222-1230). Beijing, China.

Stelzner, M. A. (2015). *2015 Social Media Marketing Industry Report*. Social Media Examiner.

Strandmark, P., and Kahl, F. (2011). Parallel and distributed graph cuts by dual decomposition. *Computer Vision and Image Understanding*, 115(12), 1721-1732.

Stutz, P., Bernstein, A., and Cohen, W. (2010). Signal/Collect: Graph Algorithms for the (Semantic) Web. *9th international semantic web conference on The semantic web (ISWC'10)*, (pp. 764-780). Shanghai, China.

35:34

S. Heidari et al.

- Suri, S., and Vassilvitskii, S. (2011). Triangles and the Curse of the Last Reducer. *20th international conference on World wide web (WWW '11)*, (pp. 607-614). Hyderabad, India.
- Szalay, A. S. (2011). Extreme Data-Intensive Scientific Computing. *Computing in Science and Engineering*, 13(6), 34-41.
- Tasci, S., and Demirbas, M. (2013). Giraphx: Parallel Yet Serializable Large-Scale Graph Processing. *19th international conference on Parallel Processing (Euro-Par '13)*, (pp. 458-469). Aachen, Germany.
- Thien Bao, N., and Suzumura, T. (2013). Towards Highly Scalable Pregel-Based Graph Processing Platform With x10. *22nd International Conference on World Wide Web (WWW '13 Companion)*, (pp. 501-508). Rio de Janeiro, Brazil.
- Tian, J., Hahner, J., Becker, C., Stepanov, I., and Rothermel, K. (2002). Graph-based mobility model for mobile ad hoc network simulation. *35th Annual Simulation Symposium*, (pp. 337-344). San Diego, California.
- Tian, Y., Balmin, A., Andreas Corsten, S., Tatikond, S., and McPherson, J. (2013). From "Think Like a Vertex" to "Think Like a Graph". *The Proceedings of the VLDB Endowment*, 7(3), 193-204.
- Tian, Y., Hankins, R. A., and Patel, J. M. (2008). Efficient aggregation for graph summarization. *ACM SIGMOD international conference on Management of data (SIGMOD '08)*, (pp. 567-580). Vancouver, BC, Canada.
- TITAN Distributed Graph Database. (2015). Retrieved August 10, 2015, from <http://thinkaurelius.github.io/titan/>
- Top500. (n.d.). Retrieved 06 25, 2016, from <https://www.top500.org/>
- Treaster, M. (2005). A Survey of Fault-Tolerance and Fault-Recovery Techniques in Parallel Systems. *ACM Computing Research Repository (CoRR)*, 1-11.
- Tsourakakis, C. E., Gkantsidis, C., Radunovic, B., and Vojnovic, M. (2014). FENNEL: Streaming Graph Partitioning for Massive Scale Graphs. *7th ACM international conference on Web search and data mining (WSDM '14)*, (pp. 333-342). New York, NY, US.
- Twitter. (2012, March 08). *Cassovary: A Big Graph-Processing Library*. Retrieved July 24, 2015, from Twitter Blog: <https://blog.twitter.com/2012/cassovary-big-graph-processing-library>
- Twitter. (2012, April 29). *Twitter/FlockDB*. Retrieved July 28, 2015, from <https://github.com/twitter/flockdb#readme>
- Valiant, L. G. (1990). A bridging model for parallel computation. *Communications of the ACM*, 33(8), 103-111.
- Vaquero, L. M., Cuadrado, F., Logothetis, D., and Martella, C. (2013). xDGP: A Dynamic Graph Processing System with Adaptive Partitioning. *4th annual Symposium on Cloud Computing (SOCC '13)*. Santa Clara, CA.
- Venkataraman, S., Bodzsar, E., Roy, I., AuYoung, A., and Schreiber, R. S. (2013). Presto: Distributed Machine Learning and Graph Processing with Sparse Matrices. *8th ACM European Conference on Computer Systems (EuroSys '13)*, (pp. 197-210). Prague, Czech Republic.
- Vicknair, C., Macias, M., Zhao, Z., Nan, X., Chen, Y., and Wilkins, D. (2010). A Comparison of a Graph Database and a Relational Database. *48th Annual Southeast Regional Conference (ACM SE '10)*. Oxford, Mississippi.
- Wang, G., Xie, W., Demers, A., and Gehrke, J. (2013). Asynchronous LargeScale Graph Processing Made Easy. *6th biennial Conference on Innovative Data Systems Research (CIDR'13)*. Asilomar, USA.
- Wang, P., Zhang, K., Chen, R., Chen, H., and Guan, H. (2014). Replication-based Fault-tolerance for Large-scale Graph Processing. *44th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)*, (pp. 562 - 573). Atlanta, GA.
- Weaver. (2015). Retrieved August 10, 2015, from <http://weaver.systems/>
- Xie, C., Chen, R., Guan, H., Zang, B., and Chen, H. (2015). SYNC or ASYNC: Time to Fuse for Distributed Graph-parallel Computation. *20th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP '15)*, (pp. 194-204). San Francisco, CA.
- Xin, R. S., Gonzalez, J. E., Franklin, M. J., and Stoica, I. (2013). GraphX: a resilient distributed graph system on Spark. *First International Workshop on Graph Data Management Experiences and Systems (GRADES '13)*. New York, NY, USA.
- Xu, N., Chen, L., and Cui, B. (2014). LogGP: A Log-based Dynamic Graph Partitioning Method. *The VLDB Endowment*, 7(14), 1917-1928.
- Xue, J., Yang, Z., Qu, Z., Hou, S., and Dai, Y. (2014). Seraph: An Efficient, Low-cost System for Concurrent Graph Processing. *23rd international symposium on High-performance parallel and distributed computing (HPDC '14)*, (pp. 227-238). Vancouver, Canada.
- Yamato, Y. (2015). Use case study of HDD-SSD hybrid storage, distributed storage and HDD storage on OpenStack. *19th International Database Engineering and Applications Symposium (IDEAS '15)*, (pp. 228-229). Yokohoma, Japan.
- Yan, D., Cheng, J., Lu, Y., and Ng, W. (2014). Blogel: A Block-Centric Framework for Distributed Computation on Real-World Graphs. *Proceedings of the VLDB Endowment*, 7(14), 1981-1992.
- Yan, D., Cheng, J., Ozsu, M., Yang, F., Lu, Y., Lui, J. C., . . . Ng, W. (2016). Quegel: A General-Purpose Query-Centric Framework for Querying Big Graphs. *The VLDB Endowment*, 9(7), 564-575.
- Yan, J., Tan, G., and Sun, N. (2013). GRE: A Graph Runtime Engine for Large-Scale Distributed Graph-Parallel Applications. *CoRR abs/1310.5603*.
- Yang, S., Yan, X., Zong, B., and Khan, A. (2012). Towards Effective Partition Management for Large Graphs. *ACM SIGMOD International Conference on Management of Data (SIGMOD '12)*, (pp. 517-528). Scottsdale, Arizona, USA.
- Yoneki, E., Nilakant, K., Dalibard, V., and Roy, A. (2014). PrefEdge: SSD Prefetcher for Large-Scale Graph Traversal. *International Conference on Systems and Storage (SYSTOR 2014)*. ACM.
- Yu, Y., Isard, M., Fetterly, D., Budiu, M., Erlingsson, U., Kumar Gunda, P., and Currey, J. (2008). DryadLINQ: A System for General-Purpose Distributed Data-Parallel Computing Using a High-Level Language. *8th USENIX conference on Operating systems design and implementation (OSDI'08)*, (pp. 1-14). San Diego, CA.
- Yuan, P., Zhang, W., Xie, C., Jin, H., Liu, L., and Lee, K. (2014). Fast Iterative Graph Computation: A Path Centric Approach. *International Conference for High Performance Computing, Networking, Storage and Analysis (SC '14)*, (pp. 401-412). New Orleans, LA.



1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Scalable Graph Processing Frameworks: A Taxonomy and Open Challenges35:35

Zaharia, M., Chowdhury, M., Das, T., Dave, A., Ma, J., McCauley, M., . . . Stoica, I. (2012). Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing. *9th USENIX conference on Networked Systems Design and Implementation (NSDI'12)*. San Jose, CA.

Zha, Z.-J., Mei, T., Wang, J., Wang, Z., and Hua, X.-S. (2009). Graph-based semi-supervised learning with multiple labels. *Journal of Visual Communication and Image Representation*, 20(2), 97-103.

Zhang, T., Zhang, J., Shu, W., Wu, M.-Y., and Liang, X. (2015). Efficient graph computation on hybrid CPU and GPU systems. *The Journal of Supercomputing*, 71(4), 1563-1586.

Zhang, Y., Gao, Q., Gao, L., and Wang, C. (2012). Accelerate large-scale iterative computation through asynchronous accumulative updates. *3rd workshop on Scientific Cloud Computing Date (ScienceCloud '12)*, (pp. 13-22). Delft, the Netherlands.

Zhang, Y., Gao, Q., Gao, L., and Wang, C. (2012). PrIter: A Distributed Framework for Prioritized Iterative Computations. *IEEE Transactions on Parallel and Distributed Systems*, 24(9), 1884-1893.

Zhang, Y., Gao, Q., Gao, L., and Wang, C. (2014). Maiter: An Asynchronous Graph Processing Framework for Delta-Based Accumulative Iterative Computation. *IEEE Transaction on Parallel and Distributed Systems*, 25(8), 2091-2100.

Zheng, D., Mhembere, D., Burns, R., Vogelstein, J., Priebe, C. E., and Szalay, A. S. (2015). FlashGraph: processing billion-node graphs on an array of commodity SSDs. *13th USENIX Conference on File and Storage Technologies (FAST'15)*, (pp. 45-58). Santa Clara, CA.

Zhong, J., and He, B. (2013). Medusa: Simplified Graph Processing on GPUs. *IEEE Transactions on Parallel and Distributed Systems (TPDS)*, 25(6), 1543 - 1552.

Zhong, J., and He, B. (2013). Towards GPU-Accelerated Large-Scale Graph Processing in the Cloud. *5th International Conference on Cloud Computing Technology and Science (CloudCom'13)* (pp. 9-16). Bristol, UK: IEEE.

Zhou, J., Bruno, N., Wu, M.-C., Larson, P.-A., Chaiken, R., and Shakib, D. (2012). SCOPE: parallel databases meet MapReduce. *The VLDB Journal — The International Journal on Very Large Data Bases*, 21(5), 611-636



Dear Editor of the ACM Computing Surveys,

We would like to bring to your attention the manuscript entitled "Scalable Graph Processing Frameworks: A Taxonomy and Open Challenges". This paper presents a comprehensive taxonomy on graph processing systems in several key aspects: computational models, programming models, coordination, execution models, and partitioning; and map the existing systems to this taxonomy. It also covers benefits and challenges of various graph processing implementation approaches along with identifying gaps which need more investigation and discussing open problems and future research directions.

This paper represents a complete new contribution, and no part of it was published before. We believe this paper will be a significant contribution to your journal.

Yours sincerely,

The authors