

# Streaming Graph Partitioning for Large Distributed Graphs

Isabelle Stanton  
University of California Berkeley  
Berkeley, CA  
isabelle@eecs.berkeley.edu

Gabriel Kliot  
Microsoft Research  
Redmond, WA  
gkliot@microsoft.com

## ABSTRACT

Extracting knowledge by performing computations on graphs is becoming increasingly challenging as graphs grow in size. A standard approach distributes the graph over a cluster of nodes, but performing computations on a distributed graph is expensive if large amount of data have to be moved. Without partitioning the graph, communication quickly becomes a limiting factor in scaling the system up. Existing graph partitioning heuristics incur high computation and communication cost on large graphs, sometimes as high as the future computation itself. Observing that the graph has to be loaded into the cluster, we ask if the partitioning can be done at the same time with a lightweight streaming algorithm.

We propose natural, simple heuristics and compare their performance to hashing and METIS, a fast, offline heuristic. We show on a large collection of graph datasets that our heuristics are a significant improvement, with the best obtaining an average gain of 76%. The heuristics are scalable in the size of the graphs and the number of partitions. Using our streaming partitioning methods, we are able to speed up PageRank computations on Spark [32], a distributed computation system, by 18% to 39% for large social networks.

## Categories and Subject Descriptors

G.2.2 [Mathematics of Computing]: Discrete Mathematics—*Graph Theory, Graph Algorithms*; D.2.8 [Software Engineering]: Metrics—*Complexity measures, Performance measures*

## General Terms

Algorithms, Experimentation

## Keywords

Graph Partitioning, Streaming Algorithms, Distributed Graphs, Experimental Evaluation

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

KDD'12, August 12–16, 2012, Beijing, China.

Copyright 2012 ACM 978-1-4503-1462-6 /12/08 ...\$10.00.

## 1. INTRODUCTION

Modern graph datasets are huge. The clearest example is the World Wide Web where crawls by large search engines currently consist of over one trillion links and are expected to exceed ten trillion within the year. Individual websites also contain enormous graph data. In Jan 2012, Facebook consisted of over 800 million active users, with hundreds of billions friend links [1]. There are over 900 million additional objects (communities, pages, events, etc.) that interact with the user nodes. In July 2009, Twitter had over 41.7 million users with over 1.47 billion social relations [21]. Since then, it has been estimated that Twitter has grown to over 200 million users. Examples of large graph datasets are not limited to the Internet and social networks - biological networks, like protein interaction networks, are of a similar size. Despite the size of these graphs, it is still necessary to perform computations over the data, such as calculating PageRank, broadcasting Twitter updates, identify protein associations, as well as many other applications.

The graphs consist of terabytes of compressed data when stored on disks and are all far too large for a single commodity type machine to efficiently perform computations. A standard solution is to split the data across a large cluster of commodity machines and use parallel, distributed algorithms for the computation. This approach introduces a host of systems engineering problems of which we focus only on the problem of data layout. For graph data, this is called *balanced graph partitioning*. The goal is to minimize the number of cross partition edges, while keeping the number of nodes (or edges) in every partition approximately even.

Good graph partitioning algorithms are very useful for many reasons. First, graphs that we encounter and care about in practice are not random. The edges display a great deal of locality, whether due to the vertices being geographically close in social networks, or related by topic or domain on the web. This locality gives us hope that good partitions, or at least partitions that are significantly better than random cuts, exist in real graphs. Next, inter-machine communication, even on the same local network, is substantially more expensive than inter-processor communication. Network latency is measured in microseconds while inter-process communication is measured in nanoseconds. This disparity substantially slows down processing when the network must be used. For large graphs, the data to be moved may border on gigabytes, causing network links to become saturated.

The primary problem with partitioning complicated graph data is that it is difficult to create a linear ordering of the data that maintains locality of the edges i.e. if it is possible

to embed the vertices of a graph into a line such that none of the edges are ‘too long’, then a good balanced cut exists in the graph. Such an ordering may not even exist at all. There is a strong connection between graph partitioning and the eigenvectors and eigenvalues of the corresponding Laplacian matrix of the graph via the Cheeger bound. This connection has inspired many spectral solutions to the problem, including ARV [11] and the many works that followed.

However, spectral methods do not scale to big data. This is in part due to the running time and in part because current formulations require full graph information. When a graph does not physically fit on one machine, maintaining a coherent view of the entire state is impossible. This has led to local spectral partitioning methods, like EvoCut [9], but local methods still require access to large portions of the graph, rely on complex distributed coordination and large computation after the data has been loaded. Thus, we look for a new type of solution. A *graph loader* is a program that reads serial graph data from a disk onto a cluster. It must make a decision about the location of each node as it is loaded. The goal is to find a close to optimal balanced partitioning with as little computation as possible. This problem is also called *streaming graph partitioning*.

For some graphs, partitioning can be entirely bypassed by using meta data associated with the vertices, e.g. clustering web pages by URL produces a good partitioning for the web. In social networks, people tend to be friends with people who are geographically nearby. When such data is available, this produces an improved cut over a node ID hashing approach. Unfortunately, this data is not always be available, and even if it is, it is not always clear which features are useful for partitioning. Our goal in this work is to find a general streaming algorithm that relies only on the graph structure and works regardless of the meta data.

## 1.1 Applications

Our motivating example for studying this problem is a large distributed graph computation system. All distributed computation frameworks, like MapReduce, Hadoop, Orleans [15] and Spark [32] have methods for handling the distribution of data across the cluster. Unfortunately, for graphs, these methods are not tuned to minimize communication complexity, and saturating the network becomes a significant barrier to scaling the system up.

The interest in building distributed systems for graph computation has recently exploded, especially within the database community. Examples of these systems include Pregel [25], GraphLab [24], InfiniteGraph, HyperGraphDB, Ne04j, and Microsoft’s Trinity [4] and Horton [3], to name but a few. Even for these graph specific systems, the graphs are laid out using a hash of the node ID to select a partition. If a good pseudorandom hash function is chosen, this is equivalent to using a random cut as graph partitioning and will result in approximately balanced partitions. However, computations on the graph run more slowly when a hash partitioning is used instead of a better partitioning, due to the high communication cost. Fortunately, these systems tend to support custom partitioning, so it is relatively easy to substitute a more sophisticated method, provided it scales to the size of the graph. As our experiments show, even using our simple streaming partitioning techniques can allow systems of this type to complete computations at least 20-40% faster.

## 1.2 Theoretical Difficulties - Lower Bound

Theoretically, a good streaming partitioning algorithm is impossible. It is easy to create graphs and orderings for any algorithm that will cause it to perform poorly. A simple example of this lower bound is a cycle. The optimal balanced 2-partition cuts only 2 edges. However, if the vertices are given in an order of ‘all even nodes then all odd nodes’, we won’t observe any edges until the odd nodes arrive. Without any information, the best an algorithm could do is to try and balance the number of vertices it has seen across the 2 partitions. This leads to an expected cut of  $\frac{n}{4}$  edges. The worst algorithm might put all even nodes in one partition leading to all edges being cut!

We can partially bypass this problem by picking the input ordering. The three popular orderings considered in the literature are adversarial, random, and stochastic. The above cycle example is an adversarial order and demonstrates that the streaming graph partitioning problem may have arbitrarily bad solutions under that input model. Given that adversarial input is unrealistic in our setting - we have control over the data - we focus on input that results from either a random ordering, or the output of a graph search algorithm. The second option is a simplification of the ordering returned by a graph crawler.

**THEOREM 1.** *One-pass streaming balanced graph partitioning with an adversarial stream order can not be approximated within  $o(n)$ .*

**THEOREM 2.** *One-pass streaming balanced graph partitioning with a random stream order can not be approximated within  $o(n)$ .*

The full proofs are elided due to space. Theorem 1 can be proved with the above cycle discussion. Theorem 2 can be proved by analyzing the expected number of vertices in a random order that arrive with no edges (a third). Each of these vertices contributes an edge to the cut in expectation. This means, in expectation, a random cut will cut  $O(n)$  edges, while an optimal partitioning cuts only 2.

We leave open the theoretical problem of analyzing the performance of our heuristics. While there are current approaches to analyzing streaming algorithms, our use of breadth-first and depth-first stream orders is novel and previous approaches can not be applied. Theorem 1 shows we should not hope to analyze any algorithm with an adversarial order, while a random ordering will always hide edges in sparse graphs until  $O(\sqrt{n})$  vertices arrive, making competitive analysis difficult.

## 1.3 The Streaming Model

We consider a simple streaming graph model. We have a cluster of  $k$  machines, each with memory capacity  $C$ , such that the total capacity,  $kC$ , is large enough to hold the whole graph. The graph is  $G = (V, E)$  where  $V$  is the vertices, and  $E$  the edges. The graph may be either directed or undirected. The vertices arrive in a stream with the set of edges where it is a member so for undirected graphs, each edge appears twice in the stream. We consider three orders: random, breadth-first search and depth-first search. As vertices arrive, a partitioner decides to place the vertex on one of the  $k$  machines. A vertex is never moved after it has been placed. In order to give the heuristics maximal flexibility, we allow the partitioning algorithm access to the entire subgraph defined by all vertices previously seen. This is a strong

assumption, but the heuristics studied in this paper use only local (depth 1) information about this subgraph. We extend the model by allowing a buffer of size  $C$  so that the partitioning algorithm may decide to place any node in the buffer, rather than the one at the front of the stream.

Our model assumes serial input and a single loader. This is somewhat unrealistic for a real system where there may be many graph loaders working in parallel on independent portions of the stream. While we will not explore this option, the heuristics we investigate can be easily adapted to a parallel setting where each loads its portion of the graph independently from the others, sharing information only through a distributed lookup table of vertices to partition IDs.

## 1.4 Contributions

We provide a rigorous, empirical study of a set of natural heuristics for streaming balanced graph partitioning. To the best of our knowledge, this is a first attempt to study this problem. We evaluate these heuristics on a large collection of graph datasets, from various domains: the World Wide Web, social networks, finite-element meshes and synthetic datasets from some popular generative models - preferential-attachment [13], RMAT [23] and Watts-Strogatz [30]. We compare the results of our streaming heuristics to both the hash based partitioning, and METIS [19], a well-regarded, fast, offline partitioning heuristic.

Our results show that some of the heuristics are good and some are surprisingly bad. Our best performing heuristic is a weighted variant of the greedy algorithm. It has a significant improvement over the hashing approach without significantly increasing the computational overhead and obtains an average gain of 76% of the possible improvement in the number of edges cut. On some graphs, with some orderings, a variety of heuristics obtain results which are very close to the offline METIS result. By using the synthetic datasets, we are also able to show that our heuristics scale with the size of the graph and the number of partitions. We demonstrate the value of the best heuristic by using it to partition both the LiveJournal and the Twitter graph for PageRank computation using the Spark cluster system [32]. These are large crawls of real social networks, and we are able to improve the running time of the PageRank algorithm by 18% to 39% by changing the data layout alone. Our experimental results motivate us to recommend that this is an interesting problem worthy of future research and is a viable preprocessing step for graph computation systems.

Our streaming partitioning is not intended to substitute for a full information graph partitioning. Certain systems or applications that need as good a partitioning as possible will still want to repartition the graph after it has been fully loaded onto the cluster. These systems can still greatly benefit from our optimization as a distributed offline partitioning algorithm started from an already reasonably partitioned graph will require less communication and may need to move fewer vertices, causing it to run faster. Our streaming partitioning algorithms can be viewed as a preprocessing optimization step that cannot hurt in exchange for a very small additional computation cost for every loaded vertex.

## 2. RELATED WORK

Graph partitioning has a rich history. It encompasses many problems and has many proposed solutions, from the very simple to the very sophisticated. We cannot hope to

cover the whole field and will only focus on the most relevant formulation - balanced  $k$ -partitioning. The goal is, given a graph  $G$  as input and a number  $k$ , to cut  $G$  into  $k$  balanced pieces while minimizing the number of edges cut. This problem is known to be NP-Hard, even if one relaxes the balanced constraint to ‘approximately’ balanced [10]. Andreev and Racke give an LP-based solution that obtains a  $O(\log n)$  approximation [10]. Even *et al.* [16] provide another LP formulation based on spreading metrics that also obtains an  $O(\log n)$  approximation. Both require full information about the graph. There are many heuristics that solve this problem with an unknown performance guarantee, like METIS [19], PMRSB [14], and Chaco [17]. In practice, these heuristics are quite effective, but many are intended for scientific computing purposes. One can recursively use any balanced 2-partitioning algorithm to approximate a balanced  $k$ -partitioning when  $k$  is a power of 2 [11].

While we are unaware of any previous work on the exact problem statement that we study - one pass balanced  $k$  partitioning - there has been much work on many related streaming problems, primarily graph sparsification in the semi-streaming model, cut projections in the streaming model as well as online algorithms in general, like online bipartite matching. This work includes both algorithms and lower bounds on space requirements.

The work on streaming graph problems where multiple passes on the stream are allowed includes estimating PageRank [29] and cut projections [28]. While PageRank has been used for local partitioning [8], the approach in [8] uses personalized PageRank vectors which does not easily generalize the approach in [29]. Additionally, cut projections do not maintain our balanced criterion.

Bahmani *et al.* [12] maintain an accurate estimate of the PageRank in one pass in the semi-streaming model, where the nodes are fixed but edges arrive in an adversarial order. In the semi-streaming model, further results are known with regards to finding minimum cuts, i.e. no balance requirement. Jin Ahn and Guha [6] give a one pass  $\tilde{O}(n/\epsilon^2)$  space algorithm that sparsifies a graph such that each cut is approximated to within a  $(1 + \epsilon)$  factor. Kelner and Levin [20] produce a spectral sparsifier with  $O(n \log n/\epsilon^2)$  edges in  $\tilde{O}(m)$  time. While sparsifiers are related to partitioning, we do not have a prespecified computation task, so we cannot be sure that a sparsified graph will give accurate answers. Zelke [33] shows lower bounds of  $o(n^2)$  space to find max and min cuts in one pass. Zelke [33] has shown that this cannot be computed in one pass with  $o(n^2)$  space. By contrast, our methods require only access to a distributed lookup table and a buffer of size  $C$ .

## 3. HEURISTICS AND STREAM ORDERS

In this paper, we examine multiple heuristics and stream orders. We now formally define each one.

### 3.1 Heuristics

The notation  $P^t$  refers to the set of partitions at time  $t$ . Each individual partition is referred to by its index  $P^t(i)$  so  $\cup_{i=1}^k P^t(i)$  is equal to all of the vertices placed so far. Let  $v$  denote the vertex that arrives at time  $t$  in the stream,  $\Gamma(v)$  refers to the set of vertices that  $v$  neighbors and  $|S|$  refers to the number of elements in a set  $S$ .  $C$  is the capacity constraint on each partition. Each of the heuristics gives an algorithm for selecting the index *ind* of the partition where

$v$  is assigned. The first seven heuristics do not use a buffer, while the last three do.

1. **Balanced** - Assign  $v$  to a partition of minimal size, breaking ties randomly:

$$ind = \arg \min_{i \in [k]} \{|P^t(i)|\}$$

2. **Chunking** - Divide the stream into chunks of size  $C$  and fill the partitions completely in order:

$$ind = \lceil t/C \rceil$$

3. **Hashing** - Given a hash function  $H : V \rightarrow \{1 \dots k\}$ , assign  $v$  to  $ind = H(v)$ . We use:

$$H(v) = (v \bmod k) + 1$$

4. **(Weighted) Deterministic Greedy** - Assign  $v$  to the partition where it has the most edges. Weight this by a penalty function based on the capacity of the partition, penalizing larger partitions. Break ties using **Balanced**.

$$ind = \arg \max_{i \in [k]} \{|P^t(i) \cap \Gamma(v)|w(t, i)\}$$

where  $w(t, i)$  is a weighted penalty function:

$w(t, i) = 1$  for unweighted greedy

$w(t, i) = 1 - \frac{|P^t(i)|}{C}$  for linear weighted

$w(t, i) = 1 - \exp\{|P^t(i)| - C\}$  for exponentially weighted

5. **(Weighted) Randomized Greedy** - Assign  $v$  according to the distribution defined by

$$Pr(i) = |P^t(i) \cap \Gamma(v)|w(t, i)/Z$$

where  $Z$  is the normalizing constant and  $w(t, i)$  is the above 3 penalty functions.

6. **(Weighted) Triangles** - Assign  $v$  according to

$$\arg \max_{i \in [k]} \left\{ \frac{|E(P^t(i) \cap \Gamma(v), P^t(i) \cap \Gamma(v))|}{|P^t(i) \cap \Gamma(v)|} w(t, i) \right\}$$

where  $w(t, i)$  is the above 3 penalty functions and  $E(S, T)$  is the set of edges between the nodes in  $S$  and  $T$ .

7. **Balance Big** - Given a way of differentiating high and low degree nodes, if  $v$  is high-degree, use **Balanced**. If it is low-degree, use **Deterministic Greedy**.

The following heuristics all use a buffer.

8. **Prefer Big** - Maintain a buffer of size  $C$ . Assign all high degree nodes with **Balanced**, and then stream in more nodes. If the buffer is entirely low degree nodes, then use **Deterministic Greedy** to clear the buffer.
9. **Avoid Big** - Maintain a buffer of size  $C$  and a threshold on large nodes. Greedily assign all small nodes in the buffer. When the buffer is entirely large nodes, use **Deterministic Greedy** to clear the buffer.
10. **Greedy EvoCut** - Use EvoCut [9] on the buffer to find small Nibbles with good conductance. Select a partition for each Nibble using **Deterministic Greedy**.

Each of these heuristics has a different motivation with some arguably more natural than others. **Balanced** and **Chunking** are simple ways of load balancing while ignoring the graph structure.

**Hashing** is currently used by many real systems [25]. The benefit of **Hashing** is that every vertex can be quickly found, from any machine in the cluster, without the need to maintain a distributed mapping table. If the IDs of the nodes are consecutive, the hash function  $H(v) = (v \bmod k) + 1$  makes **Balanced** and **Hashing** equivalent. More generally, a pseudorandom hash function should be used, making **Hashing** equivalent to a random cut.

The greedy approach is standard, although the weighted penalty is inspired by analysis of other online algorithms. The randomized versions of these algorithms were explored because adding randomness can often be shown to theoretically improve the worst-case performance.

The **(Weighted) Triangles** heuristic exploits work showing that social networks have high clustering coefficients by finding triangles completed triangles among the vertices neighbors in a partition and overweighting their importance.

Heuristics **Balance Big**, **Prefer Big**, and **Avoid Big** assume we have a way to differentiate high and low degree nodes. This assumption is based on the fact that many graphs have power law degree distributions. These three heuristics propose different treatments for the small number of high degree nodes and the large number of low degree nodes.

**Balance Big** uses the high degree nodes as seeds for the partitions to ‘attract’ the low degree nodes. The buffered version, **Prefer Big**, allows the algorithm more choice in finding these seeds. **Avoid Big** explores the idea that the high degree nodes form the expander portion of the graph, so perhaps the low degrees nodes can be partitioned after the high degree nodes have been removed.

The final heuristic, **Greedy EvoCut**, uses EvoCut [9], a local partitioning algorithm, on the buffer. This algorithm has very good theoretical guarantees with regards to the found cuts, and the amount of work spent to find them, but the guarantees do not apply to the way we use it.

**Edge Balancing:** While our experiments focus on partitions that are node balanced, in part because this is what our comparator algorithm METIS produces, there is nothing that prevents these heuristics from being used to produce edge-balanced partitions instead, i.e. each partition holds at most  $C = (1 + \epsilon)|E|/k$  edges. An edge-balanced partition may be preferable for power-law distributed graphs when the computation to be performed has complexity in terms of the number of edges and not the number of vertices. In fact, in our second set of experiments with the PageRank algorithm in Section 6 we used the edge-balanced versions of the algorithms instead, for the above reason.

### 3.2 Stream Orders

In a sense, the stream ordering is the key to having a heuristic perform well. A simple example is **Chunking**, where, if we had an optimal partitioning, and then created an ordering consisting of all nodes in partition 1, then all nodes in partition 2 and so on, **Chunking** would also return an optimal partition. For each heuristic, we can define optimal orderings, but, unfortunately, actually generating them reduces to solving balanced graph partitioning so we must settle for orderings that are easy to compute.

We consider the following three stream orderings:

- *Random* - This is a standard ordering in streaming literature and assumes that the vertices arrive in an order given by a random permutation of the vertices.
- *BFS* - This ordering is generated by selecting a starting node from each connected component of the graph uniformly at random and is the result of a breadth-first search that starts at the given node. If there are multiple connected components, the component ordering is done at random.
- *DFS* - This ordering is identical to the *BFS* ordering except that depth-first search is used.

Each of these stream orderings has a different justification. The random ordering is a standard assumption when theoretically analyzing streaming algorithms. While we generate these orderings by selecting a random permutation of the vertices, one could view this as a special case of a generic ordering that does not respect connectivity of the graph. The benefit of a random ordering is that it avoids adversarially bad orderings. The downside is that it does not preserve any locality in the edges so we expect it to do poorly for statistical reasons like the Birthday paradox. Via the Birthday paradox, we can argue that for sparse graphs, we expect to go through  $O(\sqrt{n})$  of the vertices before we find a first edge.

Both BFS and DFS are natural ways of linearizing graphs and are highly simplified models of a web crawler. In practice, web crawlers are a combination of local search approaches - they follow links, but fully explore domains and sub-domains before moving on. This is breadth-first search between domains, and depth-first search within. The main benefit of both orderings is that they guarantee that the partitioner sees edges in the stream immediately. Additionally, they maintain some locality. Each has their drawbacks, but it should be noted that BFS is a subroutine that is often used in partitioning algorithms to find a good cut, particularly for rounding fractional solutions to LPs [16].

## 4. EVALUATION SETUP

We conducted extensive experimental evaluation to discover the performance and trends of stream partitioning heuristics on a variety of graphs. The questions we ask are: Which of these heuristics are reasonable? Can we recommend a best heuristic, restricted to graph type? Do these heuristics scale to larger graphs? Our intent is to use this style of solution for graphs that include trillions of edges, yet in our initial experiments our largest graph has 1.4 million edges. We address this last question by using synthetic datasets to show that the heuristics scale and in Section 6 use our heuristics on two larger social networks successfully.

### 4.1 Datasets

We used several sources to collect multiple datasets for our experiments. From the SNAP [22] archive, we used soc-Slashdot0811, wiki-Vote and web-NotreDame. From the Graph Partitioning Archive [5] we used : 3elt, 4elt, and vibrobox. We also used: Astrophysics collaborations (astro-ph) [27], C. Elegans Neural Network (celegans) [30], and the Marvel Comics social network [7]. We used two large social networks (LiveJournal [26] and Twitter [21]) to evaluate our heuristics in a real system in Section 6.

We created synthetic datasets using popular generative models, preferential attachment (BA) [13], Watts-Strogatz

Name	$ V $	$ E $	Type	Source
3elt	4720	13,722	FEM	[5]
4elt	15606	45,878	FEM	[5]
vibrobox	12,328	165,250	FEM	[5]
celegans	297	2,148	Protein	[30]
astro-ph	18,772	396,160	Citation	[27]
Slashdot0811	77,360	504,230	Social	[22]
wiki-Vote	7,115	99,291	Social	[22]
Marvel	6,486	427,018	Social	[7]
web-ND	325,729	1,497,134	Web	[22]
BA	1,000	9,900	Synth.	[13]
BA	10,000	129,831	Synth.	[13]
BA	50,000	1,249,375	Synth.	[13]
RMAT	1,000	9,175	Synth.	[23]
RMAT	10,000	129,015	Synth.	[23]
RMAT	50,000	1,231,907	Synth.	[23]
WS	1,000	5,000	Synth.	[30]
WS	10,000	120,000	Synth.	[30]
WS	50,000	3,400,000	Synth.	[30]
PL	1,000	9,878	Synth.	[18]
PL	10,000	129,763	Synth.	[18]
PL	50,000	1,249,044	Synth.	[18]
LiveJournal	$4.6 \cdot 10^6$	$77.4 \cdot 10^6$	Social	[26]
Twitter	$41.7 \cdot 10^6$	$1.468 \cdot 10^9$	Social	[21]

Table 1: Graph datasets summary

(WS) [30], the RMAT generator [23], and a power-law graph generator with clustering (PL) [18]. Three of the synthetic datasets, BA, WS, and PL were created with the NetworkX python package. For each model, we created a degree distribution with average degree  $O(\log n)$  (average degree of 10 edges for 1,000 nodes, 13 for 10,000, and 25 for 50,000). This fully specifies the BA model. For WS and PL we used .1 as the rewiring probability. The RMAT datasets were created with the Python Web Graph Generator, a variant of the RMAT generator [2]. The RMAT or Kronecker parameters used by this implementation are [0.45,0.15;0.15,0.25].

The datasets were chosen to balance both size and variety. All are small enough so that we can find offline solutions with METIS so that our results are reproducible, while still big enough to capture the asymptotic behavior of these graph types. The collection captures a variety of real graphs, focusing on finite-element meshes (FEM) and power-law graphs. FEMs are used for scientific computing purposes to model simulations like the flow over a wing, while power-law (and other heavy-tailed) distributions capture nearly all ‘natural’ graphs, like the World Wide Web, social networks, and protein networks. In general, it is known that FEMs have good partitions because their edges are highly local, while natural graphs are more difficult to partition because they have high expansion and low diameter. The basic statistics about each graph, as well as its type and source are in Table 1.

### 4.2 Methodology

We examined all the combinations of datasets, heuristics and stream orders and ran each experiment 5 times on each combination. The *Random* ordering is a random permutation of the vertices, while *BFS* and *DFS* were created by sampling a random vertex to be the root of the BFS or DFS algorithm. Each of the heuristics was run on the same ordering. We ran each experiment on 2, 4, 8, and 16 partitions and fixed the imbalance such that no partition held more

than 5% more vertices than its share. The imbalance was chosen as a reasonable setting of this parameter in practice.

## 5. EVALUATION RESULTS

In all of the following figures, the y-axis has been scaled to zoom in on the data. The ordering of the heuristics in the figures is the one given in Table 2.

### 5.1 Upper and lower bounds

In order to evaluate the quality of our heuristics, we must establish good upper and lower bounds for the performance. A natural upper bound is the approach currently used in practice - hashing the node ID and mapping it to a partition. This approach completely ignores the edges so its expected performance is cutting a  $\frac{k-1}{k}$  fraction of edges for  $k$  partitions. This bound is marked by the upper black line in our figures. We expect Balanced and Hashing to always perform at this level, as well as Chunking on a random order.

The lower bound can be picked in many more ways. Finding an optimal lower bound is NP-hard, so we focus on more realistic approaches. We compare against a practical and fast approach, the partition produced by METIS v4.0.3. While METIS has no theoretical guarantees, it is widely respected and produces good cuts in practice, and is thus a good offline comparison for our empirical work. This METIS value is marked as the lower black line in our figures. Note that METIS is given significantly more information than the streaming heuristics, so we would not expect them to produce partitioning of the same quality. Any heuristic between these two lines is an improvement.

### 5.2 Performance on three graph types

We have included figures of the results for three of the graphs, a synthetic graph, a social network graph and a FEM with the goal of covering all three major types of graphs.

Figure 1 depicts the performance on the PowerLaw Clustered graph [18] of size 1,000 with 4 partitions. This is one of our synthetic graphs where the model is intended to capture power law graphs observed in nature. The lower bound provided by METIS is 58.9% of the edges cut, while the upper bound for 4 partitions is 75%. The first heuristic, Avoid Big, is worse than a random cut. Linear Deterministic Greedy and Balance Big both perform very well for all 3 stream orderings. These each had a best average performance of 61.7% and 63.2% of the edges cut respectively, corresponding to 82% and 73% of the possible gain in performance. This gain was calculated as the fraction of edges cut by the random minus the fraction cut by the heuristic, divided by the fraction cut by a random cut minus the fraction cut by METIS ( $\frac{\text{random} - \text{heuristic}}{\text{random} - \text{METIS}}$ ).

Figure 2 is our results for a social network, the Marvel Comics network [7], with 8 partitions. The Marvel network is synthetic, as it is the result of character interactions in books, but studies have shown it is similar to real social networks. The lower bound from METIS cuts only 32.2% of edges. The upper bound is  $7/8 = 87.5\%$  cut. The two heuristics at the upper bound level are Balanced and Hashing, with Chunking on the random order also performing poorly, as expected. Again, the best heuristic is Linear Deterministic Greedy, with 48%, 48.7% and 50.8% edges cut for the *BFS*, *DFS* and *Random* orderings respectively. This constitutes a gain of 71.3%, 70% and 66%.

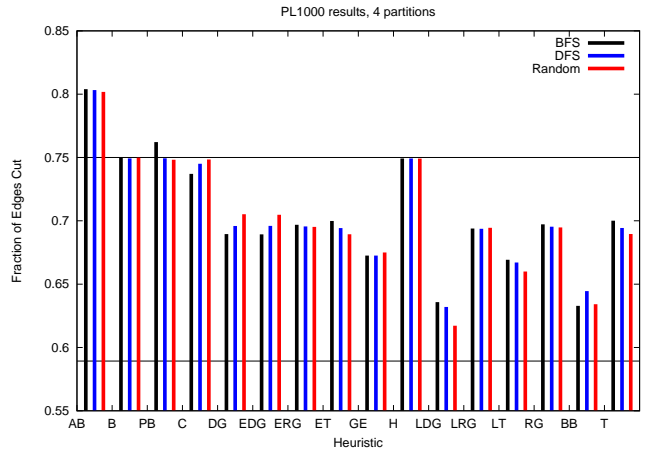


Figure 1: PL1000 results. The top line is the cost of a random cut and the bottom line is METIS. The best heuristic is Linear Deterministic Greedy. The figures are best viewed in color.

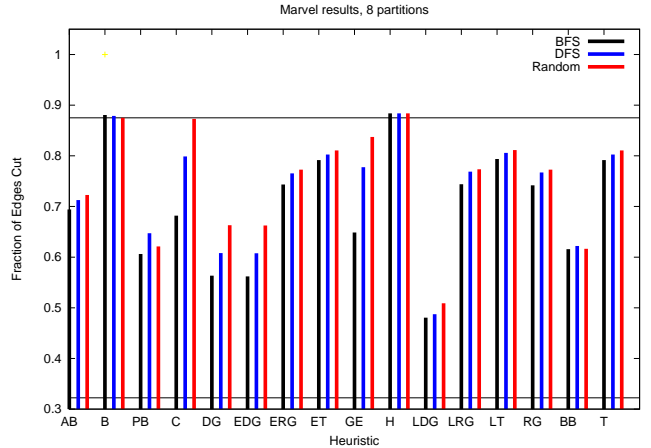


Figure 2: Marvel results. The top line is the cost of a random cut and the bottom line is METIS. The best heuristic is Linear Deterministic Greedy.

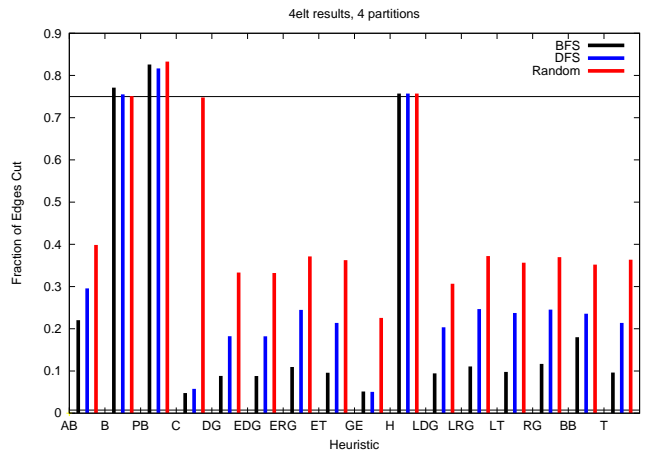


Figure 3: 4elt results. The top line is a random cut and the bottom line is METIS (0.7% edges cut).

Heuristic		<i>BFS</i>	<i>DFS</i>	<i>Random</i>
Avoid Big	AB	-27.3	-38.6	-46.4
Balanced	B	-1.5	-1.3	-0.2
Prefer Big	PB	-9.5	-18.6	-23.1
Chunking	C	37.6	35.7	0.7
Deterministic Greedy	DG	57.7	54.7	45.4
Exp. Det. Greedy	EDG	59.4	56.2	47.5
Exp. Rand. Greedy	ERG	45.6	45.6	38.8
Exp. Triangles	ET	50.7	49.3	41.6
Greedy EvoCut	GE	60.3	58.6	43.1
Hashing	H	-1.9	-2.1	-1.7
Linear Det. Greedy	LDG	76	73	75.3
Linear Rand. Greedy	LRG	46.4	44.9	39.1
Linear Triangles	LT	55.4	54.6	49.3
Randomized Greedy	RG	45.5	44.9	38.7
Balance Big	BB	67.8	68.5	63.3
Triangles	T	49.7	48.4	40.2

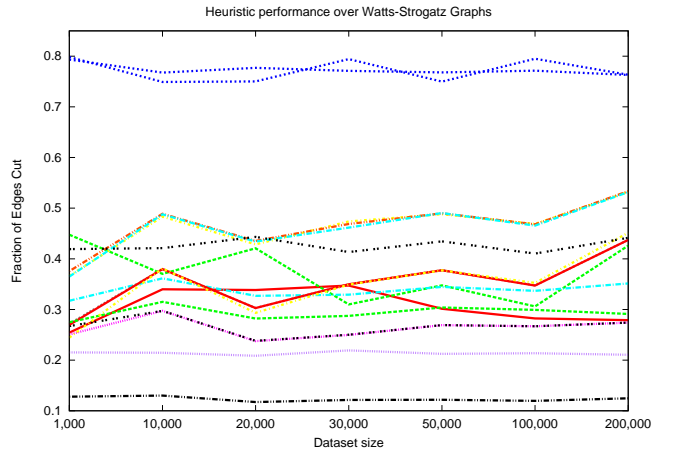
**Table 2: The average gain of each heuristic over all of our datasets and partitions sizes.**

Figure 3 contains the results for a FEM, 4elt [5], with 4 partitions. The change in graph structure gives us quite different results, not the least of which is that METIS now cuts only 0.7% of edges. The upper bound remains at 75%, providing a huge range for improvement. Surprisingly, Chunking performs extremely well here for the *BFS* and *DFS* orders, at 4.7% and 5.7% cut respectively. Translating into gain provides 94.7% and 93.3% of the optimal improvement. Chunking performs poorly on the *Random* order as expected. The other heuristic that performs well is Greedy EvoCut, obtaining 5.1% and 5% cuts for *BFS* and *DFS* respectively. Linear Deterministic Greedy obtains 9.4%, 20.3%, and 30.6% cuts for *BFS*, *DFS* and *Random* respectively. In fact, all of the heuristics beyond Balanced and Hashing are vast improvements. The *BFS* ordering is also a strict improvement for all approaches over the *DFS* and *Random* orderings.

### 5.3 Performance on all graphs: discussion

We present the gain in the performance of each heuristic in Table 2, averaged over all datasets from Table 1 (except for LiveJournal and Twitter) and all runs, for each ordering. The best heuristic is Linear Deterministic Greedy for all orderings, followed by Balance Big. Greedy EvoCut is also successful on the *BFS* and *DFS* orderings, but is computationally much more expensive than the other two approaches. Note that Balance Big is a combination of the Greedy and Balanced strategies, assigned based on node degree. There are universally bad heuristics, namely Prefer Big and Avoid Big. Both of these are significantly worse than Hashing.

We further restrict the results by type of graphs. As stated earlier, FEMs have good balanced edge cuts. For these types of graphs, no heuristic performed worse than the Hashing approach, and most did significantly better. For the *BFS* ordering, Linear Deterministic Greedy had an average 86.6% gain, with Deterministic Greedy closely behind at 84.2%. For the *DFS* ordering, the Greedy EvoCut approach performed best at 78.8%, with all 3 deterministic greedy approaches closely behind at 74.9% (exp), 74.8% (unweighted) and 75.8% (linear). Finally, as always, the *Random* ordering was the hardest, but Linear Deterministic Greedy was also the best with 63% improvement. No other method achieved more than 56%. The surprising result for FEMs is how well the Chunking heuristic performed: an 80% improvement for



**Figure 4: *BFS* with 4 partitions. Each line is a heuristics performance over 7 sizes of WS graph. The bottom line is METIS. The bottom purple line is Linear Deterministic Greedy. Best viewed in color.**

*BFS* and 72% for *DFS*. This is a huge improvement for such a simple heuristic, although it is due to the topology of the networks and the fact that *BFS* is used in partitioning algorithms to find good cuts. When given a *Random* ordering, Chunking had only a 0.2% average improvement.

The social networks results were more varied. Here, the Prefer Big and Avoid Big both have large negative improvements, meaning both should never be used for power law degree networks with high expansion. For all three orderings, Linear Deterministic Greedy was clearly the superior approach with 71% improvement for *BFS* and 70% for *DFS*. The second best performance was from both Exponential Deterministic Greedy and Deterministic Greedy at 60.5% for *BFS* and 52.9% for *DFS*. Finally, for a *Random* ordering, Linear Deterministic Greedy achieved a 64% improvement, with the other greedy approaches at only 42%.

Given that the Linear Deterministic Greedy algorithm performed so well, even compared with the other variants, one may ask why. At a high level, the penalty functions form a continuum. The unweighted version has a very strict cut-off - the penalty only applies when the partition is full and gives no indication that this restriction is approaching. The exponential penalty function has similar performance to the unweighted version because while the exponential function does not indicate that the partition is nearly full until it is very close to the limit. The linear weighting optimally balances the greedy choices with preferring less loaded partitions. Since  $1 - x \approx e^{-x}$  when  $0 < x < 1$ , the linear weighting can be seen as a normalized exponential weighting. This normalization term allows the penalty to take effect much earlier in the process and smooths the information by preventing the size of the partition from affecting the prediction. As this is a continuum, this parameter could be further fine-tuned for different types of graphs. Additionally, the implementation of the unweighted greedy algorithm in this paper breaks ties lexicographically. Breaking ties by load is equivalent to an indicator penalty function and its performance is very close to the linear penalty function.

### 5.4 Scalability in the graph size

All of our datasets discussed so far are tiny when compared with graphs used in practice. While the above results



are promising, it is important to understand whether the heuristics scale with the size of the graph. We used the synthetic datasets in order to control for the variance in different graphs. The key assumption is that using the same generative model with similar parameter settings will guarantee similar graph statistics while allowing the number of edges and nodes to vary. We began by looking at the results for the four generative models, BA, RMAT, WS, and PL. For each of these we had 3 data points: 1,000 vertices, 10,000 vertices, and 50,000 vertices. In order to get a better picture, we created additional graphs with 20,000, 30,000, 100,000 and 200,000 vertices. We will present only the results for the Watts-Strogatz graphs, but all other graphs exhibit quite similar results. Note that these results will scale to any size graph created by the same generative model because of the statistical properties of the generated graphs.

The labels in Figure 4 have been elided for clarity of the image. The bottom black line is METIS. This shows that our idea that the fraction of edges cut should scale with the size of the graph holds - it is approximately 12% for each graph. Next, there is clearly a best heuristic for this type of graph, the purple line. It corresponds to the Linear Deterministic Greedy heuristic. It has an average edge cut of 21% over all sizes of the graphs. Finally, all of the lines are approximately constant. The noise in the performance of each algorithm is due to the random nature of the orderings, and would decrease with further trials.

### 5.5 Scalability in the number of partitions

The other question is how the partitioning quality scales with the number of partitions. The fraction of edges cut must necessarily increase as we increase the number of partitions. Also, we are not trying to find an optimal number of partitions for the graph. As before, we only present data on one graph in Figure 5, the 50,000 node PowerLaw Clustered graph, but all graphs have similar characteristics. The heuristics performance closely tracks that of METIS.

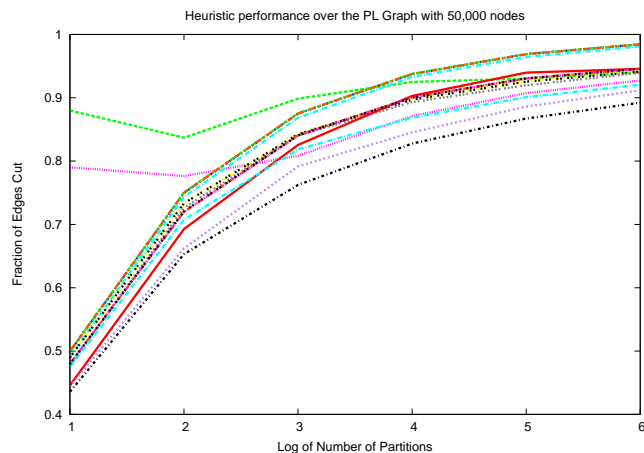


Figure 5: *BFS* with 2-64 partitions. Each line connects a heuristics performance over the 6 partition sizes. The bottom line is METIS. The bottom purple line is Linear Deterministic Greedy.

## 6. RESULTS ON A REAL SYSTEM

After evaluating the performance of the partitioning algorithms, we naturally ask whether the improvement in the

partitioning makes any measurable difference for real computations. To evaluate our partitioning scheme in a cluster application, we used an implementation of PageRank in Spark [32], a cluster computing framework for iterative applications. Spark provides the ability to keep the working set of the application (the graph topology and PageRank values) in memory across iterations, so the algorithm is primarily limited by communication between nodes. Other recent frameworks for large-scale graph processing, like Pregel [25] and GraphLab [24], also keep data in memory and are expected to exhibit similar performance characteristics.

There are many graph algorithms implemented for the Spark system, but we chose PageRank for two reasons. One is the popularity of this specific algorithm, and the other is its generality. PageRank is a specialized matrix multiplication, and many graph algorithms can be expressed similarly. Additionally, Spark has two implementations of PageRank: a naïve version that sends a message on the network for each edge, and a more sophisticated combiner version that aggregates all messages between each partition [31].

We used Linear Deterministic Greedy, as it performed best in our previous experiments. We tried both a vertex balanced version and an edge-balanced version. However, our datasets are social networks and follow a power-law degree distribution. For PageRank, the quantity that should be balanced is the number of edges in each partition as this controls the amount of computation performed in sparse matrix multiplication and we want this to be equal for all partitions. The existence of very high degree nodes means that some partitions contain many more edges than others, resulting in unbalanced computation times between the different cluster machines. We therefore modified Linear Deterministic Greedy to use the number of edges in a partition for the weight penalty. We used two datasets, LiveJournal [26] with 4.6 million nodes and 77.4 million edges, and Twitter [21] with 41.7 million nodes and 1.468 billion edges. While neither are Internet scale, they are both realistic for medium sized web systems and large enough to show the effects of reduced communication on a distributed computation.

### LiveJournal.

We used 100 partitions, with imbalance of 2% and the stream order provided by the authors of the dataset which is an unknown ordering. Linear Deterministic Greedy reduced the number of edges cut to 47,361,254 edges compared with 76,234,872 for Hashing. We ran 5 iterations of both versions of PageRank, and repeated this experiment 5 times. With the improved partitioning, naïve PageRank was 38.7% faster than the hashed partitioning version. The timing information, along with standard deviations, is summarized in Table 3. We used 10 “large” machines (7.5GB memory and 2 CPUS) on Amazon’s EC2. The combiner version with our partitioning was 28.8% faster than the hashed version. This reduction in computation time is obtained entirely by laying out the data in a slightly smarter way.

### Twitter.

We repeated the experiment for Twitter [21]. This graph is one of the largest publicly available datasets. Twitter was partitioned into 400 pieces with a maximum imbalance of 2%. Linear Deterministic Greedy cut 1.341 billion edges, while Hashing cut 1.464 billion. We used 50 “large” machines with 100 cores. The total computation time is much longer



	LJ Hash	LJ Streamed
Naïve PR Mean	296.2s	181.5s
Naïve PR STD	5.5 s	2.2 s
Combiner PR Mean	155.1 s	110.4 s
Combiner PR STD	1.5 s	0.8 s
	Twitter Hash	Twitter Streamed
Naïve PR Mean	1199.4 s	969.3 s
Naïve PR STD	81.2 s	16.9 s
Combiner PR Mean	599.4 s	486.8 s
Combiner PR STD	14.4 s	5.9 s

**Table 3: Timing data (mean and standard deviation) for 5 iterations of PageRank computation on Spark for LiveJournal and Twitter graphs, Hashing vs. Linear Deterministic Greedy.**

due to the increase in size. The naïve PageRank was 19.1% faster with our partitioning while the combiner version was 18.8% faster. For both graphs, there was additional time associated with loading the graph, about 200 seconds for Twitter and 80 seconds for LiveJournal, but this was not affected by the partitioning method.

These results show that with very little engineering effort, a simple preprocessing step that considers the graph edges can yield a large improvement in the running time. The best heuristic can be computed for each arriving node in time that is linear in the number of edges, given access to the distributed lookup table for the cluster and knowledge of the current loads of the machines. The improvement in running time is entirely due to the reduced network communication.

## 7. CONCLUSIONS AND FUTURE WORK

We have demonstrated that simple, one-pass streaming graph partitioning heuristics can dramatically improve the edge-cut in distributed graphs. Our best performing heuristic is the linear weighted variant of the greedy algorithm. This is a simple and effective preprocessing step for large graph computation systems, as the data must be loaded onto the cluster any way. One might need to perform a full graph partitioning once the graph has been fully loaded, however, as it will be re-partitioning an already partitioned graph, there will be less communication cost and it potentially may need to move fewer vertices, and will be faster. Using our approach as preprocessing step can only benefit any future computation while incurring only small cost.

There are several future directions for our work. First, there is the theoretical work. A framework should be developed for proving the performance of these heuristics. The main complications are the addition of the *BFS* and *DFS* stream orderings, and the fact that the offline optimal solution is NP-hard to compute. However, the best heuristic also performed well on the *Random* ordering, so it may be possible to prove a bound on its performance with additional assumptions like the graph is generated by a specific model.

The second direction is to address using parallel loaders. In reality, using a single machine to load terabytes of data will be far too slow. We expect the performance of running parallel loaders on independent portions of the graph stream to be similar to our experiments, requiring only access to the distributed lookup table.

## 8. REFERENCES

[1] <http://facebook.com/press/info.php?statistics>, Jan 2012.

[2] <http://pywebgraph.sourceforge.net>.

[3] <http://research.microsoft.com/ldg>, Jan 2012.

[4] <http://research.microsoft.com/trinity>, Jan 2012.

[5] <http://staffweb.cms.gre.ac.uk/~wc06/partition>.

[6] K. Jin Ahn and S. Guha. Graph sparsification in the semi-streaming model. *ICALP*, 2009.

[7] R. Alberich, J. Miro-Julia, and F. Rossello. Marvel universe looks almost like a real social network. *arXiv*, 2002.

[8] R. Andersen, F. Chung, and K. Lang. Local graph partitioning using pagerank vectors. In *FOCS*, 2006.

[9] R. Andersen and Y. Peres. Finding sparse cuts locally using evolving sets. In *STOC*, pages 235–244, 2009.

[10] K. Andreev and H. Racke. Balanced graph partitions. *Theory of Computing Systems*, 39:929–939, 2006.

[11] S. Arora, S. Rao, and U. Vazirani. Expander flows, geo-metric embeddings and graph partitioning. *J.ACM*, 2009.

[12] B. Bahmani, A. Chowdhury, and A. Goel. Fast incremental and personalized pagerank. *PVLDB*, 4(3):173–184, 2010.

[13] A-L Barabasi and R. Albert. Emergence of scaling in random networks. *Science*, 286(5439):509–512, 1999.

[14] S. Barnard. PMRSB: Parallel multilevel recursive spectral bisection. In *Supercomputing*, 1995.

[15] S. Bykov, A. Geller, G. Kliot, J. Larus, R. Pandya, and J. Thelin. Orleans: Cloud computing for everyone. In *ACM Symposium on Cloud Computing*, 2011.

[16] G. Even, J. Naor, S. Rao, and B. Schieber. Fast approximate graph partitioning algorithms. *SIAM J. Comput*, 28(6):2187–2214, 1999.

[17] B. Hendrickson and R. Leland. A multilevel algorithm for partitioning graphs. In *Supercomputing*, 1995.

[18] P. Holme and B. J. Kim. Growing scale-free networks with tunable clustering. *Phys. Rev. E*, 2002.

[19] G. Karypis and V. Kumar. Multilevel graph partitioning schemes. In *ICPP*, pages 113–122, 1995.

[20] J. Kelner and A. Levin. Spectral sparsification in the semi-streaming setting. *STACS*, pages 440–451, 2011.

[21] H. Kwak, C. Lee, H. Park, and S. Moon. What is Twitter, a social network or a news media? In *WWW*, 2010.

[22] J. Leskovec. <http://snap.stanford.edu/snap>, 2012.

[23] J. Leskovec, D. Chakrabarti, J. Kleinberg, C. Faloutsos, and Z. Ghahramani. Kronecker graphs: An approach to modeling networks. *JMLR*, 2010.

[24] Y. Low, J. Gonzalez, A. Kyrola, D. Bickson, C. Guestrin, and J. Hellerstein. GraphLab: A new framework for parallel machine learning. In *UAI*, 2010.

[25] G. Malewicz, M. Austern, A. Bik, J. Dehnert, I. Horn, N. Leiser, and G. Czajkowski. Pregel: a system for large-scale graph processing. *PODC*, 2009.

[26] A. Mislove, M. Marcon, K. Gummadi, P. Druschel, and B. Bhattacharjee. Measurement and analysis of online social networks. In *ACM/USENIX IMC*, 2007.

[27] M. E. J. Newman. The structure of scientific collaboration networks. *Natl Acad Sci*, 98:404–9, 2001.

[28] A. Das Sarma, S. Gollapudi, and R. Panigrahy. Sparse cut projections in graph streams. In *ESA*, pages 480–491, 2009.

[29] A. Das Sarma, S. Gollapudi, and R. Panigrahy. Estimating pagerank on graph streams. *J. ACM*, 58(3):13, 2011.

[30] D. J. Watts and S. H. Strogatz. Collective dynamics of ‘small-world’ networks. *Nature*, 393:440–442, 1998.

[31] M. Zaharia, M. Chowdhury, T. Das, A. Dave, J. Ma, M. McCauley, M.J. Franklin, S. Shenker, , and I. Stoica. Resilient distributed datasets: A fault-tolerant abstraction for in-memory cluster computing. *NSDI*, 2012.

[32] M. Zaharia, M. Chowdhury, M. J. Franklin, S. Shenker, and I. Stoica. Spark: Cluster Computing with Working Sets. In *HotCloud*, 2010.

[33] M. Zelke. Intractability of min- and max-cut in streaming graphs. *IPL*, 111(3):145 – 150, 2011.