

Article

# Distributed Balanced Partitioning via Linear Embedding <sup>†</sup>

Kevin Aydin <sup>‡</sup>, MohammadHossein Bateni <sup>\*,‡</sup> and Vahab Mirrokni <sup>‡</sup>

Google Research, 76 Ninth Ave, New York, NY 10011, USA

\* Correspondence: bateni@google.com

<sup>†</sup> This article is an extended version of our paper published in Proceedings of the Ninth ACM International Conference on Web Search and Data Mining, San Francisco, CA, USA, 22–25 February 2016.<sup>‡</sup> Authors contributed equally to this work.

Received: 18 July 2019; Accepted: 7 August 2019; Published: 10 August 2019



**Abstract:** Balanced partitioning is often a crucial first step in solving large-scale graph optimization problems, for example, in some cases, a big graph can be chopped into pieces that fit on one machine to be processed independently before stitching the results together, leading to certain suboptimality from the interaction among different pieces. In other cases, links between different parts may show up in the running time and/or network communications cost, hence the desire to have small cut size. We study a distributed balanced-partitioning problem where the goal is to partition the vertices of a given graph into  $k$  pieces so as to minimize the total cut size. Our algorithm is composed of a few steps that are easily implementable in distributed computation frameworks such as MapReduce. The algorithm first embeds nodes of the graph onto a line, and then processes nodes in a distributed manner guided by the *linear embedding* order. We examine various ways to find the first embedding, for example, via a hierarchical clustering or Hilbert curves. Then we apply four different techniques including local swaps, and minimum cuts on the boundaries of partitions, as well as contraction and dynamic programming. As our empirical study, we compare the above techniques with each other, and also to previous work in distributed graph algorithms, for example, a label-propagation method, FENNEL and Spinner. We report our results both on a private map graph and several public social networks, and show that our results beat previous distributed algorithms: For instance, compared to the label-propagation algorithm, we report an improvement of 15–25% in the cut value. We also observe that our algorithms admit scalable distributed implementation for any number of partitions. Finally, we explain three applications of this work at Google: (1) Balanced partitioning is used to route multi-term queries to different replicas in Google Search backend in a way that reduces the cache miss rates by  $\approx 0.5\%$ , which leads to a double-digit gain in throughput of production clusters. (2) Applied to the Google Maps Driving Directions, balanced partitioning minimizes the number of cross-*shard* queries with the goal of saving in CPU usage. This system achieves load balancing by dividing the *world graph* into several “shards”. Live experiments demonstrate an  $\approx 40\%$  drop in the number of cross-shard queries when compared to a standard geography-based method. (3) In a job scheduling problem for our data centers, we use balanced partitioning to evenly distribute the work while minimizing the amount of communication across geographically distant servers. In fact, the hierarchical nature of our solution goes well with the layering of data center servers, where certain machines are closer to each other and have faster links to one another.

**Keywords:** cut minimization; embedding to line; imbalance; local improvement; MapReduce; maps; partitioning; social networks

## 1. Introduction

Graph partitioning is a crucial first step in developing a tool for mining big graphs or solving large-scale optimization problems. In many applications, the partition sizes have to be balanced or almost balanced so that each part can be handled by a single machine to ensure the speedup of parallel computing over different parts. Such a balanced graph partitioning tool is applicable throughout scientific computation for dealing with distributed computations over massive data sets, and serves as a tradeoff between the local computation and total communication amongst machines.

In several applications, a graph serves as a substitute for the computational domain [1]. Each vertex denotes a piece of information and edges denote dependencies between information. Usually, the goal is to partition the vertices such that no part is too large, and the number of edges across parts is small. Such a partition implies that most of the work happens within a part and only minimal work or communication takes place between parts. In these applications, links between different parts may show up in the running time and/or network communications cost. Thus, a good partitioning minimizes the amount of communication during a distributed computation. In other applications, a big graph can be carefully chopped into parts that fit on one machine to be processed independently before finally stitching the individual results together, leading to certain suboptimality arising from the interaction between different parts. To formally capture these situations, we study a balanced partitioning problem where the goal is to partition the vertices of a given graph into  $k$  parts so as to minimize the total cut size. We emphasize that our main focus is on the size of the resulting cut and not the resource consumption of our algorithm, although it is fully distributable, uses linear space and runs for all the relevant experiments in reasonable amount of time (more or less comparable to the state of the art). Unfortunately, we cannot reveal the exact running times due to corporate restrictions. Nevertheless, we report in Section 6.4 relative running times of our algorithm on synthetic data of varying sizes to demonstrate its scalability.

This is a challenging problem that is computationally hard even for medium-size graphs [2] as it captures the graph bisection problem [3], hence all attempts at tackling it necessarily relies on heuristics. While the topic of large-scale balanced graph-partitioning has attracted significant attention in the literature [4–7], the large body of previous work study large-scale but non-distributed solutions to this problem. Several motivations necessitate the quest for a distributed algorithm for these problems: (i) first of all, huge graphs with hundreds of billions of edges that do not fit in memory are becoming increasingly common [4,5]; (ii) in many distributed graph-processing frameworks such as Pregel [8], Apache Giraph [9], and PEGASUS [10], we need to partition the graph into different pieces, since the underlying graph does not fit on a single machine, and therefore for the same reason, we need to partition the graph in a distributed manner; (iii) the ability to run a distributed algorithm using a common distributed computing framework such as MapReduce on a distributed computing platform consisting of commodity hardware makes an algorithm much more widely applicable in practice; and (iv) even if the graph fits in memory of a super-computer, implementing some algorithms requires superlinear memory which is not feasible on a single machine. The need for distributed algorithms has been observed by several practical and theoretical research papers [4,5].

For streaming and some distributed optimization models, Stanton [11] has shown that achieving formal approximation guarantees is information theoretically impossible, say, for a class of random graphs. Given the hardness of this problem, we explore several distributed heuristic algorithms. Our algorithm is composed of a few logical, simple steps, which can be implemented in various distributed computation frameworks such as MapReduce [12].

### 1.1. Motivation

Following the publication of the extended abstract of this work at WSDM 2016 [13], we have deployed this algorithm in several applications from various domains (such as Maps, Search and Infrastructure). However, we only mention two motivating applications in this section, in order to

keep the reader focused. The interested reader is encouraged to visit Section 7 for more information about applications.

#### 1.1.1. Google Maps Driving Directions

As one of the applications of our balanced-partitioning algorithm, we study a balanced-partitioning problem in the Google Maps Driving Directions. This system computes the optimal driving route for any given *source-destination query* (pairs of locations on Google Maps). A plausible approach to answer such pairwise queries in this large-scale graph application is to obtain a suitable partitioning of the graph beforehand so that each server handles its own section of the graph, called a *shard*. In this application, ideally we seek to minimize the number of cross-shard queries, for which the source and destination are not in the same shard, that is, the total cross-shard traffic is as low as possible, since handling cross-shard queries is much more costly and entails extra communication and coordination across multiple servers. In order to reduce the cross-shard query traffic we want to minimize the cut size between the partitions, with the expectation that less directions will be produced where the source and the destination are in different partitions. Reducing the cut size in practice results in denser regions in the same partition, therefore most queries will be contained within one shard. We confirm this observation with extensive studies using historical query data and live experiments (see Section 7.1 for more details).

Conveniently, a linear embedding provides a very simple way of specifying partitions by merely providing two numbers for each shard's boundaries. Therefore it is straightforward and fast to identify the shard a given query point belongs to. As we will discuss later, our study of linear embedding of graphs while minimizing the cut size is partly motivated by these reasons.

#### 1.1.2. Treatment and Control Groups for Randomized Experiments

Another application where balanced partitioning has historically been found very effective is in the context of causal inference in randomized experiments [14,15]. In most settings there is an underlying interference between experiment units (for example, advertisers). Interference is often modeled by a graph in which nodes are treatment units and there is an edge between two nodes if their potential outcomes depend on each others' treatment assignment. The graph is then partitioned into clusters of approximately equal size that minimize the total edge weights crossing the clusters. Finding such clusters is equivalent to the balanced-partitioning problem, which is the focus of our study. Once the clusters are generated, a set of clusters can be chosen at random to set up the treatment portion of the experiment.

### 1.2. Our Contributions

Our contributions in this paper can be divided into five categories.

1. We introduce a multi-stage distributed optimization framework for balanced partitioning based on first embedding the graph into a line and then optimizing the clusters guided by the order of nodes on the line. (Embedding into a line is also essential in graph compression; see, e.g., References [16,17]). This framework is not only suitable for distributed partitioning of a large-scale graph, but also it is directly applicable to balanced-partitioning applications like those in Google Maps Driving Directions described above. More specifically, we develop a general technique that first embeds nodes of the graph onto a line, and then post-processes nodes in a distributed manner guided by the order of our linear embedding.
2. As for the initialization stage, we examine several ways to find the first embedding: (i) by naïvely using a random ordering, (ii) for map graphs by the Hilbert-curve embedding, and (iii) for general graphs by applying hierarchical clustering on an edge-weighted graph where the edge weights are based either on the number of common neighbors of the nodes in the graph or on the inverse of the distance of the corresponding nodes on the map. Whereas the Hilbert-curve

embedding and our hierarchical clustering based on node distances may only be applied to maps graphs, the hierarchical clustering based on the number of common neighbors is applicable to all graphs. While all these methods prove useful, to our surprise, the latter (most general) initial embedding technique is the best initialization technique even for maps. We later explain that all these methods have very efficient and scalable distributed implementations. See the discussion in Section 3.3 regarding the efficiency of the hierarchical clustering, and consult Reference [18] regarding the challenge to construct the similarity metric. (It might be worth studying spectral embedding methods or standard embedding techniques into  $\ell_1$ , but we did not try them due to not having a scalable distributed implementation, and also since the hierarchical clustering-based embedding worked pretty well in practice. We leave this for future research.)

3. As a next step, we apply four methods to postprocess the initial ordering and produce an improved cut. The methods include a metric ordering optimization method, a local improvement method based on random swaps, another based on computing minimum cuts in the boundaries of partitions, and finally a technique based on contracting the minimum-cut-based clusters, and applying dynamic programming to compute optimal cluster boundaries. Our final algorithm (called Combination) combines the best of our initialization methods, and iterates on the various postprocessing methods until it converges. The resulting algorithm is quite scalable: it runs smoothly on graphs with hundreds of millions of nodes and billions of edges.
4. In our empirical study, we compare the above techniques with each other, and more importantly with previous work. In particular, we compare our results to prior work in (distributed) balanced graph-partitioning, including a label propagation-based algorithm [4], FENNEL [5], Spinner [19] and METIS [20]. Some works such as Reference [21] are indirectly compared to, and others such as Reference [22] did not report cut sizes. We relied for these comparisons on available cut results for a host of large public graphs. We report our results on both a large private map graph, and also on several public social networks studied in previous work [4,5,19,22].
  - First, we show that our distributed algorithm consistently beats the label-propagation algorithm by Ugander and Backstrom [4] (on LiveJournal) by a reasonable margin for all values of  $k$ . For example, for  $k = 20$ , we improve the cut by 25% (from 37% to 27.5% of total edge weight), and for  $k = 100$ , we improve the cut by 15% (from 49% to 41.5%). The clustering outputs can be found in Reference [23].
  - In addition, for  $k > 2$  partitions, we show that our algorithm beats METIS and FENNEL (on Twitter) by a reasonable factor. For  $k = 2$ , the results that we obtain beats the output of METIS but it is slightly inferior to the result reported by FENNEL.
  - For both graphs, our results are consistently superior to that of Spinner [19] with a wide margin.
  - We also note that our algorithms admit scalable distributed implementation for small or large number of partitions. More specifically, changing  $k$  from 2 to tens of thousands does not change the running time significantly.
  - As for comparing various initialization methods, we observe that while geographic Hilbert-curve techniques outperform the random ordering, the methods based on the hierarchical clustering using the number of common neighbors as the similarity measure between nodes outperform the geography-based initial embeddings, even for map graphs. Note that the computation for the number of common neighbors is needed only once to obtain the initial graph, and the new edge weights in the subsequent rounds are computed via aggregation. Further note that approximate triangle counting can be done efficiently in MapReduce; see, for example, Reference [18].
  - As for comparing other techniques, we observe the random swap techniques are effective on Twitter, and the minimum-cut-based or contract and dynamic program techniques are very

effective for the map graphs. Overall, we realize that these techniques complement each other, and combining them in a Combination algorithm is the most effective method.

5. As mentioned earlier, we apply our results to the Google Maps Driving Directions application, and deploy two linear-embedding-based algorithms on the world's map graph. We first examine the best imbalance factor for our cut-optimization technique, and observe that we can reduce 21% of cross-shard queries by increasing the imbalance factor from 0% to 10%. The two methods that we examined via live experiments were (i) a baseline approach based on the Hilbert-curve embedding, and (ii) one method based on applying our cut-optimization post-processing techniques. In live experiments on the real traffic, we observe the number of multi-shard queries from our cut-optimization techniques is 40% less compared to the baseline Hilbert embedding technique. This, in turn, results in less CPU usage in response to queries. (The exact decrease in CPU usage depends on the underlying serving infrastructure which is not our focus and is not revealed due to company policies. We note that there are several other algorithmic techniques and system tricks that are involved in setting up the distributed serving infrastructure for Google Maps driving directions. This system is handled by the Maps engineering team, and is not discussed here as it is not the focus of this paper.)

### 1.3. Other Related Work

Balanced partitioning is a challenging problem to approximate within a constant factor [2,24,25] or to approximate within any factor in several distributed or streaming models [11]. As for heuristic algorithms for this problem in the distributed or streaming models, a number of recent papers have been written about this topic [4,5,21]. Our algorithms are different from the ones studied in these papers. The most similar related work are the label propagation-based methods of Ugander and Backstrom [4] and Martella et al. (Spinner) [19] which develop a scalable distributed algorithm for balanced partitioning. Our random-swap technique is similar in spirit to the label-propagation algorithm studied in Reference [4], however, we also examine three other methods as a postprocessing stage and find out that these methods work well in combination. Moreover, Reference [4] studied two different methods for their initialization, a random initialization, and a geographic initialization. We also examine random ordering, and a Hilbert-curve ordering which is similar to the geographic initialization in Reference [4]. However, we examine two other initialization techniques and observe that even for map-based geographic graphs, the initialization methods based on hierarchical clustering outperform the geography-based initial ordering. Overall, we compare our algorithm directly on a LiveJournal public graph (the only public graph reported in Reference [4]), and improve the cut values achieved in References [4,19] by a large margin for all values of  $k$ . In addition, algorithms developed in References [5,21] are suitable for the streaming model but one can implement variants of those algorithms in a distributed manner. We also compare our algorithm directly to the numbers reported on the large-scale Twitter graph by FENNEL [5], and show that our algorithm compares favorably with the FENNEL output, hence indirectly comparing against Reference [21] because of the comparison provided in Reference [5].

Motivated by a variety of big data applications, distributed clustering has attracted significant attention over the literature, both for applications with no strict size constraints on the clusters [26–28] or with explicit or implicit size constraints [29]. A main difference between such balanced-clustering problems and the balanced graph-partitioning problems considered here is that in the graph-partitioning problems a main objective function is to minimize the cut size, whereas in those clustering problems, the main goal is minimize the maximum or average distance of nodes to the centers of their clusters.

## 2. Preliminaries

For an integer  $n$ , we define  $[n] = \{1, 2, \dots, n\}$ . We also slightly abuse the notation and, for a function  $f : A \mapsto \mathbb{R}$ , define  $f(A') = \sum_{a \in A'} f(a)$  if  $A' \subseteq A$ . For a function  $f : S \mapsto T$ , we let  $f^{-1}(t)$



for  $t \in T$  denote the set of elements in  $S$  mapped to  $t$  via  $f$ ; more precisely,  $f^{-1}(t) = \{s \mid f(s) = t\}$ . For a permutation  $\pi$ , let  $\pi_i$  denote, for  $i \in [n]$ , the element in position  $i$  of  $\pi$ . Moreover, let  $\pi(i \rightarrow j) = \{\pi_i, \pi_{i+1}, \dots, \pi_j\}$  for  $i \in [n], j \in [n] \cup \{0\}$ . Note that, in particular,  $\pi(i \rightarrow j) = \emptyset$  if  $j < i$ .

### 2.1. Problem Definition

Given is a graph  $G(V, E)$  of  $n$  vertices with nonnegative edge lengths and nonnegative node weights, as well as an integer  $k$ , and a real number  $\alpha \geq 0$ . Let us denote the edge lengths by  $d : E \mapsto \mathbb{R}$ , and the vertex weights by  $w : V \mapsto \mathbb{R}$ . A partition of vertices of  $G$  into  $k$  parts  $\{V_i : i \in [k]\}$  is said to be  $\alpha$ -balanced if and only if

$$(1 - \alpha) \frac{w(V)}{k} \leq w(V_i) \leq (1 + \alpha) \frac{w(V)}{k}.$$

In particular, a zero-balanced (or fully balanced) partition is one where all partitions have the same weight. The cut length of the partition is the total sum of all edges whose endpoints fall in different parts:

$$\sum_i \sum_{j < i} \sum_{u \in V_i} \sum_{v \in V_j: (u,v) \in E(G)} w(u, v).$$

Our goal is to find an  $\alpha$ -balanced partition whose cut size is (approximately) minimized.

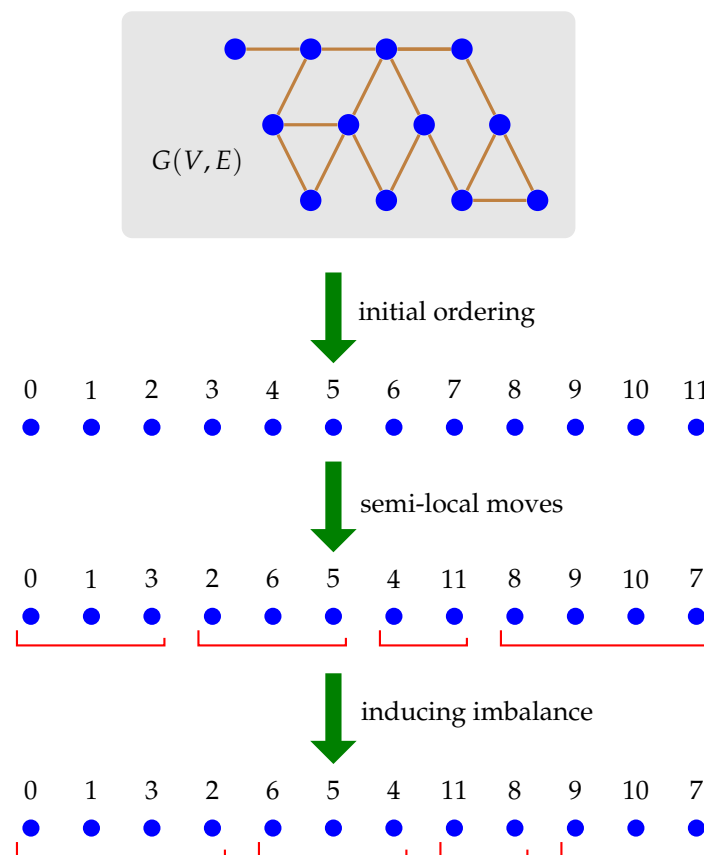
The problem is NP-hard as the case of  $\alpha = 0, k = 2$  is equivalent to the minimum bisection problem [3]. For arbitrary  $k$ , we get the minimum balanced  $k$ -cut problem, that is known to be inapproximable within any finite factor [2]; the best known approximation factor for it is  $O(\log^{1.5} n)$  (if  $\alpha > 0$  is constant).

### 2.2. Our Algorithm

Our algorithm consists of three main parts (depicted in Figure 1):

1. We first find a suitable mapping of the vertices to a line. This gives us an ordering of the vertices that presumably places (most) neighbors close to each other, therefore somewhat reduces the minimum-cut partitioning problem to an almost local optimization one.
2. We next attempt to improve the ordering mainly by swapping vertices in a semilocal manner. These moves are done so as to improve certain metrics (perhaps, the cut size of a fully balanced partition).
3. Finally, we use local postprocessing optimization in the “split windows” (i.e., a small interval around the equal-size partitions cut points taking into account permissible imbalance) to improve the partition’s cut size.

Note that one implementation may use only some of the above, and clearly any implementation will pick one method to perform each task above and mix them to get a final almost balanced partition. Indeed, we report some of our results (specially those in comparison to the previous work) based on *Combination*, which uses *AffCommNeigh* to get the initial ordering and then iteratively applies techniques from the second and third stages above (e.g., *Metric*, *Swap*, *MinCut*) until the result converges. (The convergence in our experiments happens in two or three rounds).



**Figure 1.** A natural flow for our algorithm (Algorithm 1). The vertices of  $G$  are first embedded on a line, then the embedding is improved via swaps and finally imbalanced clusters are created according to the ordering of vertices. As we explain later on, iterating on these steps yields better results.

---

**Algorithm 1** Combination( $G, k, \alpha$ )

---

**Input:** Graph  $G$ , number of parts  $k$ , imbalance parameter  $\alpha$

**Output:** A partition of  $V(G)$  into  $k$  parts

```

1: for all  $(u, v) \in E(G)$  do
2:    $w(u, v) \leftarrow$  number of common neighbors of  $u, v$ 
3:  $\pi \leftarrow \text{AffinityOrdering}(G, w)$ 
4: for  $i = 0$  to  $k$  do
5:    $q(i) \leftarrow \lfloor \frac{in}{k} \rfloor$  {fully balanced split points}
6: repeat
7:    $(\pi', q') \leftarrow (\pi, q)$ 
8:   Use any (semilocal or imbalance-inducing) postprocessing technique on  $(G, \pi', q')$  to obtain  $(\pi, q)$ 
9: until  $\pi = \pi'$  and  $q = q'$ 
10: for all  $i \in [k]$  do
11:    $V_i \leftarrow \pi(q(i) + 1 \rightarrow q(i + 1))$ 
12: return  $\{V_1, V_2, \dots, V_k\}$ 

```

---

### 3. Initial Mapping to Line

#### 3.1. Random Mapping

The easiest method to produce an ordering for the vertices is to randomly permute them. Though very fast, this method does not seem to lead to much progress towards our goal of finding

good cuts. In particular, if we then turn this ordering into a partition by cutting contiguous pieces of equal size, we end up with a random partition of the input (into equal parts) which is almost surely a bad cut: a standard probabilistic argument shows that the cut has expected ratio  $1 - \frac{1}{k}$ . Nevertheless, the next stage of the algorithm (i.e., the semilocal optimization by swapping) can generate a relatively good ordering with this naïve starting point.

### 3.2. Hilbert Curve Mapping

For certain graphs, geographic/geometric information is available for each vertex. It is fairly easy, then, to construct an ordering using a space-filling curve—prime examples are Peano, Morton and Hilbert curves but we focus on the latter in this work. These methods are known to capture proximity well: nodes that are close in the space are expected to be placed nearby on the line.

There has been extensive study of the applicability of these methods in solving large-scale optimization problems in parallel; see, for example, References [30,31].

Not only can this algorithm be used on its own without any cut guarantees but it can also be employed to break a big instance down into smaller ones that we can afford to run more intensive computations on. Both applications were known previously.

The previous works do not offer any theoretical guarantees on the quality of the cut generated from Hilbert curves. However, certain assumptions on the distributions of edge lengths and node positions let us bound the resulting cut ratio and show that it is significantly less than the result of random ordering [32–34]. This is observed in our experiments, too.

### 3.3. Affinity-Based Mapping

One drawback of the Hilbert-curve cover (even when geographic coordinates are available) is that it ignores the actual edges in the graph. For an illustration, consider using the Hilbert curve to map certain points in an archipalego (or just a small number of islands or peninsulas). The Hilbert curve, unaware of the connectivities, traverses the vertices in a semirandom order, therefore, it may jump from island to island without covering the entire island first.

To address this issue, we use an agglomerative hierarchical clustering method, usually called average-linkage clustering. We call the method “affinity clustering” since it takes into account the affinity of vertices. Informally, every node starts in a singleton cluster of its own and then in several stages we group vertices that are closely connected, hence building a tree of these connections. More precisely, at every stage, each node turns on the connection to its closest neighbor, after which the connected components form clusters. The similarities used in the next level between constructed clusters are computed via some function of the similarities between the elements forming the clusters; in particular, we take the average function for this purpose. Typically attributed to Reference [35], this is reminiscent of a parallel version of Borůvka’s algorithm for minimum spanning tree where the min operation is replaced by average [36].

The final ordering is produced by sorting the vertex labels produced as follows. Let the label for each vertex be the concatenation of vertex ID strings from the root to the corresponding leaf. Sorting the constructed labels places the vertices under each branch in a contiguous piece on the line. The same guarantee holds recursively, hence the (edge) proximity is preserved well. A pseudocode for this procedure is given in Algorithm 2: Line 15 uses notation  $a\#b$  to denote the concatenation of labels for  $a$  and  $b$ . Notice that this pseudocode explains a sequential algorithm, however, the procedure can be efficiently implemented in a few rounds of MapReduce. In fact, our algorithm is very similar to that in Reference [37] but uses Distributed Hash Tables as explained in Reference [38]. In fact, as mentioned therein, an implementation of the connected-components algorithm within the same framework and using the same techniques provides a 20–40 times improvement in running time compared to best previously known algorithms.



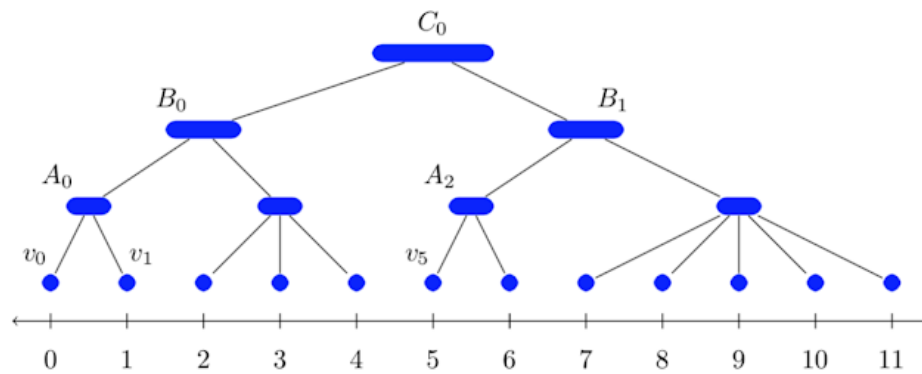
**Algorithm 2** AffinityOrdering( $G, w$ )**Input:** Graph  $G$ , (partial) similarity function  $w$ **Output:** A permutation of vertices

```

1:  $\mathcal{C}^0 \leftarrow \{\{v\} \mid v \in V(G)\}$ 
2: for  $v \in V(G)$  do
3:    $\ell(v) \leftarrow v$ 
4:  $i \leftarrow 0$ 
5: repeat
6:   for  $s \in \mathcal{C}^i$  do
7:      $p(s) \leftarrow \arg \max_{t \in \mathcal{C}^i} w(s, t)$ 
8:    $\mathcal{C}^{i+1} \leftarrow \emptyset$ 
9:   for  $t \in \{p(s) \mid s \in \mathcal{C}^i\}$  do
10:     $S \leftarrow p^{-1}(t)$ 
11:     $S' \leftarrow \bigcup_{s \in S} s$ 
12:     $\mathcal{C}^{i+1} \leftarrow \mathcal{C}^{i+1} \cup S'$ 
13:     $r \leftarrow \arg \min_{r' \in S'} s$ 
14:    for  $s' \in S'$  do
15:       $\ell(s') \leftarrow r \# \ell(s')$ 
16:    $i \leftarrow i + 1$ 
17: until  $|\mathcal{C}^i| = |\mathcal{C}^{i-1}|$ 
18:  $\pi \leftarrow$  permutation of vertices sorted according to  $\ell$ 
19: return  $\pi$ 

```

Intuitively, we expect vertices connected in lower levels (farther from the root) to be closer than those connected in later stages. (See Figure 2 for illustration). In particular, we observed that a graph with several connected components gives an advantage to affinity tree ordering over the ordering produced by Hilbert curve.



**Figure 2.** Illustration of Affinity tree ordering. The nodes at the bottom row correspond to the vertices of the graph and the nodes at the higher levels are formed by merging nodes from lower levels. The vertices are finally sorted in the picture according to their labels. For instance, the label for  $v_0$  is  $C_0 \# B_0 \# A_0 \# v_0$  whereas  $v_5$  is labeled  $C_0 \# B_1 \# A_2 \# v_5$ .

For this approach to work, we require meaningful distances/similarities between vertices of the graph. Conveniently it does not matter whether we have access to distances or similarity values. However, though theoretically possible to run the affinity clustering algorithm on an unweighted graph with distances 1 and  $\infty$  for edges and non-edges respectively, this leads to a lot of arbitrary

tie-breaks that makes the result irrelevant. Practically the worse outcome is a highly unbalanced hierarchical partitioning because arbitrary tie-breaks favor the “first” cluster.

Fortunately, there are standard ways to impose a metric on an unweighted graph: for example, common-neighbors ratio or personalized page rank. We focus on the former metric and compute for every pair of neighbors the ratio of the number of their common neighbors over the total number of their distinct neighbors. Computing the number of common neighbors can be implemented using standard techniques in MapReduce, however, one can use sampling when dealing with high-degree vertices to improve the running time and memory footprint (leading to an approximate result). As alluded earlier in the text, there are very fast distributed algorithms for computing this metric; see, for example, Reference [18].

#### 4. Improve Ordering Using SemiLocal Moves

Given an initial ordering of vertices for example by AffinityOrdering, we can further improve the cut size by applying semilocal moves. The motivation here is that these simple techniques provide us with highly parallelizable algorithms that can be iteratively applied (e.g., by a MapReduce pipeline) so the result converges to high-quality partitions. In the following we discuss two such approaches.

##### 4.1. Minimum Linear Arrangement

Minimum Linear Arrangement (MinLA) is a well-studied NP-hard optimization problem [3], which admits  $O(\log \log n)$  approximation on planar metrics [39] and  $O(\sqrt{\log n} \log \log n)$  approximation on general graph metrics [40,41]. Given an undirected graph  $G = (V, E)$  (with edge weights  $w_{uv}$ ) we seek a one-to-one function  $\phi : V \rightarrow \{0, 1, \dots, n-1\}$  that minimizes

$$\sum_{(u,v) \in E} |\phi(v) - \phi(u)| w_{uv}. \quad (1)$$

Optimizing the sum above results in a linear ordering of vertices along a line such that neighbor nodes are near each other. The end result is that dense regions of the graph will be isolated into clusters, which makes it easier to identify cut boundaries on the line.

We have implemented a very simple distributed algorithm (in MapReduce) that directly optimizes MinLA. The algorithm simply iterates on the following MapReduce phases until it converges (or runs up to a specified number of steps):

- **MR Phase 1:** Each node computes its optimal rank as the weighted median of its neighbors and outputs its new rank. Note that, all else fixed, moving a node to the weighted median of its neighbors—a standard computation in MapReduce—optimizes Equation (1).
- **MR Phase 2:** Assigns final ranks to each node (i.e., handles duplicate resolution, for example, by simple ID-based ordering).

Note that since Phase 1 is done in parallel for each node and independently of other moves, there will be side effects, hence it is an optimistic algorithm with the hope that it will converge to a stable state after few rounds. Our experimental results show that in practice this algorithm indeed converges for all the graph types that we tried, although the number of rounds depends on the underlying graph.

##### 4.2. Rank Swap

Given an existing linear ordering of vertices in a graph we can further improve the cut size (of the fully balanced partition based on it) via semilocal swaps. Note that unlike MinLA heuristic discussed in Section 4.1, this process depends on the pre-selected cut boundaries, that is, the number of final partitions  $k$ . Notice that one expects that the semilocal swap operations will be effective once a good initial ordering provided by either the Affinity-based mapping and/or MinLA. Indeed our experimental results show that this procedure is extremely effective and produces cut sizes better than the competition on some public social graphs we tried.

RankSwap can be implemented in a straightforward way on a distributed (for example, MapReduce) framework as outlined by Algorithm 3.

---

**Algorithm 3** RankSwap( $G, k, r$ )
 

---

**Input:** Graph  $G$ , number of partitions  $k$ , number of intervals per partition  $r$

**Output:** A partition of  $V(G)$  into  $k$  parts

– **Controller:**

- 1: Pair nearby partitions
- 2: Split each partition into  $r$  intervals  $I_0, \dots, I_{r-1}$  and randomly pair intervals between paired partitions

– **Map**  $\langle I_i; I_j \rangle$ :

- 3: **repeat**
  - 4: Pick node pair  $u \in I_i, v \in I_j$  with the best cut improvement
  - 5: Swap  $u$  and  $v$
  - 6: **until** no pair improves the cut
  - 7: **for all**  $u \in I_i \cup I_j$  **do**
  - 8: Emit  $\langle u; r_u \rangle$ .  $\{r_u$  is the new rank of  $u\}$
  - **Reduce**  $\langle u; r_u \rangle$ :
  - 9: Emit  $\langle u; r_u \rangle$ .  $\{\text{Identity Reducer}\}$
- 

Subdividing each partition into intervals is important, especially for small number of partitions  $k$ , to achieve parallelism and allowing possibly more time- and memory-consuming swap operations. Pairing the intervals between two partitions can be done in various ways. One simple example is random pairing. Our experiments showed that this is almost as effective as more complicated ones, hence we do not discuss the other methods in detail here.

The swap operation between two intervals (handled by each Mapper task) can be done in various ways. (The version described in Algorithm 3 is Method 3 below).

Each approach below starts with computing the cut size reduction as a result of moving a node from one interval to the other. Once these values are computed for each node in the two intervals, the following are the alternatives we tried (from the simplest to more complicated):

- **Method 1:** Sort the nodes in each interval in descending order by their cut size reduction. Then simply do a pairwise swap between entries at  $i$ -th place in each interval provided that the swap results in a combined reduction. This method is very simple and fast, however, it may have side effects within the interval (and also outside) since it does not account for those.
- **Method 2:** This is a modified version of Method 1 with the addition that after each swap is performed, we also update the cut reduction values of other nodes. This method does only one pass from top to bottom and stops when there is no further improvement.
- **Method 3:** In this method we iterate on Method 2 until we converge to a stable state. In this way we reach a local optimum between two intervals. The other improvement is that instead of pairing entries at  $i$ -th position, we find the best pair in the second interval and swap with that. This is because the top entries in each interval can be neighbors of each other and therefore they may not be the best pair to swap.

Our extensive studies with these methods confirmed that Method 3 gives the best results in terms of quality at the expense of an acceptable additional cost. It runs fast in practice due to the relatively small sizes of the intervals because of the parallelism we achieve. We can also afford an additional cost in the swap method because each interval pair is handled by a separate Mapper task.

## 5. Imbalance-Inducing Postprocessing

We have now an ordering of the vertices that presumably places similar vertices close to each other. Instead of cutting the sequence at equidistant positions to obtain a fully balanced partitioning, we describe in this section several postprocessing techniques that allow us to take advantage of the permissible imbalance and produce better cuts.

First we sketch how the problem can be solved optimally if the number of vertices is not too large. This approach is based on dynamic programming and requires the entire graph to live in memory, which is a barrier for using this technique even if the running time were almost linear. In fact, the procedure proposed here runs in time  $O(n^3 \log k)$  but the running time can be improved to almost linear. Despite the memory requirement, this algorithm is usable with the caveat that one first needs to group together blocks of contiguous vertices and contract them into supernodes, such that the graph of supernodes fits in memory and is small enough for the algorithm to handle. The fact that nearby vertices are similar to each other implies that this sort of contraction should not hurt the subsequent optimization significantly. Indeed trying two different block sizes (resulting in 1000 and 5000 blocks) produced similar cuts in our experiments.

The other two approaches work on the concept of “windows” and perform only local optimizations within each. Let  $\pi$  denote the permutation of vertices at the start of the local optimization and let  $\pi_i$  be the  $i$ -th vertex in the permutation for  $1 \leq i \leq n$ . Given that we want to get  $k$  parts of (almost) the same size, we consider the ideal split points  $q_j = \lfloor \frac{jn}{k} \rfloor$  for  $1 \leq j < k$ . Though there are only  $k - 1$  real split points, we define two dummy split points at either end to make the notation cleaner:  $q_0 = 0$  and  $q_k = n$ . Then, the ideal partitions, in terms of size, will be formed of  $\pi(q_j + 1 \rightarrow q_{j+1})$  for different  $0 \leq j < k$ . Given permissible imbalance  $\alpha > 0$ , a window is defined around each (real) split point allowing for  $\alpha/2$  imbalance on either side. More precisely,  $W_j = \pi(\pi_{q_j - \alpha'} \rightarrow \pi_{q_j + \alpha'})$  for  $\alpha' = \lceil \frac{\alpha n}{2k} \rceil$ . Essentially we are free to shift the boundaries of partitions within each window without violating the balance requirements. We propose two methods to take advantage of this opportunity. One finds the optimal boundaries within windows, while the other also permutes the vertices in each window arbitrarily to find a better solution.

### 5.1. Dynamic Program to Optimize Cuts

The partitioning problem becomes more tractable if we fix an ordering  $\pi$  of vertices and insist on each partition consisting of consecutive vertices (perhaps, with wrap-around). Let us for now assume that no wrap-around is permissible; that is, each part corresponds to a contiguous subset of vertices on  $\pi$ .

In the dynamic programming (DP) framework, instead of solving the given instance, we solve several instances (all based on the input and closely related to one another) and do so in a carefully chosen sequence such that the instances already solved make it easier to quickly solve newer instances, culminating in the solution to the real instance.

Let  $A_{i,j,q}$  for  $1 \leq i, j \leq n, 1 \leq q \leq k, j \geq i - 1$  denote the smallest cut size achievable for the subgraph induced by vertices  $\pi(i \rightarrow j)$  if exactly  $q$  partitions thereof are desired. Once all these entries are computed,  $A_{1,n,k}$  yields the solution to the original problem. As is customary in dynamic programs, we only discuss how to find the “cost”—here, the cut size—of the optimal solution; finding the actual solution—here, the partition—can be achieved in a straight-forward manner by adding certain information to the DP table.

We fill in the table  $A_{i,j,q}$  in the order of increasing  $q$ . For the case of  $q = 1$ , we clearly have  $A_{i,j,q} = 0$  if and only if

$$w(\pi(i \rightarrow j)) \leq (1 + \alpha) \frac{w(V)}{k},$$

that is, vertices  $\pi(i \rightarrow j)$  can be placed in one part without violating the balanced property. Otherwise, it is defined as infinity.

We use a recursive formula to compute  $A_{i,j,q}$  for  $q > 1$ . The formula depends on  $A_{i',j',q'}$  entries where  $q' < q$ , hence these entries have all been already computed. Let us introduce some notation before presenting the recurrence. We define  $C(i, j, k)$  as the total length of edges going from  $\pi(i \rightarrow j)$  to  $\pi(j+1 \rightarrow k)$ . We now present the recursive formula for computing  $A_{i,j,q}$  as

$$A_{i,j,q} = \min_{i-1 \leq k \leq j} \left[ A_{i,k,1} + A_{k+1,j,q-1} + C(i, k, j) \right]. \quad (2)$$

Notice that the first term in the minimization is only used to signify whether  $\pi(i \rightarrow k)$  is a valid part in the intended  $\alpha$ -balanced partition: its value is either zero or infinity.

**Lemma 1.** Equation (2) is a valid recursion for  $A_{i,j,q}$ , which coupled with the initialization step given above yields a sound computation for all entries in the DP table.

**Proof.** The argument proceeds by mathematical induction. The initialization step is clearly sound as it is simply verifying feasibility of single parts.

For the inductive part, it suffices to show that the best solution with a part  $\pi(i \rightarrow k)$  has a cut size  $A_{i,k,1} + A_{k+1,j,q-1} + C(i, k, j)$ . Then, the first term in (2) accounts for the case when we use one empty part (i.e., we partition  $\pi(i \rightarrow j)$  using only  $q-1$  parts) and the minimization considers the cases when the first part consists of  $\pi(i \rightarrow k)$ . To see the cost of the latter is as in (2), notice that any cut edge in the partition of  $\pi(i \rightarrow j)$  is either completely inside  $\pi(i \rightarrow k)$  (i.e., not contributing to the cut size), inside  $\pi(k+1 \rightarrow j)$  (which is accounted for in  $A_{k+1,j,q-1}$  or connects  $\pi(i \rightarrow k)$  to  $\pi(k+1 \rightarrow j)$  (in which case is taken into account by  $C(i, k, j)$ ).  $\square$

The above procedure can be implemented to run in time  $O(n^3k)$  and consume space  $O(n^2k)$ . Next we show how to improve this further. See Algorithm 4. The key idea is to avoid computing  $A_{i,j,q}$  for all values  $1 \leq q \leq k$ . Rather we focus on a limited subset of such  $q$ 's of size  $O(\log k)$ .

---

**Algorithm 4** ImbalanceDP( $G, k, \pi$ )

---

**Input:** Graph  $G$ , number of partitions  $k$ , fixed ordering  $\pi$  of vertices

**Output:** A partition of  $V(G)$  into  $k$  parts

---

```

1: Let  $n = |V(G)|$  be number of vertices in  $G$ 
2: for all  $1 \leq i \leq j \leq n$  do
3:   Compute  $w(\pi(i \rightarrow j))$ , the total weight of interval  $\pi(i \rightarrow j)$ 
4: for all  $1 \leq i \leq j \leq n$  do
5:   if  $w(\pi(i \rightarrow j)) \leq (1 + \alpha) \frac{w(V)}{k}$  then
6:      $A_{i,j,1} \leftarrow 0$ 
7:   else
8:      $A_{i,j,1} \leftarrow \infty$ 
9: for all  $q = 2$  to  $k$  do
10:  if  $q = \lfloor k/2^l \rfloor$  or  $q = \lceil k/2^l \rceil$  for some  $l \geq 0$  then
11:    for all  $1 \leq i \leq j \leq n$  do
12:      for all  $i-1 \leq p < j$  do
13:        Compute  $C(i, p, j)$ , the total length of edges going from  $\pi(i \rightarrow j)$  to  $\pi(j+1 \rightarrow k)$ 
14:         $v \leftarrow A_{i,p,\lfloor \frac{q}{2} \rfloor} + A_{p+1,j,\lceil \frac{q}{2} \rceil} + C(i, p, j)$ 
15:        if  $A_{i,j,q} > v$  then
16:           $A(i, j, q) \leftarrow v$ 
17: return  $A_{1,n,k}$ 

```

---

We start by rewriting the recurrence as

$$A_{i,j,q} = \min_{i-1 \leq k < j} \left[ A_{i,k,\lfloor \frac{q}{2} \rfloor} + A_{k+1,j,\lceil \frac{q}{2} \rceil} + C(i,k,j) \right]. \quad (3)$$

Then, computing the desired value  $A_{1,n,k}$  only requires the computation of  $A_{i,j,q}$  where  $q$  is either  $\lfloor \frac{k}{2^l} \rfloor$  or  $\lceil \frac{k}{2^l} \rceil$  for some  $l \geq 0$ . This reduces the running time to  $O(n^3 \log k)$  and a bottom-up DP computation needs no more than three different values for  $q$ , hence the memory requirement is  $O(n^2)$ .

### 5.2. Scalable Linear Boundary Optimization

Recall the notion of “windows” defined at the beginning of Section 5. We now focus on a window  $W = W_j$ . This is small enough so that all information on the vertices in the window (including all their edges) fit in the memory (of one Mapper or Reducer). The linear optimization postprocessing finds a new split point in each window, so that the total weight of edges crossing it is minimized.

Let  $B$  and  $A$  denote the set of vertices appearing before and after  $W$ , respectively, in the ordering  $\pi$ . The edges going from  $B$  to  $A$  are irrelevant to the local optimization since those edges are necessarily cut no matter what split point we choose. The other edges have at least one endpoint in  $W$ . Then, there is a simple algorithm of running time  $O(|V_W| \cdot |E_W|)$  to find the best split point in  $W$  where  $V_W$  and  $E_W$  denote the set of vertices and edges, respectively, corresponding to  $W$ : look at each candidate split point and go over all relevant edges to determine the weight of the associated cut.

This can be done more efficiently to run in time  $O(|V_W| + |E_W|)$ , too. Let  $s_v$  for  $v \in W$  be the cut value (ignoring the effect of the edges between  $B$  and  $A$ ). Let  $s_0 = s_{\pi_{q_j - \alpha'}}$  (i.e., the cut value for the first split point in  $W$ ), which is the total weight of edges between  $B$  and  $W$  and itself can be computed in time  $O(|V_W| + |E_W|)$  (without hurting the overall runtime guarantee). We scan the candidate split points from left to right and compute  $s_v - s_0$ , however, the additive term  $-s_0$  does not matter in comparing the different split point candidates.

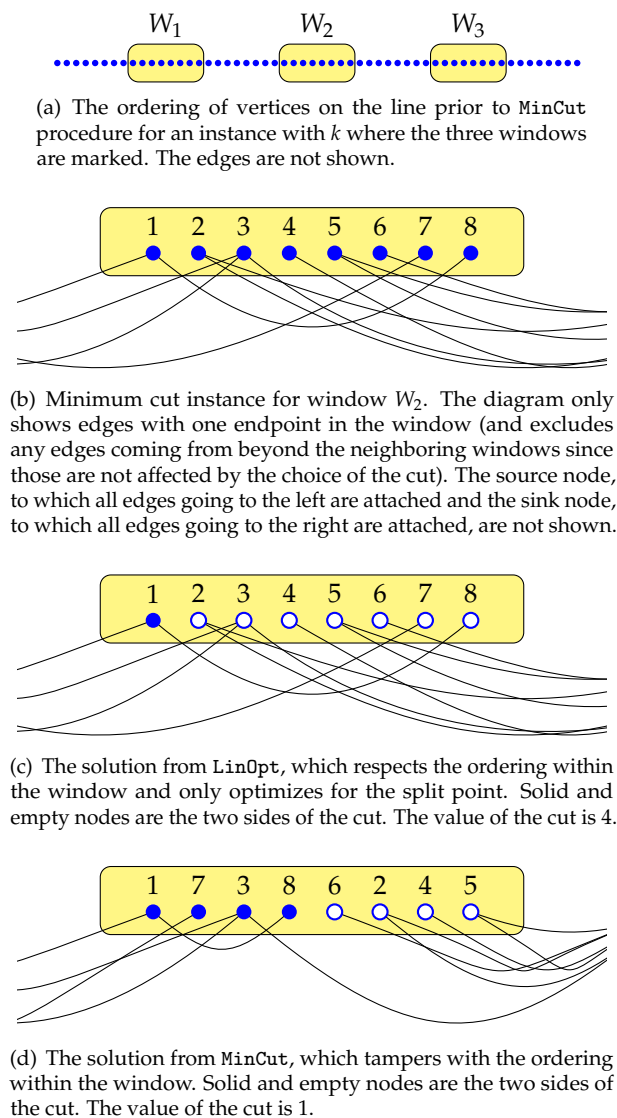
For each vertex  $v$  following  $u$ , start with  $s_u - s_0$ , remove the edges between  $v$  and the vertices to its left (these already showed up in the cut value) and add up the edges between  $u$  and the vertices to its right; this gives the cut value  $s_v - s_0$ .

Since the windows are relatively small compared to the entire graph, and their corresponding subproblems can be solved independently, there is a very efficient MapReduce implementation for window-based postprocessing methods.

### 5.3. Scalable Minimum-Cut Optimization

Once again we focus on a window  $W = W_j$  and denote by  $B$  and  $A$  the set of vertices belonging to previous and following vertices. Recall that the linear optimization of the previous section finds the best split point in the window while respecting the ordering of the vertices given in the permutation in  $\pi$ . The minimum-cut optimization, though, finds the best way to partition  $W$  into two pieces  $W_L$  and  $W_R$  so that the total weight of edges going from  $B \cup W_L$  to  $A \cup W_R$  is minimized. See Figure 3 for an illustration of this method.





**Figure 3.** Illustration of MinCut postprocessing.

#### 5.4. Parallel Implementation

All the components of our algorithm (described in Sections 3–5) may be implemented in a distributed fashion, using a framework such as MapReduce [12]. Indeed the empirical study of Section 6 is based on an internal MapReduce implementation. An immediate advantage of this is that we can run our algorithm on very large graphs (with billions of nodes), as the implementation easily scales horizontally by

- dividing the graph into pieces which can be hosted in a distributed file system and
- launching multiple worker processes that operate on each piece independently.

It is worthwhile to note that in some cases workers might need access to other pieces (for example, the Affinity-based mapping described in Section 3.3). A typical solution that scales well in such cases is to temporarily use Distributed Hash Tables (DHT) [38]. The disadvantage of using a DHT is that it requires additional system resources. Some parts of the algorithm, especially the post-processing operations described in Sections 4 and 5, may optionally be implemented using a DHT.

In a follow-up work [42], we present more details about the inner workings of the hierarchical clustering algorithm and provide a theoretical backing for its performance. In fact, a recent blog post [43] briefly summarizes the said work and the current paper.

## 6. Empirical Studies

First we describe the different datasets used in our experiments. Next we compare our results to previous work and we also compare our different methods to each other: that is, we demonstrate through experiments how much value each ingredient of the algorithm adds to the solution. Finally we discuss the scalability of our approach by running the algorithm on instances of varying sizes (but of the same type).

### 6.1. Datasets

We present our results on three datasets: World-roads, Twitter and LiveJournal (as well as publish our output on Friendster). As the names suggest, the first one is a geographic dataset while the other three are social graphs. All are big graphs, representatives of maps and social networks and we test the quality and scalability of our algorithm on them.

World-roads	A subset of the entire world road network with hundreds of millions of vertices and over a billion edges. (Due to near-planarity, the average degree is small.) Edges do not have weights but vertices have longitude/latitude information. Our algorithms run smoothly on this graph in reasonable time using a small number of machines, demonstrating their scalability.
Twitter	The public graph of tweets, with about 41 million vertices (twitter accounts) and 2.4 billion (directed) edges (denoting followership) [44]. This graph is unweighted, too. We run all our algorithms on the undirected underlying graph.
LiveJournal	The undirected version of this public social graph (snapshot from 2006) has 4.8 million vertices and 42.9 million edges [4].
Friendster	The undirected version of this public social graph (snapshot from 2006) has 65.6 million vertices and 1.8 billion [45].

Table 1 presents a summary for these datasets. Reference [5] contains graphs other than Twitter, which are either small or not public.

**Table 1.** Statistics about datasets used. (Numbers for World-roads are approximate).

Dataset	V	E	Max Degree
LiveJournal	4,847,571	42,851,237	20,333
Twitter	41,652,230	1,202,513,046	2,997,487
Friendster	65,608,366	1,806,067,135	3615
World-roads	$>10^8$	$>10^9$	$<20$

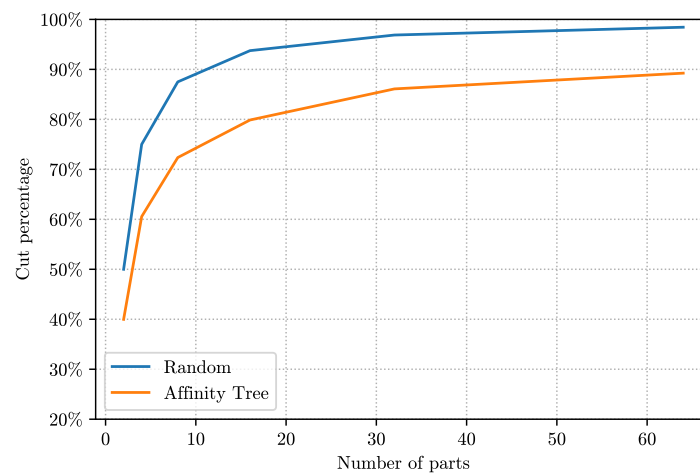
### 6.2. Comparison of Different Techniques

#### 6.2.1. Initial Ordering

We have four methods to obtain an initial ordering for the geographic graphs and two for non-geographic graphs. Here we compare these methods to each other based on the value of the cut obtained by chopping the resulting order into equal contiguous pieces. In order to make this meaningful, we report all the cut sizes as the fraction of cut edges to the total number of edges in the graph.

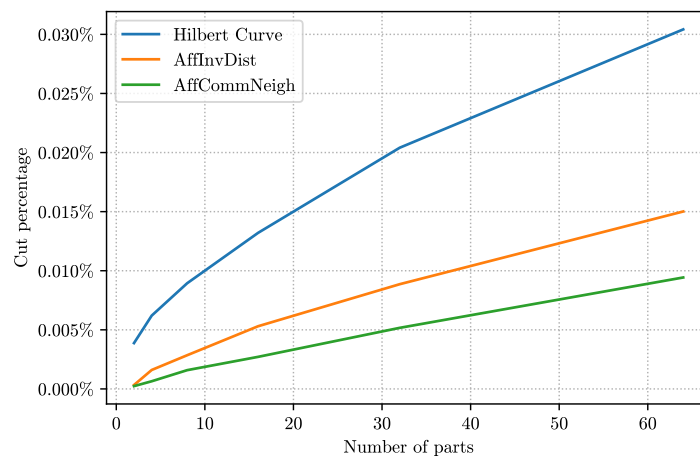
Figure 4 compares the results of two methods: RandInit produces a random permutation of the vertices, hence, as we observed previously, gets a cut size of approximately  $1 - \frac{1}{k}$ ; AffCommNeigh (using the AffinityOrdering with the similarity oracle based on the number of common neighbors) clearly produces better solutions with improvements ranging from 20% to 10% (more improvement for smaller number of parts). The AffCommNeigh tree in this case was pretty shallow (much lower than the

theoretical  $\log n$  upper bound) with only four levels. Whenever we report results for different number of partitions in one plot,  $k = 2, 4, 8, 16, 32, 64$  are used for experiments.



**Figure 4.** Comparison of the (fully balanced) cut size for RandInit and AffCommNeigh on Twitter.

Figure 5 compares the results of three methods, each producing a cut of size at most 3% and significantly improving upon the  $1 - \frac{1}{k}$  result of RandInit: Compared to Hilbert, the two other methods AffInvDist (using the inverse of geographic distance as the similarity oracle) and AffCommNeigh, respectively, obtain 90% to 50% and 95% to 70% improvement; once again the easier instances are those with smaller  $k$ . The corresponding trees, with 13 and 11 levels, respectively, are not as shallow as the one for Twitter. The quality of the cuts correlates with the runtime complexity of the algorithms: RandInit (fastest and worst), Hilbert, AffInvDist and AffCommNeigh (best and slowest). Note that although the implementation of Hilbert seems more complicated, it is significantly faster in practice.



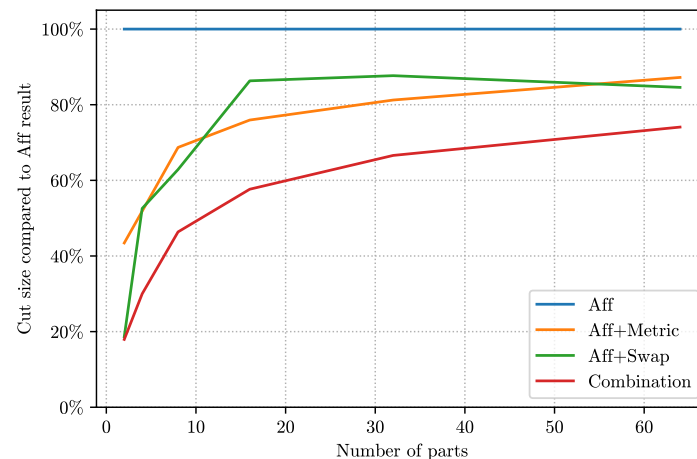
**Figure 5.** Comparison of the (fully balanced) cut size for Hilbert, AffInvDist and AffCommNeigh on World-roads.

For the following experiments, we focus on AffCommNeigh since it consistently produces better results for different values of  $k$  as well as both for Twitter and World-roads. Specially in diagram legends, this initialization step may be abbreviated as Aff.

### 6.2.2. Improvements and Imbalance

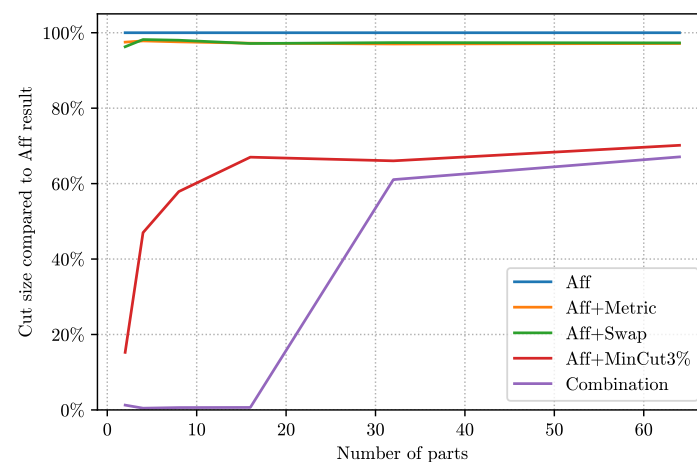
The performance of the semilocal improvement methods, Aff+Metric and Aff+Swap, as well as our main algorithm, Combination, on Twitter is reported in Figure 6. The diagram shows

the percentage of improvement of each over the baseline AffCommNeigh for different values of  $k = 2, 4, 8, 16, 32, 64$ . Except for  $k = 16, 32$ , Aff+Swap outperforms Aff+Metric; their performance varies from 7% to 75%, with best performance for smaller  $k$ . Combining the two methods with the imbalance-inducing techniques to get Combination yields results significantly better than every single ingredient, with the relative improvement sometimes as large as 50%. The final results have very little imbalance and indeed cuts of comparable quality can be obtained without any imbalance.



**Figure 6.** Comparison of the (fully balanced) cut size for three semilocal operations on Twitter. The results are compared to AffCommNeigh, set as 100%.

Figure 7 also sets the results of various postprocessing methods against one another and depicts their improvements compared to the baseline of AffCommNeigh when run on World-roads. The first two consists of semilocal optimization methods, Aff+Metric and Aff+Swap, which also appeared in the discussion for Twitter. These do not yield significant improvement over their AffCommNeigh starting point: the results compared to the starting point (baseline) ranges from 95% to 98%. The other two algorithms use imbalance-inducing ideas and yield significantly better results. Our main algorithm, which we call Combination, clearly outperforms all the others specially for small values of  $k$ .

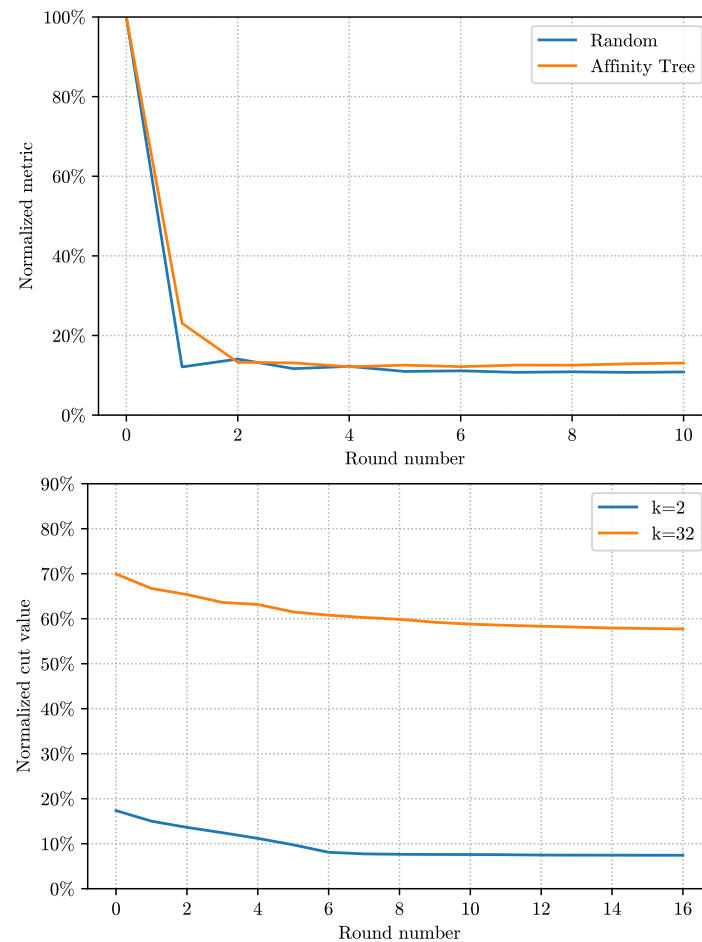


**Figure 7.** Cut value obtained by semilocal optimization methods on World-roads as compared to their baseline, AffCommNeigh.

All the experiments assume 3% imbalance in partition sizes—each partition may have size within range  $(1 \pm 0.03)\frac{n}{k}$ . However, only Aff+MinCut and Combination use this to improve the quality of the partition.

### 6.2.3. Convergence Analysis

Next we consider the convergence rate of the two semilocal improvement methods on two graphs. For Twitter, the rate of convergence for metric optimization is the same whether we start with AffCommNeigh or RandInit (See Figure 8). The convergence happens essentially in three or four rounds.



**Figure 8.** Convergence rate of Aff+Metric (top) and Aff+Swap (bottom) on Twitter.

The convergence for the rank swap methods happens in 10–15 steps for Twitter and in 5–10 steps for World-roads; more steps required for larger values of  $k$ .

The convergence for the metric optimization on World-roads happens in 3–4 steps if the starting point is AffCommNeigh. However, if we start with a random ordering, it may take up to 20 steps for the metric to converge. See Figure 9.

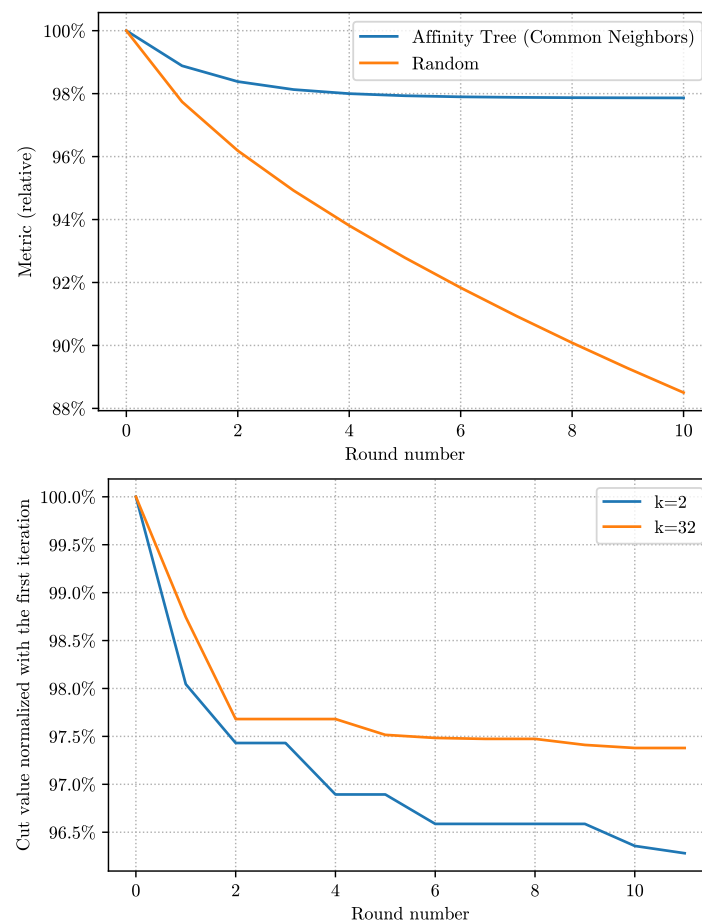


Figure 9. Convergence rate of Aff+Metric (top) and Aff+Swap (bottom) on World-roads.

### 6.3. Comparison to Previous Work

The most relevant among the previous work are the scalable label propagation-based algorithms of Ugander and Backstrom [4] and Martella et al. [19]. The former reports their results on an internal graph (for Facebook) and a public social graph, LiveJournal, while the latter reports results on LiveJournal and Twitter. Table 2 compares our results to prior work on LiveJournal. The cut sizes in the table for  $k = 40, 60, 80$  for Reference [4] and all the numbers for Reference [19] are approximates taken from the graph provided in Reference [4] as they do not report the exact values. The numbers in parentheses show the maximum imbalance in partition sizes, and the cut sizes are reported as fractions of total number of edges in the graph.

**Table 2.** Comparison to prior work on LiveJournal. The main entries are the ratio of edges cut by the clustering over the total number of edges in the graph. The number of clusters is denoted by  $k$ . Numbers in parentheses are permissible imbalance (Best results in each row are typeset in bold for easier spotting).

$k$	[4] (−0.5 mm) (5%)	Spinner (5%)	Aff (0%)	Combination (0%)
20	37%	38%	35.71%	<b>27.50%</b>
40	43%	40%	40.83%	<b>33.71%</b>
60	46%	43%	43.03%	<b>36.65%</b>
80	47.5%	44%	43.27%	<b>38.65%</b>
100	49%	46%	45.05%	<b>41.53%</b>



Even our initial linear embedding (obtained by constructing the hierarchical clustering based on a weighted version of the input, where edge weights denote the number of common neighbors between two vertices) is consistently better than the best previous result. Our algorithm with the postprocessing obtains 15% to 25% improvement over the previous results (better for smaller  $k$ ) and it only requires a couple of postprocessing rounds to converge. We take note that the results of Reference [4] allow 5% imbalance whereas our results produce almost perfectly balanced partitions.

Twitter is another public graph, for which the results of state-of-the-art minimum-cut partitioning is available. In Table 3 we report and compare the results of our main algorithm, Combination, for this graph to the best previous methods. (Numbers for Reference [21] are quoted from Reference [19]).

**Table 3.** Comparison to prior work on Twitter. The main entries are the ratio of edges cut by the clustering over the total number of edges in the graph. The number of clusters is denoted by  $k$ . Numbers in parentheses denote the permissible imbalance for each algorithm (Best results in each row are typeset in bold for easier spotting).

$k$	[21] (4–15%)	Spinner (5%)	FENNEL (10%)	METIS (3%)	Combination (3%)
2	34%	15%	<b>6.8%</b>	11.98%	7.43%
4	55%	31%	29%	24.39%	<b>18.16%</b>
8	66%	49%	48%	35.96%	<b>33.55%</b>
16	76%	61%	59%	N/A	<b>46.18%</b>
32	80%	69%	67%	N/A	<b>57.67%</b>

Algorithms developed in References [5,21] are suitable for the streaming model but one can implement variants of those algorithms in a distributed manner. Our implementation of a natural distributed version of the FENNEL algorithm does not achieve the same results as those reported for the streaming implementation reported in Reference [5]. However, we compare our algorithm directly to the numbers reported on Twitter by FENNEL [5].

Finally we report in Table 4 the cut sizes for another public graph, Friendster, so others can compare their results with it. The running time of our algorithms are not affected with large  $k$ ; the running time difference for  $k = 2$  and tens of thousands is less than 1%, well within the noise associated with the distributed system. In fact, construction of the initial ordering is independent of  $k$ , and the post-processing steps may only take advantage of the increased parallelism possible for large  $k$ .

**Table 4.** Cut sizes produced by our algorithm on Friendster (ratio of cut edges to the total number of edges), where  $k$  denotes the number of clusters.

$k$	2	10	$10^2$
Cut size	11.9%	41.4%	59.8%

We also make our partitions for Twitter, LiveJournal and Friendster publicly available [23].

#### 6.4. Scalability

We noted above that the choice of  $k$  does not affect the running time of the algorithm significantly: the running time for two and tens of thousands of partitions differed less than 1% for Friendster.

As another measure of its scalability, we run the algorithm with  $k = 2$  on a series of random graphs (that are similar in nature) of varying sizes. In particular, we use RMAT graphs [46] with parameter 20, 22, 24, 26 and 28, whose node and edge count is given in Table 5. In addition, the last column gives the normalized running time of our algorithm on these graphs. Note that the size of the graph almost quadruples from one graph to the next.

**Table 5.** General statistics (size) about the RMAT graphs as well as the relative running times of balanced-partitioning algorithm on these graphs.

Graph	V	E	Max Degree	Running Time
RMAT 20	650 K	31 M	65 K	100%
RMAT 22	2.4 M	130 M	160 K	110%
RMAT 24	8.9 M	525 M	400 K	133%
RMAT 26	32.8 M	2.1 B	1 M	160%
RMAT 28	120 M	8.5 B	2.5 M	402%

## 7. Applications

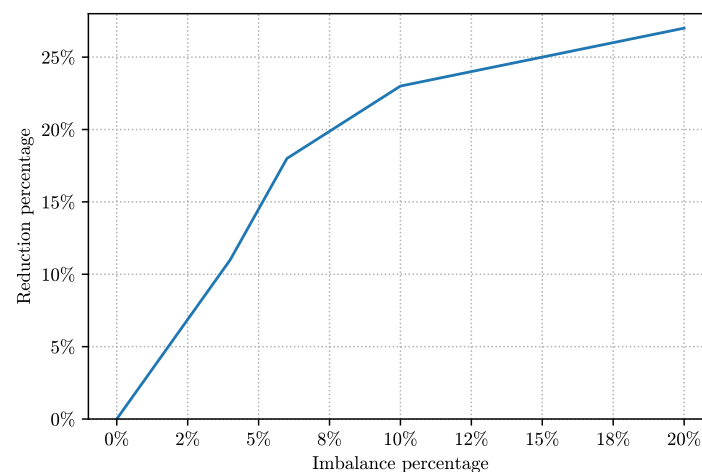
In the past five years, we have applied this algorithm to various applications within Google, including those in Maps, Search and Infrastructure domains. In this section, we briefly explain three examples.

### 7.1. Google Maps Driving Directions

As discussed in Section 1, we apply our results to the Google Maps Driving Directions application. Note that the evaluation metric for balanced partitioning of the world graph is the expected percentage of cross-shard queries over the total number of queries. More specifically, this objective function can be captured as follows: we first estimate the number of times each edge of graph may be part of the shortest path (or the driving direction) from the source to destination across all the query traffic. This produces an edge-weighted graph in which the weight of each edge is proportional to the number of times we expect this edge appearing in a driving direction and our goal is to partition the graph into a small number of pieces and minimize the total weight of the edges cut. We can estimate edge weights based on historical data and use it as a proxy for the number of times the edge appears on a driving direction on the real data. This objective is aligned with minimizing the weighted cut and we can apply our algorithms to solve this problem.

Before deciding about our live experiments, we realized it might be suitable to use an imbalance factor in dividing the graph and as a result, we first examined the best imbalance factor for our cut-optimization technique. The result of this study is summarized in Figure 10. In particular, we observe that we can reduce cross-shard queries by 21% when increasing the imbalance factor from 0% to 10%.

The two methods that we examined via live experiments were (i) a baseline approach based on the Hilbert-curve embedding and (ii) one method based on applying our cut-optimization postprocessing techniques. Note that both these algorithms first compute an embedding of nodes into a line, which results in a much simpler system to identify the corresponding shards for the source and destination of each query at serving time. Finally, by running live experiments on the real traffic, we observe that the number of multi-sharded queries from our cut-optimization techniques is almost 40% less than the number of multi-sharded queries compared to the naïve Hilbert embedding technique.



**Figure 10.** The percentage of improvement in cut value as we increase the allowable imbalance factor.

### 7.2. Load Balancing in Google Search Backend

Another application of this large-scale balanced-partitioning algorithm is in cache-aware load balancing of data centers [47]. While this is not the focus of the present work, we present a brief overview to demonstrate a different application.

The load-balancer of a big search engine distributes search queries among several identical serving replicas. Then the server pulls, for each term in the query, the index of all pages associated with the search term. As significant time and network traffic is consumed in loading these index tables, a cache is employed on the serving replicas. This leads to better utilization if each search term is usually sent to a specific machine. The complication, though, originates from the fact search queries often consist of several search terms, each with its own index table.

The proposed solution is a voting table to select the right replica for each query. The voting table is constructed via a multi-stage algorithm, yet the initial voting table—a decent solution in and of itself—comes from a balanced-partitioning formulation: We aim to partition the search terms in a balanced way, where the sizes denote lengths of index tables, so as to minimize the expected cache miss rate. This is done by carefully building a term-query bipartite graph, whose vertex weights correspond to index table sizes and whose edge weights are related to cache miss rates. See Reference [47] for more information.

Deployed in 2015, the said solution decreases cache miss rate by about 0.5% and increases end-to-end throughput by double-digit percentage.

### 7.3. Co-Location of Tasks in a Data Center

Yet another application of our balanced-partitioning algorithm is in scheduling of computation tasks within data centers. The cost of intra-network bandwidth in a data center is nonnegligible, and adds up over time. Besides, the underlying network topology may be hierarchical: the machines at a lower-level subtree can have much higher bandwidth and may even be organized as a clique. Consequently the communication cost between machines across different subtrees can be much higher.

One way to reduce network bandwidth usage is to co-locate related tasks within closer physical proximity and therefore optimize the network traffic. As communication patterns of many applications stay stable over time, one can look into historical logs to construct a graph of communication among tasks. In this graph, each node is an application (or task) and the edge weight between two nodes is the communication cost between them. Machines have limited size or memory, thus we naturally cast the problem as balanced partitioning: We cluster the tasks (subject to a physical limit) and minimize the communication cost across clusters.

Our initial implementation of this technique (within a subset of the overall traffic) resulted in a six-fold increase of co-located tasks.

#### 7.4. Randomized Experimental Design via Geographic Clustering

Recently, the authors of Reference [48] used balanced partitioning as a way to cluster geographic regions in order to set up randomized experiments over them. This is especially useful in web-based services, since it does not require tracking individual users or browser cookies. Given the fact that users, especially mobile ones, may issue queries from different geographical locations, it cannot be assumed that geographic regions are independent and interference may still be present in the experiments. These characteristics of the problem make it suitable for applying balanced partitioning where the geographic regions are modeled as a graph and the edge weights across output clusters are minimized while producing clusters of almost the same size.

### 8. Conclusions

We develop a three-stage algorithm for balanced partitioning of a graph, with the main objective of minimizing the overall weight of cut edges.

The first two stages construct an ordering of the vertices based on the neighborhoods in the graph—nearby vertices are expected to be close to each other in the ordering.

The third stage considers the permissible imbalance and moves the boundaries of parts so as to optimize the cut value. It may also permute vertices that are very close to the part boundaries.

Though we give a few variants for each stage, it is not necessary to run all the three stages; some can be replaced by a dummy routine. To get the full potential of the algorithm, one had better use the best variants of every stage, however, even using random ordering in the first stage, for instance, produces a nontrivial result. Similarly, only using a natural balanced partitioning based on the hierarchical agglomerative clustering produces decent results. Interestingly, this outperforms the Hilbert cover-based ordering, even for the geographic graphs where the latter technique is feasible.

For  $k > 2$ , our algorithm beats METIS and FENNEL, and is more scalable than both of them. For  $k = 2$ , that is, the minimum bisection problem, while our algorithm is far superior in terms of running time and scalability, the results that we obtain are on par with (or slightly worse than) those of the other two.

Our algorithm is pretty scalable: for example, it runs smoothly on a portion of the world graph with half a billion vertices and more than a billion edges.

Furthermore, we study or deploy this solution in various applications at Google.

- A simple version of our algorithm (using Hilbert curves for the initial ordering, no optimization in stage two and dynamic programming in the third stage) yields a 40% improvement for the Google Maps Driving Directions application: the improvement is both in terms of overall CPU usage as well as the number of multishard queries.
- In another work [47], we show how to use balanced partitioning to optimize the load balancing component of Google Search backend. Deployed in 2015, the said solution decreases cache miss rate by about 0.5% and increases end-to-end throughput by double-digit percentage.
- This algorithm has been effective in scheduling jobs in data centers with the objective of minimizing communication.
- It is also demonstrated in a separate work [48] that experiment design can benefit from a balanced partitioning of users based on their geography. This helps reduce the dependency between various experiment arms.

**Author Contributions:** Conceptualization, K.A., M.B. and V.M.; Data curation, K.A.; Methodology, V.M.; Software, K.A. and M.B.; Supervision, V.M.; Visualization, K.A. and M.B.; Writing—Original Draft, M.B.; Writing—Review & Editing, K.A., M.B. and V.M.

**Funding:** This research received no external funding.

**Acknowledgments:** The authors wish to thank Aaron Archer and Raimondas Kiveris for fruitful discussions.

**Conflicts of Interest:** The authors declare no conflict of interest.

## References

1. Dongarra, J.; Foster, I.; Fox, G.; Gropp, W.; Kennedy, K.; Torczon, L.; White, A. *The Sourcebook of Parallel Computing*; Morgan Kaufmann Publishers Inc.: San Francisco, CA, USA, 2003.
2. Andreev, K.; Räcke, H. Balanced Graph Partitioning. *Theory Comput. Syst.* **2006**, *39*, 929–939. [\[CrossRef\]](#)
3. Garey, M.R.; Johnson, D.S. *Computers and Intractability: A Guide to the Theory of NP-Completeness*; W. H. Freeman and Company: New York, NY, USA, 1979.
4. Ugander, J.; Backstrom, L. Balanced label propagation for partitioning massive graphs. In Proceedings of the Sixth ACM International Conference on Web Search and Data Mining (WSDM '13), Rome, Italy, 4–8 February 2013; pp. 507–516.
5. Tsourakakis, C.E.; Gkantsidis, C.; Radunovic, B.; Vojnovic, M. FENNEL: Streaming graph partitioning for massive scale graphs. In Proceedings of the Seventh ACM International Conference on Web Search and Data Mining (WSDM 2014), New York, NY, USA, 24–28 February 2014; pp. 333–342. [\[CrossRef\]](#)
6. Delling, D.; Goldberg, A.V.; Razenshteyn, I.; Werneck, R.F.F. Graph Partitioning with Natural Cuts. In Proceedings of the 25th IEEE International Symposium on Parallel and Distributed Processing (IPDPS 2011), Anchorage, AK, USA, 16–20 May 2011; pp. 1135–1146. [\[CrossRef\]](#)
7. Delling, D.; Goldberg, A.V.; Razenshteyn, I.; Werneck, R.F.F. Exact Combinatorial Branch-and-Bound for Graph Bisection. In Proceedings of the 14th Meeting on Algorithm Engineering & Experiments (ALENEX 2012), Kyoto, Japan, 16 January 2012; pp. 30–44. [\[CrossRef\]](#)
8. Malewicz, G.; Austern, M.H.; Bik, A.J.; Dehnert, J.C.; Horn, I.; Leiser, N.; Czajkowski, G. Pregel: A system for large-scale graph processing. In Proceedings of the 2010 ACM SIGMOD International Conference on Management of Data (SIGMOD '10), Indianapolis, IN, USA, 6–10 June 2010.
9. Ching, A.; Kunz, C. Giraph: Large-Scale Graph Processing on Hadoop. 2011. Available online: <https://www.youtube.com/watch?v=l4nQjAG6fac> (accessed on 8 August 2019).
10. Kang, U.; Tsourakakis, C.E.; Faloutsos, C. PEGASUS: A Peta-Scale Graph Mining System. In Proceedings of the Ninth IEEE International Conference on Data Mining (ICDM 2009), Miami, FL, USA, 6–9 December 2009; pp. 229–238. [\[CrossRef\]](#)
11. Stanton, I. Streaming Balanced Graph Partitioning Algorithms for Random Graphs. In Proceedings of the Twenty-Fifth Annual ACM-SIAM Symposium on Discrete Algorithms (SODA 2014), Portland, OR, USA, 5–7 January 2014; pp. 1287–1301. [\[CrossRef\]](#)
12. Dean, J.; Ghemawat, S. MapReduce: Simplified Data Processing on Large Clusters. In Proceedings of the 6th Symposium on Operating Systems Design and Implementation (OSDI), San Francisco, CA, USA, 6–8 December 2004; pp. 137–150.
13. Aydin, K.; Bateni, M.; Mirrokni, V.S. Distributed Balanced Partitioning via Linear Embedding. In Proceedings of the Ninth ACM International Conference on Web Search and Data Mining, San Francisco, CA, USA, 22–25 February 2016; pp. 387–396. [\[CrossRef\]](#)
14. Eckles, D.; Karrer, B.; Ugander, J. Design and analysis of experiments in networks: Reducing bias from interference. *J. Causal Inference* **2017**, *5*. [\[CrossRef\]](#)
15. Pouget-Abadie, J.; Mirrokni, V.; Parkes, D.C.; Airolidi, E.M. Optimizing Cluster-based Randomized Experiments Under Monotonicity. In Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD '18), London, UK, 19–23 August 2018; pp. 2090–2099.
16. Boldi, P.; Rosa, M.; Santini, M.; Vigna, S. Layered label propagation: A multiresolution coordinate-free ordering for compressing social networks. In Proceedings of the 20th International Conference on World Wide Web (WWW '11), Hyderabad, India, 28 March–1 April 2011; pp. 587–596. [\[CrossRef\]](#)
17. Chierichetti, F.; Kumar, R.; Lattanzi, S.; Mitzenmacher, M.; Panconesi, A.; Raghavan, P. On compressing social networks. In Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '09), Paris, France, 28 June–1 July 2009; pp. 219–228. [\[CrossRef\]](#)
18. Tsourakakis, C.E.; Kolountzakis, M.N.; Miller, G.L. Approximate Triangle Counting. *arXiv* **2009**, arXiv:0904.3761.

19. Martella, C.; Logothetis, D.; Siganos, G. Spinner: Scalable Graph Partitioning for the Cloud. *arXiv* **2014**, arXiv:1404.3861.
20. Karypis, G.; Kumar, V. A Fast and High Quality Multilevel Scheme for Partitioning Irregular Graphs. *SIAM J. Sci. Comput.* **1998**, *20*, 359–392. [[CrossRef](#)]
21. Stanton, I.; Kliot, G. Streaming graph partitioning for large distributed graphs. In Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '12), Beijing, China, 12–16 August 2012; pp. 1222–1230. [[CrossRef](#)]
22. Yan, D.; Cheng, J.; Lu, Y.; Ng, W. Blogel: A Block-Centric Framework for Distributed Computation on Real-World Graphs. *PVLDB* **2014**, *7*, 1981–1992. [[CrossRef](#)]
23. Aydin, K.; Bateni, M.; Mirrokni, V. Public Data for Balanced Partitioning Paper. Available online: <http://goo.gl/okvwpa> (accessed on 8 August 2019).
24. Goldschmidt, O.; Hochbaum, D.S. Polynomial Algorithm for the  $k$ -Cut Problem. In Proceedings of the 29th Annual Symposium on Foundations of Computer Science (FOCS 1988), White Plains, NY, USA, 24–26 October 1988; pp. 444–451. [[CrossRef](#)]
25. Feige, U.; Krauthgamer, R. A Polylogarithmic Approximation of the Minimum Bisection. *SIAM J. Comput.* **2002**, *31*, 1090–1118. [[CrossRef](#)]
26. Ene, A.; Im, S.; Moseley, B. Fast clustering using MapReduce. In Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '11), San Diego, CA, USA, 21–24 August 2011; pp. 681–689.
27. Bahmani, B.; Moseley, B.; Vattani, A.; Kumar, R.; Vassilvitskii, S. Scalable  $K$ -Means++. *PVLDB* **2012**, *5*, 622–633. [[CrossRef](#)]
28. Balcan, M.F.; Ehrlich, S.; Liang, Y. Distributed  $k$ -means and  $k$ -median clustering on general communication topologies. In Proceedings of the NIPS, Lake Tahoe, NV, USA, 5–10 December 2013.
29. Bateni, M.; Bhaskara, A.; Lattanzi, S.; Mirrokni, V. Distributed Balanced Clustering via Mapping Coresets. In Proceedings of the Advances in Neural Information Processing Systems 27: 28th Annual Conference on Neural Information Processing Systems 2014, Montreal, QC, Canada, 8–13 December 2014.
30. Strongin, R.G.; Sergeyev, Y.D. Global Optimization with Non-Convex Constraints: Sequential and Parallel Algorithms. In *Nonconvex Optimization and Its Applications*; Kluwer Academic Publishers: Dordrecht, The Netherlands; Boston, MA, USA; London, UK, 2000.
31. Moon, B.; Jagadish, H.V.; Faloutsos, C.; Saltz, J.H. Analysis of the clustering properties of the Hilbert space-filling curve. *IEEE Trans. Knowl. Data Eng.* **2001**, *13*, 2001. [[CrossRef](#)]
32. Gotsman, C.; Lindenbaum, M. On the metric properties of discrete space-filling curves. *IEEE Trans. Image Process.* **1996**, *5*, 794–797. [[CrossRef](#)] [[PubMed](#)]
33. Niedermeier, R.; Reinhardt, K.; Sanders, P. Towards optimal locality in mesh-indexings. *Discret. Appl. Math.* **2002**, *117*, 211–237. [[CrossRef](#)]
34. Fishburn, P.C.; Tetali, P.; Winkler, P. Optimal linear arrangement of a rectangular grid. *Discret. Math.* **2000**, *213*, 123–139. [[CrossRef](#)]
35. Sokal, R.R.; Michener, C.D. A statistical method for evaluating systematic relationships. *Univ. Kans. Sci. Bull.* **1958**, *38*, 1409–1438.
36. Borůvka, O. Příspěvek k řešení otázky ekonomické stavby elektrovodních sítí (Contribution to the solution of a problem of economical construction of electrical networks). *Elektron. Obz.* **1926**, *15*, 153–154. (In Czech)
37. Rastogi, V.; Machanavajjhala, A.; Chitnis, L.; Sarma, A.D. Finding connected components in map-reduce in logarithmic rounds. In Proceedings of the 29th IEEE International Conference on Data Engineering (ICDE 2013), Brisbane, QLD, Australia, 8–12 April 2013; pp. 50–61. [[CrossRef](#)]
38. Kiveris, R.; Lattanzi, S.; Mirrokni, V.; Rastogi, V.; Vassilvitskii, S. Connected Components in MapReduce and Beyond. In Proceedings of the Fifth ACM Symposium on Cloud Computing, SOCC 2014, Seattle, WA, USA, 3–5 November 2014.
39. Rao, S.; Richa, A.W. New Approximation Techniques for Some Linear Ordering Problems. *SIAM J. Comput.* **2004**, *34*, 388–404. [[CrossRef](#)]
40. Feige, U.; Lee, J.R. An improved approximation ratio for the minimum linear arrangement problem. *Inf. Process. Lett.* **2007**, *101*, 26–29. [[CrossRef](#)]
41. Charikar, M.; Hajiaghayi, M.T.; Karloff, H.J.; Rao, S.  $l_2^2$  Spreading Metrics Vertex Ordering Probl. *Algorithmica* **2010**, *56*, 577–604. [[CrossRef](#)]



42. Bateni, M.; Behnezhad, S.; Derakhshan, M.; Hajiaghayi, M.; Kiveris, R.; Lattanzi, S.; Mirrokni, V.S. Affinity Clustering: Hierarchical Clustering at Scale. In Proceedings of the Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems, Long Beach, CA, USA, 4–9 December 2017; pp. 6867–6877.
43. Aydin, K.; Bateni, M. Balanced Partitioning and Hierarchical Clustering at Scale. Available online: <https://ai.googleblog.com/2018/03/balanced-partitioning-and-hierarchical.html> (accessed on 8 August 2019).
44. Kwak, H.; Lee, C.; Park, H.; Moon, S. What is Twitter, a social network or a news media? In Proceedings of the 19th International Conference on World Wide Web (WWW '10), Raleigh, NC, USA, 26–30 April 2010; ACM: New York, NY, USA, 2010; pp. 591–600. [CrossRef]
45. Yang, J.; Leskovec, J. Defining and Evaluating Network Communities Based on Ground-Truth. In Proceedings of the 12th IEEE International Conference on Data Mining (ICDM 2012), Brussels, Belgium, 10–13 December 2012; pp. 745–754. [CrossRef]
46. Chakrabarti, D.; Zhan, Y.; Faloutsos, C. R-MAT: A Recursive Model for Graph Mining. In Proceedings of the Fourth SIAM International Conference on Data Mining, Lake Buena Vista, FL, USA, 22–24 April 2004; pp. 442–446. [CrossRef]
47. Archer, A.; Aydin, K.; Bateni, M.; Mirrokni, V.S.; Schild, A.; Yang, R.; Zhuang, R. Cache-aware load balancing of data center applications. *PVLDB* **2019**, *12*, 709–723. [CrossRef]
48. Rolnick, D.; Aydin, K.; Pouget-Abadie, J.; Kamali, S.; Mirrokni, V.; Najmi, A. Randomized Experimental Design via Geographic Clustering. In Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, Anchorage, AK, USA, 4–8 August 2019.



© 2019 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>).