



ELSEVIER

Contents lists available at [ScienceDirect](#)

# Parallel Computing

journal homepage: [www.elsevier.com/locate/parco](http://www.elsevier.com/locate/parco)

## Incremental closeness centrality in distributed memory

Ahmet Erdem Sariyüce<sup>a,b,\*</sup>, Erik Saule<sup>d</sup>, Kamer Kaya<sup>e,a</sup>, Ümit V. Çatalyürek<sup>a,c</sup><sup>a</sup> Dept. Biomedical Informatics, The Ohio State University, United States<sup>b</sup> Dept. Computer Science and Engineering, The Ohio State University, United States<sup>c</sup> Dept. Electrical and Computer Engineering, The Ohio State University, United States<sup>d</sup> Dept. Computer Science, University of North Carolina at Charlotte, United States<sup>e</sup> Dept. Computer Science and Engineering, Sabanci University, Turkey

### ARTICLE INFO

#### Article history:

Available online 22 January 2015

#### Keywords:

Closeness centrality  
Incremental centrality  
BFS  
Parallel programming  
Cluster computing

### ABSTRACT

Networks are commonly used to model traffic patterns, social interactions, or web pages. The vertices in a network do not possess the same characteristics: some vertices are naturally more connected and some vertices can be more important. Closeness centrality (CC) is a global metric that quantifies how important is a given vertex in the network. When the network is dynamic and keeps changing, the relative importance of the vertices also changes. The best known algorithm to compute the CC scores makes it impractical to recompute them from scratch after each modification. In this paper, we propose *STREAMER*, a distributed memory framework for incrementally maintaining the closeness centrality scores of a network upon changes. It leverages pipelined, replicated parallelism, and SpMM-based BFSs, and it takes NUMA effects into account. It makes maintaining the closeness centrality values of real-life networks with millions of interactions significantly faster and obtains almost linear speedups on a 64 nodes 8 threads/node cluster.

© 2015 Elsevier B.V. All rights reserved.

## 1. Introduction

How central is a vertex in a network? Which vertices are more important during an entity dissemination? Centrality metrics have been used to answer such questions. They have been successfully used to carry analysis for various purposes such as power grid contingency analysis [16], quantifying importance in social networks [23], analysis of covert networks [18], decision/action networks [8], and even for finding the best store locations in cities [26]. As the networks become large, efficiency becomes a crucial concern while analyzing these networks. The algorithm with the best asymptotic complexity to compute the closeness and betweenness metrics [4] is believed to be asymptotically optimal [17]. The research on fast centrality computation have focused on approximation algorithms [7,9,24] and high performance computing techniques [22,32,20]. Today, the networks to be analyzed can be quite large, and we are always in a quest for faster techniques which help us to perform centrality-based analysis.

Many of today's networks are dynamic. And for such networks, maintaining the exact centrality scores is a challenging problem which has been studied in the literature [10,19,27]. The problem can also arise for applications involving static networks such as the power grid contingency analysis and robustness evaluation of a network. The findings of such analyses and evaluations can be very useful to be prepared and take proactive measures; for instance if there is a natural risk or a possible adversarial attack that can yield undesirable changes on the network topology in the future. Similarly, in some applications,

\* Corresponding author at: 250 Lincoln Tower, 1800 Canon Drive Columbus, OH 43210, United States. Tel.: +1 614 688 9637; fax: +1 614 688 6600.

E-mail addresses: [sariyuce.1@osu.edu](mailto:sariyuce.1@osu.edu) (A.E. Sariyüce), [esaule@uncc.edu](mailto:esaule@uncc.edu) (E. Saule), [kaya@sabanciuniv.edu](mailto:kaya@sabanciuniv.edu) (K. Kaya), [umit@bmi.osu.edu](mailto:umit@bmi.osu.edu) (Ü.V. Çatalyürek).

one might be interested in trying to find the minimal topology modifications on a network to set the centrality scores in a controlled manner. (Applications include speeding-up or containing the entity dissemination, and making the network immune to adversarial attacks).

Offline closeness centrality (CC) computation can be expensive for large-scale networks. Yet, one could hope that the incremental graph modifications can be handled in an inexpensive way. Unfortunately, as Fig. 1 shows, the effect of a local topology modification can be global. In a previous study, we proposed a sequential incremental closeness centrality algorithm which is orders of magnitude faster than the best offline algorithm [27]. Still, the algorithm was not fast enough to be used in practice. In a previous work, we proposed STREAMER [29] to parallelize these incremental algorithms. In this paper, we present an improved version of STREAMER, to efficiently parallelize the incremental CC computation on high-performance clusters.

The best available algorithm for the offline centrality computation is pleasingly parallel (and scalable if enough memory is available) since it involves  $n$  independent executions of the single-source shortest path (SSSP) algorithm [4]. In a naive distributed framework for the offline case, one can distribute the SSSPs to the nodes and gather their results. Here the computation is static, i.e., when the graph changes, the previous results are ignored and the same  $n$  SSSPs are re-executed. On the other hand, in the online approach, the graph modifications can arrive at any time even while the centrality scores for a previous modification are still being computed. Furthermore, the scores which need to be recomputed (the SSSPs that need to be executed) change depending on the modification of the graph. Finding these SSSPs and distributing them to the nodes is not a straightforward task. To be able to do that, the incremental algorithms maintain complex information such as the biconnected component decomposition of the current graph [27]. Hence, after each edge insertion/deletion, this information needs to be updated. There are several (synchronous and asynchronous) blocks in the online approach. And it is not trivial to obtain an efficient parallelization of the incremental algorithm. As our experiments will show, the dataflow programming model and pipelined parallelism are very useful to achieve a significant overlap among these computation/communication blocks and yield a scalable solution for the incremental centrality computation.

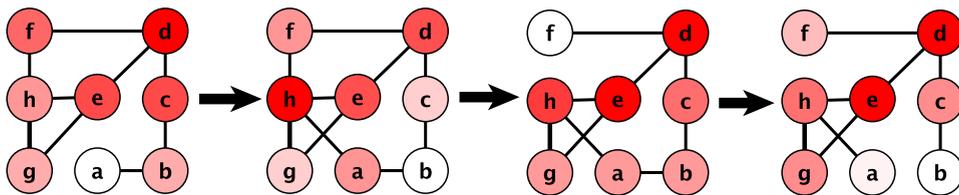
In this paper, we extend STREAMER that we introduced in [29]. Our contributions in [29] can be summarized as follows:

1. We proposed the first distributed-memory framework STREAMER for the incremental closeness centrality computation problem which employs pipelined parallelism to achieve computation–computation and computation–communication overlap [29].
2. The worker nodes we used in the experiments have 8 cores. In addition to the distributed-memory parallelization, we also leveraged the shared-memory parallelization and take NUMA effects into account [29].
3. The framework scales linearly: when 63 worker nodes (8 cores/node) are used, STREAMER obtains almost linear speedups compared to a single worker node–single thread execution [29].

In addition to above contributions, we introduce new extensions in this paper as follows:

1. The STREAMER framework is modular which makes it easily extendable. When the number of used nodes increases, the computation inevitably reaches a bottleneck on the extremities of the analysis pipeline which are not parallel. In [29], this effect appeared on one of the graph (*web-NotreDame*). Here, we show how the computation can be made parallel by leveraging the modularity of dataflow middleware.
2. Using an SpMM-based BFS formulation, we significantly improved the incremental CC computation performance and show that the dataflow programming model makes STREAMER highly modular and easy to enhance with novel algorithmic techniques.
3. These new techniques provide an improvement of a factor between 2.2 and 9.3 times compared to the techniques presented in [29].

The paper is organized as follows: Section 2 introduces the notation, formally defines the closeness centrality metric, and describes the incremental approach in [27]. Section 3 presents DataCutter [3], our in-house distributed memory dataflow middleware leveraged in this work. Section 4 describes the proposed distributed framework for incremental centrality computations in detail. The experimental analysis is given in Section 5, and Section 6 concludes the paper.



**Fig. 1.** A toy network with eight vertices and three consecutive edge (ah, fh, and ab, respectively) insertions/deletions. The vertices are colored with respect to their relative CC scores where red implies a higher closeness score. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

## 2. Incremental closeness centrality

Let  $G = (V, E)$  be a network modeled as a simple undirected graph with  $n = |V|$  vertices and  $m = |E|$  edges where each node is represented by a vertex in  $V$ , and a node–node interaction is represented by an edge in  $E$ . Let  $\Gamma_G(v)$  be the set of vertices which share an edge with  $v$ .

A graph  $G' = (V', E')$  is a *subgraph* of  $G$  if  $V' \subseteq V$  and  $E' \subseteq E$ . A *path* is a sequence of vertices such that there exists an edge between consecutive vertices. Two vertices  $u, v \in V$  are *connected* if there is a path from  $u$  to  $v$ . If all vertex pairs are connected we say that  $G$  is *connected*. If  $G$  is not connected, then it is *disconnected* and each maximal connected subgraph of  $G$  is a *connected component*, or a component, of  $G$ . We use  $d_G(u, v)$  to denote the length of the shortest path between two vertices  $u, v$  in a graph  $G$ . If  $u = v$  then  $d_G(u, v) = 0$ . If  $u$  and  $v$  are not connected  $d_G(u, v) = \infty$ .

Given a graph  $G = (V, E)$ , a vertex  $v \in V$  is called an *articulation vertex* if the graph  $G - v$  has more connected components than  $G$ .  $G$  is *biconnected* if it is connected and it does not contain an articulation vertex. A maximal biconnected subgraph of  $G$  is a *biconnected component*.

### 2.1. Closeness centrality

The *farness* of a vertex  $u \in V$  in a graph  $G = (V, E)$  is defined as  $\text{far}[u] = \sum_{\substack{v \in V \\ d_G(u, v) \neq \infty}} d_G(u, v)$ . And the closeness centrality of  $u$  is defined as  $\text{cc}[u] = \frac{|V|}{\text{far}[u]}$ . If  $u$  cannot reach any vertex in the graph, then  $\text{cc}[u] = 0$ . (All the notations are also given in Table 1).

For a graph  $G = (V, E)$  with  $n$  vertices and  $m$  edges, the complexity of the best  $\text{cc}$  algorithm is  $\mathcal{O}(n(m + n))$  (Algorithm 1). For each vertex  $s \in V$ , it executes a Single-Source Shortest Paths (SSSP), i.e., initiates a breadth-first search (BFS) from  $s$  and computes the distances to the connected vertices. As the last step, it computes  $\text{cc}[s]$ . Since a BFS takes  $\mathcal{O}(m + n)$  time, and  $n$  SSSPs are required in total, the complexity follows.

---

#### Algorithm 1. Offline centrality computation

---

**Data:**  $G = (V, E)$

**Out:**  $\text{cc}[\cdot]$

**for each**  $s \in V$  **do**

    ▷SSSP( $G, s$ ) with centrality computation

$Q \leftarrow$  empty queue

$d[v] \leftarrow \infty, \forall v \in V \setminus \{s\}$

$Q.\text{push}(s), d[s] \leftarrow 0$

$\text{far}[s] \leftarrow 0$

**while**  $Q$  is not empty **do**

$v \leftarrow Q.\text{pop}()$

**for all**  $w \in \Gamma_G(v)$  **do**

**if**  $d[w] = \infty$  **then**

$Q.\text{push}(w)$

$d[w] \leftarrow d[v] + 1$

$\text{far}[s] \leftarrow \text{far}[s] + d[w]$

$\text{cc}[s] = \frac{|V|}{\text{far}[s]}$

**return**  $\text{cc}[\cdot]$

---

### 2.2. Incremental closeness centrality

Algorithm 1 is an offline algorithm: it computes the CC scores from scratch. But today's networks are dynamic and their topologies are changing through time. Centrality computation is an expensive task, and especially for large scale networks,

**Table 1**  
Notations.

Notation	Meaning
$G$	Graph
$n$	Number of vertices
$m$	Number of edges
$\Gamma_G(v)$	Neighbors of $v$
$d_G(u, v)$	Length of the shortest path between $u$ and $v$
$\text{far}[u]$	$\sum_{\substack{v \in V \\ d_G(u, v) \neq \infty}} d_G(u, v)$
$\text{cc}[u]$	$\frac{n}{\text{far}[u]}$
$\Pi(uv)$	Biconnected component edge $uv$ belongs to

an offline algorithm cannot cope with the changing network topology. Hence, especially for large-scale, dynamic networks, online algorithms which do not perform the computation from scratch but only update the required scores in an incremental fashion are required. In a previous study, we used a set of techniques such as *level-based work filtering* and *special-vertex utilization* to reduce the centrality computation time for dynamic networks [27]. Here we remind them briefly and direct the interested readers to [27] for proofs and details.

### 2.3. Level-based work filtering

The level-based filtering aims to reduce the number of SSSPs in Algorithm 1. Let  $G = (V, E)$  be the current graph and  $uv$  be an edge to be inserted. Let  $G' = (V, E \cup \{uv\})$  be the modified graph. The centrality definition implies that for a vertex  $s \in V$ , if  $d_G(s, t) = d_{G'}(s, t)$  for all  $t \in V$  then  $cc[s] = cc'[s]$ . The following theorem is used to filter the SSSPs of such vertices.

**Theorem 2.1** (Saryüce et al. [27]). *Let  $G = (V, E)$  be a graph and  $u$  and  $v$  be two vertices in  $V$  s.t.  $uv \notin E$ . Let  $G' = (V, E \cup \{uv\})$ . Then  $cc[s] = cc'[s]$  if and only if  $|d_G(s, u) - d_G(s, v)| \leq 1$ .*

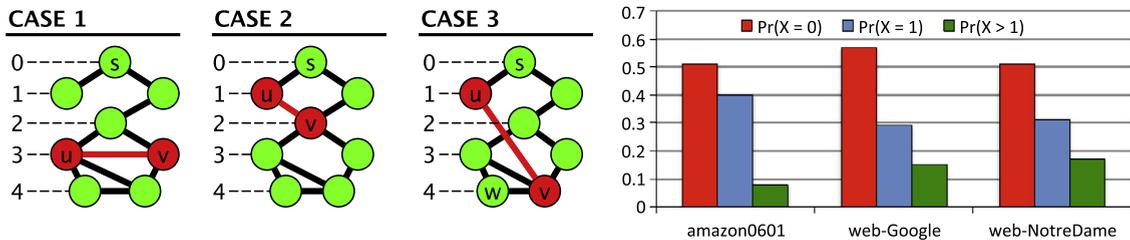
Many interesting real-life networks are scale free. The diameters of a scale-free network is small, and when the graph is modified with minor modifications, it tends to stay small. These networks also obey the power-law degree distribution. The level-based work filter is particularly efficient on these kind of networks. Fig. 2 (left) shows the three cases while an edge  $uv \in E$  is being added to  $G$ :  $d_G(s, u) = d_G(s, v)$ ,  $|d_G(s, u) - d_G(s, v)| = 1$ , and  $|d_G(s, u) - d_G(s, v)| > 1$ . Due to Theorem 2.1, an SSSP is required in Algorithm 1 only for the last case, since for the first two cases, the closeness centrality of  $s$  does not change. As Fig. 2 (right) shows, the probability of the last case is less than 20% for three social networks used in the experiments. Hence, more than 80% of the SSSPs are avoided by using level-based filtering.

Although Theorem 2.1 yields to a filter only in case of edge insertions, the same idea can easily be used for edge deletions. When an edge  $uv$  is inserted/deleted, to employ the filter, we first compute the distances from  $u$  and  $v$  to all other vertices. Detailed explanation can be found in [27].

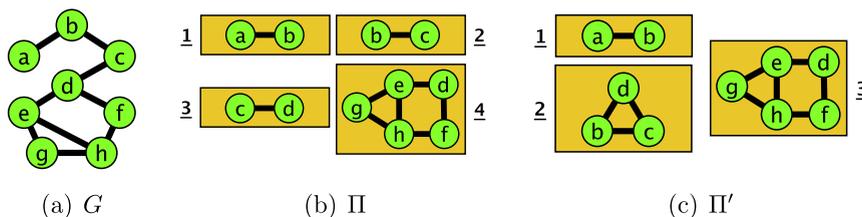
### 2.4. Special-vertex utilization

The work filter can be assisted by employing and maintaining a biconnected component decomposition (BCD) of  $G$ . A BCD is a partitioning  $\Pi$  of the edge set  $E$  where  $\Pi(e)$  is the component of each edge  $e \in E$ . A toy graph and its BCDs before and after an edge insertion are given in Fig. 3.

Let  $uv$  be the edge inserted to  $G = (V, E)$  and the final graph be  $G' = (V, E' = E \cup \{uv\})$ . Let  $f_{ar}$  and  $f_{ar}'$  be the fairness scores of all the vertices in  $G$  and  $G'$ . If the intersection  $\{\Pi(uw) : w \in \Gamma_G(u)\} \cap \{\Pi(vw) : w \in \Gamma_G(v)\}$  is not empty, there must be only one element in it (otherwise  $\Pi$  is not a valid BCD),  $cid$ , which is the id of the biconnected component of  $G'$  containing



**Fig. 2.** Three possible cases when inserting  $uv$ : for each vertex  $s$ , one of the following is true: (1)  $d_G(s, u) = d_G(s, v)$ , (2)  $|d_G(s, u) - d_G(s, v)| = 1$ , or (3)  $|d_G(s, u) - d_G(s, v)| > 1$  (left). The bars show the distribution of random variable  $X = |d_G(w, u) - d_G(w, v)|$  into three cases while an edge  $uv$  is being added to  $G$  (right). For each network, the probabilities are computed by using 1000 random edges from  $E$ . For each edge  $uv$ , we constructed the graph  $G = (V, E \setminus \{uv\})$  by removing  $uv$  from the final graph and computed  $|d_G(s, u) - d_G(s, v)|$  for all  $s \in V$ .



**Fig. 3.** A graph  $G$  (left), its biconnected component decomposition  $\Pi$  (middle), and the updated  $\Pi'$  after the edge  $bd$  is inserted (right). The articulation vertices before and after the edge insertion are  $\{b, c, d\}$  and  $\{b, d\}$ , respectively. After the addition, the second component contains the new edge, i.e.,  $cid = 2$ . This component is extracted first, and the algorithm performs updates only for its vertices  $\{b, c, d\}$ . It also initiates a fixing phase to make the CC scores correct for the rest of the vertices.

$uv$ . In this case, updating the BCD is simple:  $\Pi'(e)$  is set to  $\Pi(e)$  for all  $e \in E$  and  $\Pi'(uv)$  is set to  $cid$ . If the intersection is empty (see the addition of  $bd$  in Fig. 3(b)), we construct  $\Pi'$  from scratch and set  $cid = \Pi'(uv)$  (e.g.,  $cid = 2$  in Fig. 3(b)). A BCD can be computed in linear,  $\mathcal{O}(m+n)$  time [14]. Hence, the cost of BCD maintenance is negligible compared to the cost of updating closeness centrality.

Let  $G'_{cid} = (V_{cid}, E'_{cid})$  be the biconnected component of  $G'$  containing  $uv$ . Let  $\mathcal{A}_{cid} \subseteq V_{cid}$  be the set of articulation vertices of  $G'$  in  $G'_{cid}$ . Given  $\Pi'$ , it is easy to find the articulation vertices since  $u \in V$  is an articulation vertex if and only if it is at least in two components in the BCD:  $|\{\Pi'(uw) : uw \in E'\}| > 1$ .

The incremental algorithm executes SSSPs only for the vertices in  $G'_{cid}$ . The contributions of the vertices in  $V \setminus V_{cid}$  are integrated to the SSSPs through their *representatives*,  $rep : V \rightarrow V_{cid} \cup \{\text{null}\}$ . For a vertex in  $V_{cid}$ , the representative is itself. And for a vertex  $v \in V \setminus V_{cid}$ , the representative is either an articulation vertex in  $\mathcal{A}_{cid}$  or  $\text{null}$  if  $v$  and the vertices of  $V_{cid}$  are disconnected. Also, for all vertices  $x \in V \setminus V_{cid}$ , we have  $\text{far}'[x] = \text{far}[x] + \text{far}'[rep(x)] - \text{far}[rep(x)]$ . Therefore, there is no need to execute SSSPs from these vertices. Detailed explanation and proofs are omitted for brevity and can be found in [27].

In addition to articulation vertices, we exploit the *identical* vertices which have the same/a similar neighborhood structure to further reduce the number of SSSPs. In a graph  $G$ , two vertices  $u$  and  $v$  are *type-I-identical* if and only if  $\Gamma_G(u) = \Gamma_G(v)$ . In addition, two vertices  $u$  and  $v$  are *type-II-identical* if and only if  $\{u\} \cup \Gamma_G(u) = \{v\} \cup \Gamma_G(v)$ . Let  $u, v \in V$  be two identical vertices. One can easily see that for any vertex  $w \in V \setminus \{u, v\}$ ,  $d_G(u, w) = d_G(v, w)$ . Therefore, if  $\mathcal{I} \subseteq V$  is a set of (type-I or type-II) identical vertices, then the CC scores of all the vertices in  $\mathcal{I}$  are equal.

We maintain the sets of identical vertices and while updating the CC scores of the vertices in  $V$ , we execute an SSSP for a *representative* vertex from each identical-vertex set. We then use the computed score as the CC score of the other vertices in the same set. The filtering is straightforward and the modifications on the algorithm are minor. When an edge  $uv$  is added/removed to/from  $G$ , to maintain the identical vertex sets, we first subtract  $u$  and  $v$  from their sets and insert them to new ones. Candidates for being identical vertices are found using a hash function and the overall cost of maintaining the data structure is  $\mathcal{O}(n+m)$  [27].

## 2.5. Simultaneous source traversal

The performance of sparse kernels is mostly hindered by irregular memory accesses. The most famous example for sparse computation is the multiplication of a sparse matrix by a dense vector (SpMV). Several techniques, like register blocking [6,33] and usage of different matrix storage formats [2,21], are proposed to regularize the memory access pattern. However, multiplying a sparse matrix by multiple vectors is the most efficient and popular technique to regularize the memory access pattern. Once the multiple vectors are organized as a dense matrix, the problem becomes the multiplication of a sparse matrix by a dense matrix (SpMM). Each nonzero of the sparse matrix causes the multiplication of a single element of the vector in SpMV, and it results in the multiplications of as many consecutive elements of the dense matrix as its number of columns in SpMM.

Accommodating that idea for closeness centrality computation turns out to be concurrently computing the multiple sources at the same time. However, as opposed to SpMV, in which the vector is dense and therefore each non-zero induces exactly one multiplication, in BFS, not all the non-zeros will induce operations. That is to say, a vertex in BFS may or may not be traversed depending on which level is currently being processed. Thus, the traditional queue-based implementation of BFS does not seem to be easily extendable to support concurrent BFSs (co-BFS) in a vector-friendly manner. We developed this method in [30,31] and present here the main idea.

### 2.5.1. An SpMV-based formulation of closeness centrality

The idea is to convert to a simpler definition of level synchronous BFS: if one of the neighbor of  $v$  is part of level  $\ell - 1$  and  $v$  is not part of any level  $\ell' < \ell$ , then vertex  $v$  is part of level  $\ell$ . This formulation is used in parallel implementations of BFS on GPU [15,25,32], on shared memory systems [1] and on distributed memory systems [5].

The algorithm is better represented using binary variables. Let  $x_i^\ell$  be the binary variable that is `true` if vertex  $i$  is part of the frontier at level  $\ell$  for a BFS. The neighbors of level  $\ell$  is represented by a vector  $y^{\ell+1}$  computed by  $y_k^{\ell+1} = \text{OR}_{j \in \Gamma(k)} x_j^\ell$ . The next level is then computed with  $x_i^{\ell+1} = y_i^{\ell+1} \text{ AND not } (\text{OR}_{\ell' < \ell} x_i^{\ell'})$ . Using these variables, one can increase the farness of the source by  $\ell$  if  $i$  is at level  $\ell$  (i.e., if  $x_i^\ell = 1$ ). One can remark that  $y^{\ell+1}$  is the result of the “multiplication” of the adjacency matrix of the graph by  $x^\ell$  in the (OR,AND) semi-ring.

### 2.5.2. An SpMM-based formulation of closeness centrality

It is easy to derive an algorithm from the formulation given above for closeness centrality computation that processes multiple sources concurrently. Instead of manipulating a single vector  $x$  and  $y$  where each element is a single bit, one can encode 32-bit vectors for 32 BFSs so that one *int* can encode the state of a single vertex across the 32 BFSs. The algorithm becomes quite efficient as it does not use more memory and process 32 BFS concurrently. All the operations become simple bit-wise `and`, `or` and `not`.

Theoretically, the asymptotic complexity changes when BFS is implemented using an SpMM approach. The complexity of the traditional queue-based BFS algorithm is  $\mathcal{O}(|E|)$ . If the adjacency matrix is stored row-wise, the SpMM-based implementation boils down to a bottom-up implementation of BFS which has a natural write access pattern. However, it becomes

impossible to only traverse the relevant nonzero of the matrix and the complexity of the algorithm becomes  $\mathcal{O}(|E| \times L)$ , where  $L$  is the diameter of the graph. Social networks have small world properties which implies that their diameter is low and we do not feel that this asymptotic factor of  $L$  will hinder performance.

Moreover, multiple BFSs are performed concurrently (here 32) which can recoup for the loss. In [30,31], the algorithm computes the impact of the sources on all the vertices of the graph. What we presented in this section does the reverse and compute the impact of all the vertices of the graph on the sources. Despite worse asymptotic complexity the performance of co-BFS outperforms traditional BFS approach [30,31]. Moreover, such algorithm is compatible with the decomposition of the graph in biconnected components [28] which can lead to further improvement. Because this algorithm computes the farness of the sources, it can be used to compute centrality incrementally.

### 3. DataCutter

STREAMER employs *DataCutter* [3], our in-house dataflow programming framework for distributed memory systems. In *DataCutter*, the computations are carried by independent computing elements, called *filters*, that have different responsibilities and operate on data passing through them. *DataCutter* follows the component-based programming paradigm which has been used to describe and implement complex applications [11–13,29] by way of components – distinct tasks with well-defined interfaces. This is also known as the filter-stream programming model [3] (a specific implementation of the dataflow programming model). A *stream* denotes a uni-directional data flow from some filters (i.e., the producers) to others (i.e., the consumers). Data flows along these *streams* in untyped *databuffers* so as to minimize various system overheads. A *layout* is a filter ontology which describes the set of application tasks, streams, and the connections required for the computation. By describing these components and the explicit data connections between them, the applications are decomposed along natural task boundaries according to the application domain. Therefore, the component-based application design is an intuitive process with explicit demarcation of task responsibilities. Furthermore, the communication patterns are also explicit; each component includes its input data requirements and outputs in its description.

Applications composed of a number of individual tasks can be executed on parallel and distributed computing resources and gain extra performance over those run on strictly sequential machines. This is achieved by specifying a *placement* which is an instance of a *layout* with a mapping of the filters onto physical processors. There are three main advantages of this scheme: first, it exposes an abstract representation of the application which is decoupled from its practical implementation. Second, the coarse-grain dataflow programming model allows *replicated parallelism* by instantiating a given filter multiple times so that the work can be distributed among the instances to improve the parallelism of the application and the system's performance. And third, the execution is pipelined, allowing multiple filters to compute simultaneously on different iterations of the work. This *pipelined parallelism* is very useful to achieve overlapping of communication and computation.

Additionally, provided the interfaces exposed by a task to the rest of the application, different implementations of tasks, possibly on different processor architectures can co-exist in the same application deployment, allowing developers to take full advantage of modern, heterogeneous supercomputers. Fig. 4 shows an example filter-stream layout and placement. In this work, we used both distributed- and shared-memory architectures. However, thanks to filter-stream programming model, many-core systems such as GPUs and accelerators can also be used easily and efficiently if desired [13].

As mentioned above, one of the *DataCutter*'s strengths is that it enables pipelined parallelism, where multiple stages of the pipeline (such as A and B in the layout in Fig. 4) can be executed simultaneously, and replicated parallelism can be used at the same time if some computation is stateless (such as filter B in the same figure). *DataCutter* makes all this parallelism possible by mapping each placed filter to a POSIX thread of the execution platform.

### 4. STREAMER

STREAMER is implemented in the *DataCutter* framework. We propose to use the four-filter layout shown in Fig. 5. *InstanceGenerator* is responsible for sending the updates to all the other components. *StreamingMaster* does the work filtering for each update, explained in Section 2, and generates the workload for following components. *ComputeCC* component executes the real work and computes the updated CC scores for each incoming update. *Aggregator* does the necessary adjustments

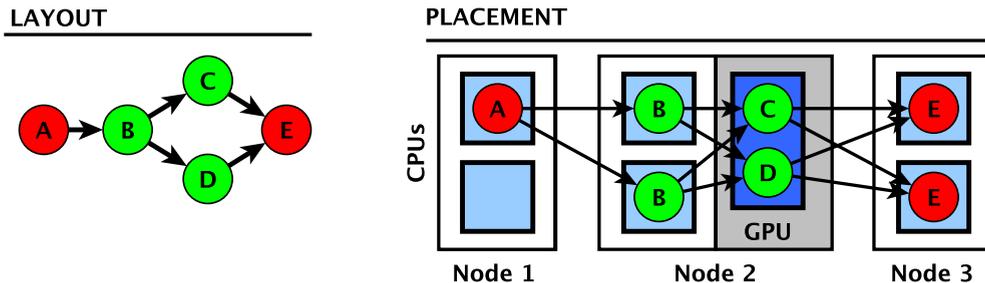


Fig. 4. A toy filter-stream application layout and its placement.

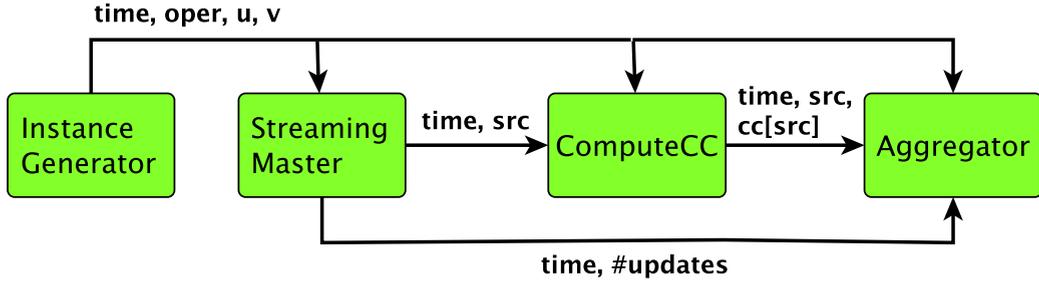


Fig. 5. Layout of STREAMER.

related to identical vertex sets and biconnected component decomposition. While computing the CC scores, the main portion of the computation comes from performing SSSPs for the vertices whose scores need to be updated. If there are many updates (we use the term “update” to refer to the SSSP operation which updates the CC score of a vertex), that part of the computation should occupy most of the machine. A typical synchronous decomposition of the application makes the work filtering of a streaming event (handling a single edge change) wait for the completion of all the work incurred by a previous streaming event. Since the worker nodes will wait for the work filtering to be completed, there can be a large waste of resources. We argue that the pipelined parallelism should be used to overlap the process of filtering the work and computing the updates on the graph. In this section, we explain each component in detail and define their responsibilities.

The first filter is the *InstanceGenerator* which first sends the initial graph to all the other filters. It then sends the streaming Events as 4-tuples  $(t, oper, u, v)$  to indicate that edge  $uv$  has been either added or removed (specified by *oper*) at a given time  $t$ . (In the following, we only explain the system for edge insertion, but it is essentially the same for an edge removal.) In a real world application, this filter would be listening on the network or on a database trigger for topology modifications; but in our experiments, all the necessary information is read from a file.

*StreamingMaster* is responsible for the work filtering after each network modification. Upon inserting  $uv$  at time  $t$ , it first computes the shortest distances from  $u$  and  $v$  to all other vertices at time  $t - 1$ . Then, it adds the edge  $uv$  into its local copy of the graph and updates the identical vertex sets as described in Section 2.4. It partitions the edges of the graph to its biconnected components by using the algorithm in [14] and finds the component containing  $uv$ . For each vertex  $s \in V$ , it decides whether its CC score needs to be recomputed by checking the following conditions: (1)  $d(s, u)$  and  $d(s, v)$  differ by at least 2 units at time  $t - 1$ , (2)  $s$  is adjacent to an edge which is also in  $uv$ 's biconnected component, (3)  $s$  is the representative of its identical vertex set. *StreamingMaster* then informs the *Aggregator* about the number of updates it will receive for time  $t$ . Finally, it sends the list of SSSP requests to the *ComputeCC* filter, i.e., the corresponding source vertex ids whose CC scores need to be updated.

*ComputeCC* performs the real work and computes the new CC scores after each graph modification. It waits for work from *StreamingMaster*, and when it receives a CC update request in the form of a 2-tuple  $(t, s)$  (update time and source vertex id), *ComputeCC* advances its local graph representation to time  $t$  by using the appropriate updates from *InstanceGenerator*. If there is a change on the local graph, the biconnected component of  $uv$  is extracted, and a concise information of the graph structure and the set of articulation vertices are updated (as described in [27]). Finally, the exact CC score  $cc[s]$  at time  $t$  is computed and sent to the *Aggregator* as a 3-tuple  $(t, s, cc[s])$ . *ComputeCC* can be replicated to fill up the whole distributed memory machine without any problem: as long as a replica reads the update requests in the order of non-decreasing time units, it is able compute the correct CC scores.

The *Aggregator* filter gets the graph at a time  $t$  from *InstanceGenerator*. Then, it obtains the number of updates for that time from *StreamingMaster*. It computes the identical vertex sets as well as the BCD. It gets the updated CC scores from *ComputeCC*. Due to the pipelined parallelism used in the system and the replicated parallelism of *ComputeCC*, it is possible that updates from a later time can be received; STREAMER stores them in a backlog for future processing. When a  $(t, s, cc[s])$  tuple is processed, the CC score of  $s$  is updated. If  $s$  is the representative of an identical vertex set, the CC scores of all the vertices in the same set are updated as well. If  $s$  is an articulation point, then the CC scores of the vertices which are represented by  $s$  (and are not in the biconnected component of  $uv$ ) are updated as well, by using the difference in the CC score of  $s$  between time  $t$  and  $t - 1$ . Since *Aggregator* needs to know the CC scores at time  $t - 1$  to compute the centrality scores at time  $t$ , the system must be bootstrapped: the system computes explicitly all the centrality scores of the vertices for time  $t = 0$ .

#### 4.1. Exploiting the shared memory architecture

The main portion of the execution time is spent by the *ComputeCC* filter. Therefore, it is important to replicate this filter as much as possible. Each replica of the filter will end up maintaining its own graph structure and computing its own BCD. Modern clusters are hierarchical and composed of distributed memory nodes where each node contains multiple processors featuring multiple cores that share the same memory space. For instance, the nodes used in our experiments are equipped with two processors, each having 4 cores.

It is a waste of computational power to recompute the data structure on each core. But it is also a waste of memory. Indeed, the cores of a processor typically share a common last level of cache and using the same memory space for all

the cores in a processor might improve the cache utilization. We propose to split the *ComputeCC* filter in two separate filters which are transparent to the rest of the system thanks to *DataCutter* being component-based. The *Preparator* filter constructs the decomposed graph for each Streaming Event it is responsible for. The *Executor* filter performs the real work on the decomposed graph. In *DataCutter*, the filters running on the same physical node act run in separate pthreads within the same MPI process making sharing the memory as easy as communicating pointers. The release of the memory associated with the decomposed graph is handled by atomically decreasing a reference counter by the *Executor*.

The decoupling of the graph management and the CC score computation allows to either creating a single graph representation on each distributed memory node or having a copy of the graph on each NUMA domain of the architecture. This is shown in Fig. 6.

#### 4.2. Parallelizing StreamingMaster

When the number of cores used for *ComputeCC* increases, the relative importance of *ComputeCC* in the total runtime decreases. Theoretically, with an infinite number of cores for *ComputeCC*, the time required by it will drop to zero. In this case, the bottleneck of the application becomes the maximum rate at which *StreamingMaster* can generate updates request and the rate at which *Aggregator* can merge the computed results. To improve these rates, we replace them with a construct that allow parallel execution.

*StreamingMaster* is decomposed in three filters which are laid out according to Fig. 7. Most of the work done by *StreamingMaster* is done by a filter (we still call it *StreamingMaster* for convenience) which supports replication. Each of the replica

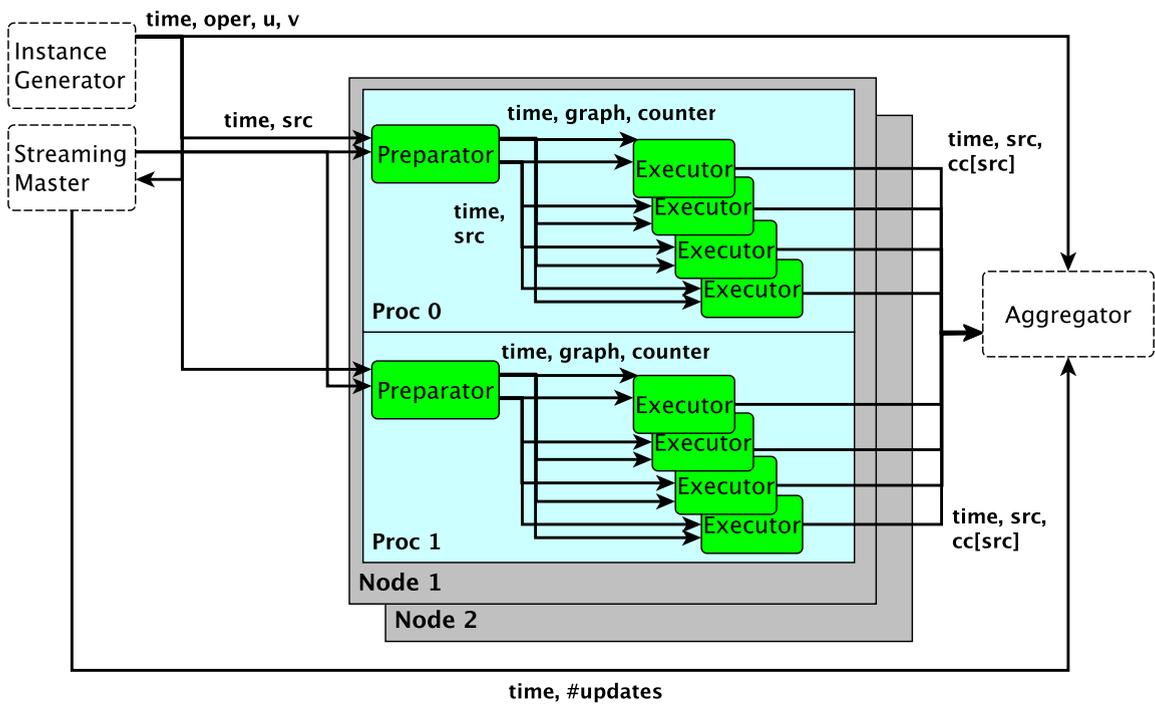


Fig. 6. Placement of STREAMER using 2 worker nodes with 2 quad-core processors. (The node 2 is hidden). The remaining filters are on node 0.

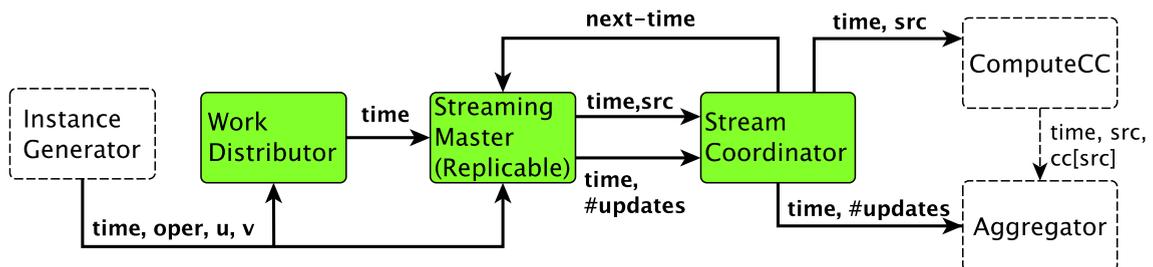


Fig. 7. Replicating StreamingMaster for a better scaling when the number of processors is large.

receives the list of edges it has to compute the filtering from a *WorkDistributor*. This *WorkDistributor* just listens the modifications on the graph and distribute the Streaming Events among different *StreamingMasters*.

It is important that *ComputeCC* receives the update requests in non-decreasing order of Streaming Events. *StreamCoordinator* is responsible for enforcing that order. *StreamCoordinator* sits between the *StreamingMaster* and the *ComputeCC* (and the *Aggregator*) and relays messages to them. The *StreamCoordinator* tells *StreamingMaster* which streaming event is the next one. In other words, before outputting the list of updates (and metadata for the *Aggregator*), the *StreamingMaster* reads from the *StreamCoordinator* whether it is time to output.

### 4.3. Parallelizing Aggregator

One of the challenges in parallelizing the *Aggregator* is that there can be only one filter that actually stores the centrality values of the network. Fortunately, most of the computation time spent by the *Aggregator* is spent in preparing the network rather than in applying the updates. We modify the layout of the *Aggregator* to match that of Fig. 8.

Therefore only a single filter, we will call *Aggregator* for the sake of simplicity, is responsible for applying the updates, and is only responsible for this. It takes three kinds of input: the updates on the graph itself, the information of how many updates will be applied for each streaming event and information on the graph (the graph itself, its biconnected decomposition and identical vertices).

The graph information is constructed by another filter called *AggregatorPreparator* which can be replicated. It listens to the Streaming Events and receive work assignments. It then computes the sets of identical vertices and the graph's biconnected component decomposition and send them through its downstream.

The work in the *AggregatorPreparator* is distributed in a way similar to the parallelization of the *StreamingMaster*. Also the graph information must reach the *Aggregator* in the order of the Streaming Event. An *AggregatorCoordinator* is used to regulate the order in which the graph information is sent. It behaves under the same principle as *StreamCoordinator*.

## 5. Experiments

STREAMER runs on the *owens* cluster in the Department of Biomedical Informatics at The Ohio State University. For the experiments, we used all the 64 computational nodes, each with dual Intel Xeon E5520 Quad-core CPUs (with 2-way Simultaneous Multithreading, and 8 MB of L3 cache per processor), 48 GB of main memory. The nodes are interconnected with 20 Gbps InfiniBand. The algorithms were run on CentOS 6, and compiled with GCC 4.8.1 using the `-O3` optimization flag. DataCutter uses an InfiniBand-aware MPI to leverage the high performance interconnect: here we used MVAPICH2 2.0b.

For testing purposes, we picked 4 large social network graphs from the SNAP dataset to perform the tests at scale. The properties of the graphs are summarized in Table 2. For simulating the addition of the edges, we removed 50 edges from the graphs and added them back one by one. The streamed edges were selected using a uniform random distribution. For comparability purposes, all the runs performed on the same graph use the same set of edges. The number of updates induced

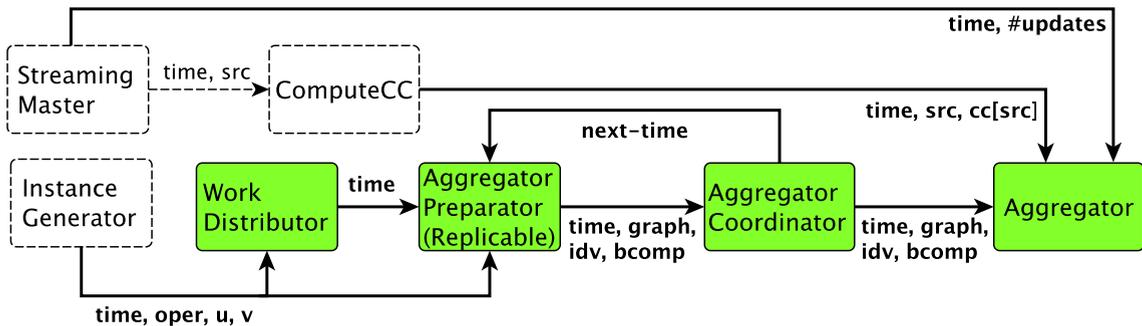


Fig. 8. Replicating Aggregator for a better scaling when the number of processors is large.

Table 2

Properties of the graphs we used in the experiments and execution time on a 64 node cluster.

Name	$ V $	$ E $	# Updates	Time (s)	Speedup w.r.t [27] seq. non-incremental	Speedup w.r.t [27] seq. incremental
web-NotreDame	325,729	1,090,008	399,420	3.29	43,237	805
amazon0601	403,394	2,443,308	1,548,288	33.16	22,471	449
web-Google	916,428	4,321,958	2,527,088	71.20	45,860	578
soc-pokec	1,632,804	30,622,464	4,924,759	816.73	–	–

by that set of edges when applying filtering using identical vertices, biconnected component decomposition, and level filtering is given in Table 2. In the experiments, the data comes from a file, and the Streaming Events are pushed to the system as quickly as possible so as to stress the system.

All the results presented in this section are extracted from a single run of STREAMER with proper parameters. As our preliminary results show, the regularity in the plots indicates there is a small variance on the runtimes, which induces a reasonable confidence in the significance of the quoted numbers. In the experiments, *StreamingMaster* and *Aggregator* run on the same node, apart from all the *ComputeCC* filters. Therefore, we report the number of worker nodes, but an extra node is always used.

To give an idea of the actual amount of computation, in the fourth column of Table 2, we report the time STREAMER spends to update the CC scores upon 50 edge insertions by using all 63 worker nodes. We also present the speedup of parallel implementation on 64 nodes with respect to sequential non-incremental computation and sequential incremental computation. The STREAMER framework is never sequential due to its distributed-memory nature and the pipelined parallelism, i.e., different filters are always handled by different threads even in the most basic setting with no filter replication. (STREAMER uses at least the four filters of Fig. 5, so at least four POSIX threads are always used.) Therefore, there is no sequential runtime for the STREAMER framework. When we mention the sequential time, it refers to our previous work [27], which runs sequentially using a single core of the same cluster. As all the execution times given in this section, the times in Table 2 do not contain the initialization time. That is the time measurement starts once STREAMER is idle, waiting to receive Streaming Events.

### 5.1. Basic performance results

Fig. 9 shows the performance and scalability of the system in different configurations with a single *StreamingMaster* and *Aggregator*. In general, the success of a streaming graph analytics operation is measured by the rate at which they can

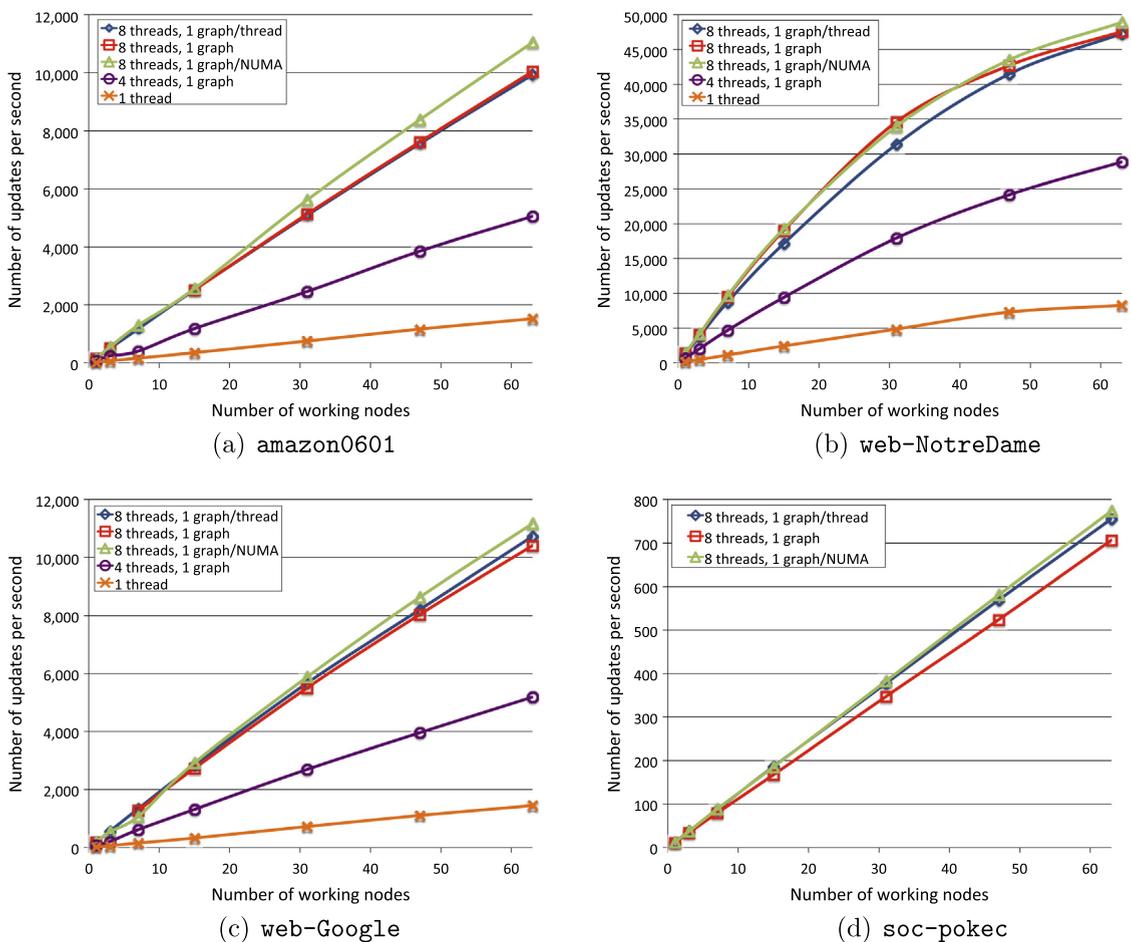


Fig. 9. Scalability: the performance is expressed in the number of updates per second. Different worker-node configurations are shown. “8 threads, 1 graph/thread” means that 8 *ComputeCC* filters are used per node. “8 threads, 1 graph” means that 1 *Preparator* and 8 *Executor* filters are used per node. “8 threads, 1 graph/NUMA” means that 2 *Preparators* per node (one per NUMA domain) and 8 *Executors* are used.

maintain the analytic. For this reason, the performance is expressed in number of (vertex's centrality) updates per second. The framework obtains up to 11,000 updates/s on `amazon0601` and `web-Google`, 49,000 updates/s on `web-NotreDame`, and more than 750 updates/s on the largest tested graph `soc-pokec`. It appears to scale linearly on the graphs `amazon0601` and `web-Google`, `soc-pokec`. For the first two graphs, it reaches a speedup of 456 and 497, respectively, with 63 worker nodes and 8 threads/node (504 *Executor* threads in total) compared to the single worker node-single thread configuration (the incremental centrality computation on `soc-pokec` with a single node and a single thread was too long to run the experiment, but the system is clearly scaling well on this graph). The last graph, `web-NotreDame`, does not exhibit a linear scaling and obtains a speedup of only 316.

Let us first evaluate the performance obtained under different node-level configurations. Table 3 presents the relative performance of the system using 31 worker nodes while using 1, 4, or 8 threads per node. When compared with the single thread configuration, using 4 threads (the second column) is more than 3 times faster, while using 8 threads (columns 3–5) per node usually gives a speedup of 6.5 or more. Overall, having multiple cores is fairly well exploited. Properly taking the shared-memory aspect of the architecture into account (column 5) brings a performance improvement between 1% to 10% (the last column). In one instance (`web-Google` with a graph for each NUMA domain), we observed that the normalized performance is more than the number of cores. This can be explained by the fact that actually 10 threads are running on each computing node (8 *Executor* and 2 *Preparator*) which can lead to a higher parallelism.

## 5.2. Execution-log analysis

Here we discuss the impact of pipelined parallelism and the sub-linear speedup achieved on `web-NotreDame`. In Fig. 10, we present the execution logs for that graph obtained while using 3, 15, and 63 worker nodes. Each log plot shows three data series: the times at which *StreamingMaster* starts to process the Streaming Events, the total number of updates sent by *StreamingMaster*, and the number of updates processed by the *Executors* collectively. The three different logs show what happens when the ratio of update produced and update consumed per second changes.

The first execution-log plot with 3 worker nodes (Fig. 10(a)) shows the amount of the updates emitted and processed as two perfectly parallel *almost straight* lines. This indicates that the runtime of the application is dominated by processing the updates. As the figure shows, the times at which the *StreamingMaster* starts processing the Streaming Events are not evenly distributed. As mentioned before, *StreamingMaster* starts filtering for the next Streaming Event as soon as it sends all the updates for the current one. In other words, the amount of updates emitted for a given Streaming Event can be read from the execution log as the difference of the y-coordinates of two consecutive “update emitted” points (the first line). In the first plot, we can see that 6 out of 50 Streaming Events (the ticks at the end of each partial tick-lines) incurred significantly more updates than the others. While these events are being processed, the two lines stay straight and parallel, because in DataCutter, writing to a downstream filter is a buffered operation. Once the buffer is full, the operation becomes blocking.

The second execution log with 15 worker nodes (Fig. 10(b)) shows a different picture. Here, the log is about 4 times shorter and the lines are not perfectly parallel. The number of updates emitted shows three plateaus for more than a second around times 0, 5, and 16 s. These plateaus exist because many consecutive Streaming Events do not generate a significant amount of updates; therefore, the *StreamingMaster* spends all its time by filtering the work for these Streaming Events.

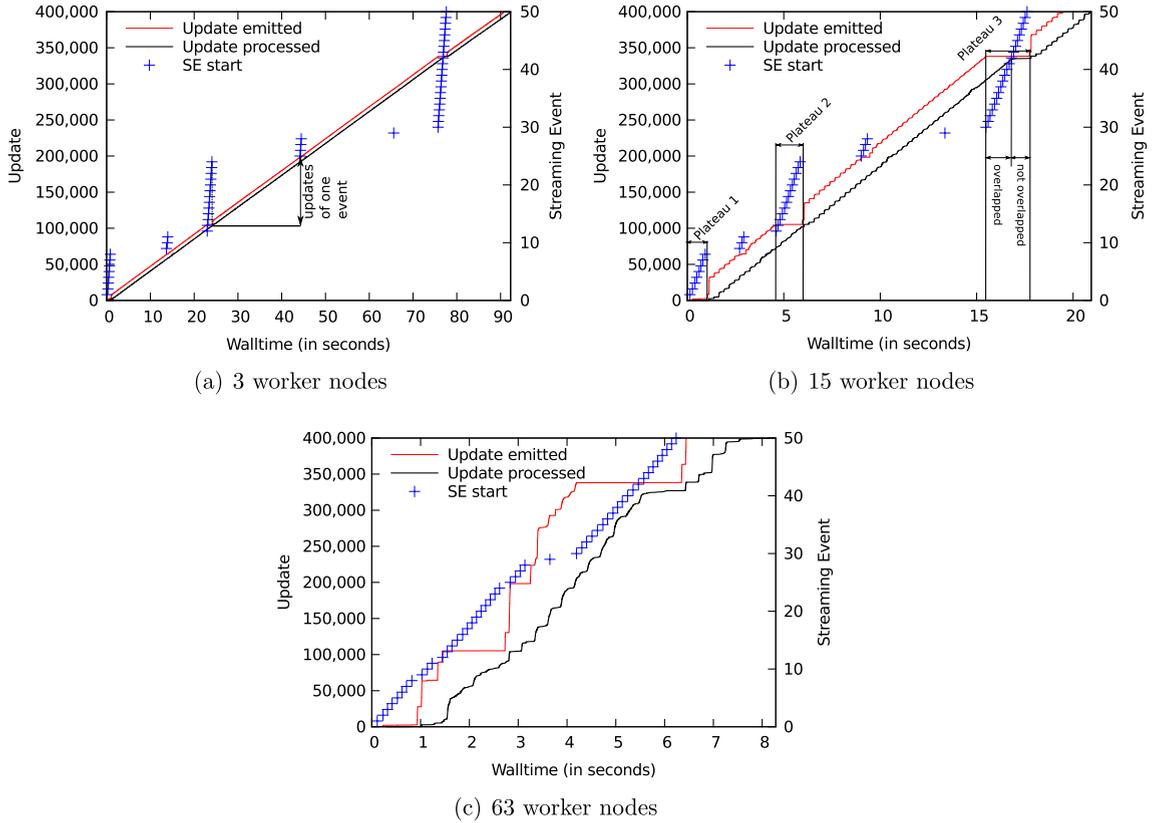
The second plateau is around time 5 s of the execution log with 15 worker nodes, it lasts 1.2 s, and less than 100 updates are sent during that interval. However, as the plot shows, the worker nodes do not run out of work and process more than 25,000 updates during the plateau. This is possible because the computation in *STREAMER* is pipelined. If the system were synchronous the worker nodes would spend most of that plateau waiting which yields a longer execution time and worse performance. In addition to the three large plateaus, cases with a few consecutive Streaming Events that lead to barely no updates are slightly visible around times 3 and 9. These two smaller cases are hidden by the pipelined parallelism. The third plateau is much longer than the second one (20 Streaming Events, 2.1 s) and the worker nodes eventually run out of work halfway through the plateau. As can be seen in Fig. 9(b), the performance does not show linear scaling at 15 worker nodes; but it is still good, thanks to the pipelined parallelism.

When 63 worker nodes are used, the execution log (Fig. 10(c)) presents another picture. With the increase on the workers' processing power, a single *StreamingMaster* is now the main bottleneck of the computation. Two additional, considerably large plateaus appeared, and *StreamingMaster* starts to spend more than half of its time with the work filtering. However,

**Table 3**

The performance of *STREAMER* with 31 worker nodes and different node-level configurations normalized to 1 thread case (performance on `soc-pokec` is normalized to 8 threads, 1 graph/thread). The last column is the advantage of shared memory awareness (ratio of columns 5 and 3).

Name	4 Threads	8 threads, 1 graph per			Shared Mem. awareness
		Thread	Node	NUMA	
<code>web-NotreDame</code>	3.69X	6.46X	7.13X	6.99X	1.08X
<code>amazon0601</code>	3.26X	6.75X	6.81X	7.45X	1.10X
<code>web-Google</code>	3.69X	7.77X	7.55X	8.06X	1.03X
<code>soc-pokec</code>	-	1.00X	0.92X	1.01X	1.01X



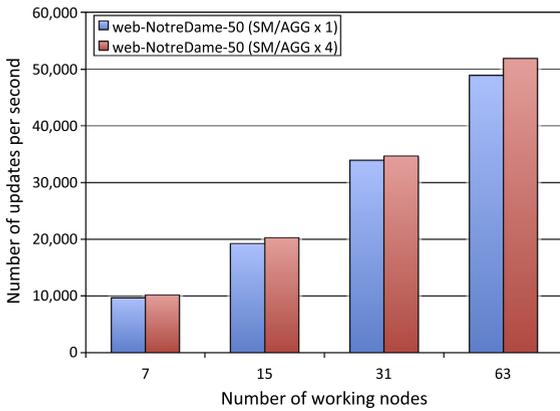
**Fig. 10.** Execution logs for `web-NotreDame` on different number of nodes. Each plot shows the total number of updates sent by *StreamingMaster* and processed by the *Executors*, respectively (the two lines), and the times at which *StreamingMaster* starts to process Streaming Events (the set of ticks).

during these times, the workers keep processing the updates, but at varied rates, due to temporary work starvation. The work filtering and the actual work are being processed mostly simultaneously showing that pipelined parallelism is very effective in this situation. Without the pipelined parallelism, the computation time would certainly be 2 s longer, and 25% worse performance would be achieved.

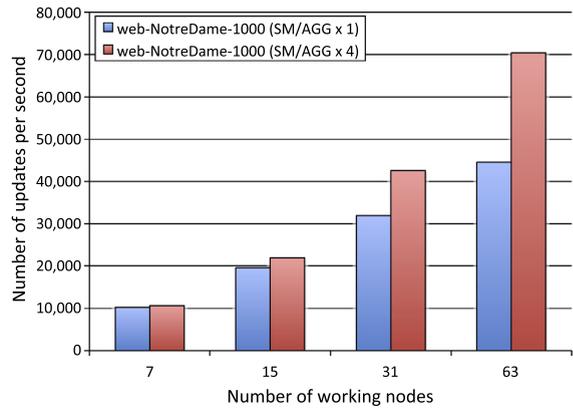
We used the techniques described in Sections 4.2 and 4.3 (Figs. 7 and 8) to replicate the *StreamingMaster* and *Aggregator* filters, respectively, and obtain a better performance when these filters becomes bottleneck throughout the incremental closeness centrality computation. The results on the `web-NotreDame` graph are given for 50 and 1000 Streaming Events in Fig. 11. As the figure shows, using four *StreamingMaster* and *Aggregator* filters instead of one yields around 6% improvement for 50 Streaming Events when 63 working nodes in the cluster are fully utilized. This small improvement is due to a lack of sufficient number of Streaming Events which generates a large amount of updates (see Fig. 10). Hence, even with a large number of *StreamingMaster* and *Aggregator* filters, due to the load balancing problem on these filters, one cannot improve the performance more with 50 Streaming Events by just replicating them. Fortunately, in practice this number is usually much higher. In Fig. 11(b), we repeated the same experiment for 1000 Streaming Events. As the figure shows, the performance significantly increases when the filters are replicated. Furthermore, the percentage of the improvement increases when more nodes are used and reaches to 58% with 63 working nodes. This is expected since, with more cores for the *Executor* filters, the time spent for *StreamingMaster* and *Aggregator* becomes (relatively) more important. When applied on the other graphs, going from one *StreamingMaster* and *Aggregator* to four have not yield significant difference since these components were not bottlenecks. Therefore, we omitted those results here.

### 5.3. Plug-and-play filters: *co-BFS*

As stated above, thanks to filter-stream programming model, different filter implementations and various hardware such as GPUs can be used easily and efficiently if desired. Here, we show that using the SpMM-based approach described in Section 2.5, one can modify the *ComputeCC* filter in Fig. 5 (or the *Executor* filters in Fig. 6) to increase the performance. For this experiment, we swapped the *Executor* filter with one that uses the *co-BFS* algorithm which computes 32 BFSs from different sources concurrently. The results of the experiments with 15, 31, and 63 working nodes are shown in Fig. 12. Using *co-BFS*

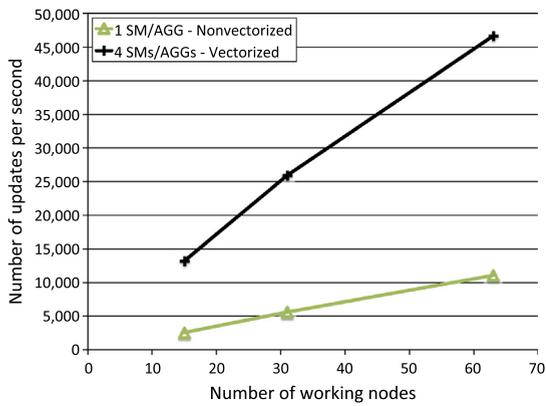


(a) 50 edge insertions on web-NotreDame

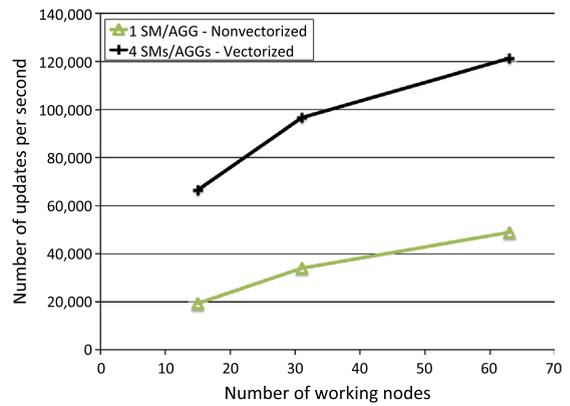


(b) 1,000 edge insertions on web-NotreDame

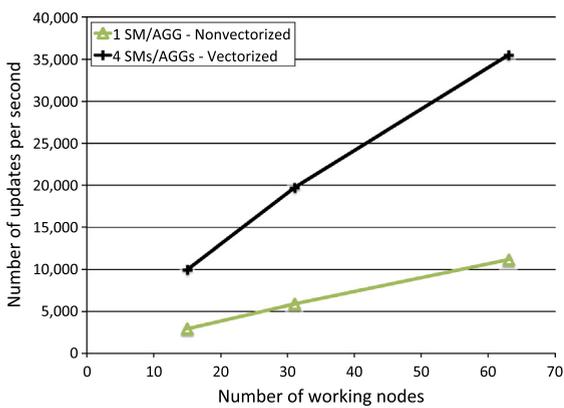
**Fig. 11.** Parallelizing *StreamingMaster* and *Aggregator*: the number of updates per second for *web-NotreDame* with 50 and 1000 Streaming Events, respectively. The best node configuration from Fig. 9, i.e., 8 threads, 1 graph/NUMA, is used for both cases.



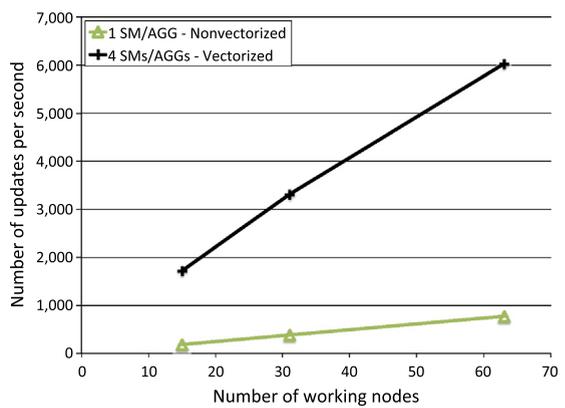
(a) amazon0601



(b) web-NotreDame



(c) web-Google



(d) soc-pokec

**Fig. 12.** co-BFS: the performance is expressed in the number of updates per second. The best worker-node configuration, “8 threads, 1 graph/NUMA”, is used for the experiments.

(and coupled with multiple *StreamingMaster* and *Aggregator*) improves the performance of the regular version by a factor ranging from 2.2 to 9.3 depending on the graph and number of working nodes.

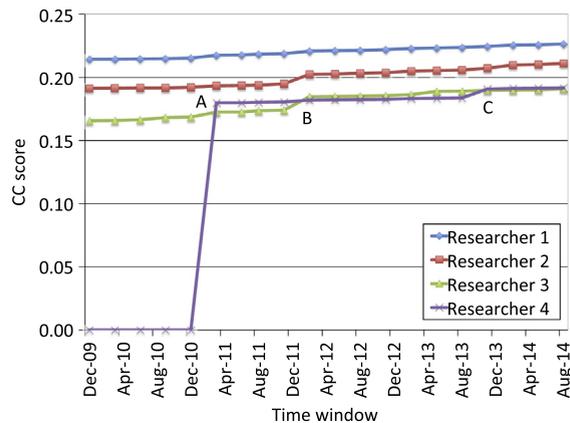


Fig. 13. Closeness centrality score evolution in DBLP coauthor network.

#### 5.4. Illustrative example for closeness centrality evolution

In this section, we present a real world example to show how the closeness centrality scores of four researchers change over time in the temporal coauthor network obtained from DBLP.<sup>1</sup> We selected the four authors of this manuscript which have different experiences and looked at their closeness centrality score evolution from December 2009 to August 2014. We report the closeness centrality scores at the end of every 3 months by our incremental algorithms.

Fig. 13 shows how the CC score changes when time passes. Researcher 4 is a PhD student who started in September 2010 and his first paper was published at the beginning of 2011. Point A shows the impact of the first paper on his CC score. Researcher 3 joined the team as a postdoc in September 2011. His first paper with the team members was published in early 2012 (Point B). We can observe that this publication increased his CC score, making him more central in the DBLP coauthor network. This publication also effected the centrality score of Researcher 2, who was another postdoc of the team at the time. Another significant point in the figure is point C, which corresponds to the publication of Researcher 4 as a result of his internship in an different institute. This publication made Researcher 4 more central since he is connected to new researchers in the DBLP coauthor network. Apart from those important milestones, we can see that there is a steady increase in CC scores of the four researchers.

#### 5.5. Summary of the experimental results

The experiments we conducted shows that *STREAMER* can scale up and efficiently utilize our entire experimental cluster. By taking the hierarchical composition of the architecture into account (64 nodes, 2 processors per node, 4 cores per processor) and not considering it as a regular distributed machine (a 512-processor MPI cluster), we enabled processing of larger graphs and obtained 10% additional improvement. Furthermore, the pipelined parallelism proved to be extremely necessary while using a large amount of nodes in a concurrent fashion.

Replicating the *ComputeCC* filter leads to significant speedup. Yet, the bottleneck eventually becomes the filters that cannot be replicated automatically. For filters where the ordering of the messages is important, we can substitute an alternative filter architecture to alleviate the bottleneck and make the whole analysis pipeline highly parallel.

The flexibility of the filter-stream programming model allows to easily substitute a component of the application by an alternative implementation. For instance, one can use modern vectorization techniques to improve the performance by a significant factor. Similarly, one could have an alternative implementation which use different type of hardware such as accelerators.

For the three of the graphs *web-NotreDame*, *amazonO601* and *web-Google*, a reference sequential time is known from [27] for both the non-incremental and the incremental cases. *STREAMER* using 63 worker nodes (8 cores per node), 4 *StreamingMaster* and 4 *Aggregators* and co-BFS computing filters improved the runtime of the incremental algorithm by a factor of 805, 449 and 578 respectively on the three graphs. Compared to a sequential non-incremental computation of the closeness centrality value, *STREAMER* improves the runtime by a factor ranging from 22,471 to 45,860. These numbers are reported in Table 2.

## 6. Conclusion

Maintaining the correctness of a graph analysis is important in today's dynamic networks. Computing the closeness centrality scores from scratch after each graph modification is prohibitive, and even sequential incremental algorithms are too expensive for networks of practical relevance. In this paper, we proposed *STREAMER*, a distributed memory framework which

<sup>1</sup> <http://dblp.uni-trier.de/xml/>

guarantees the correctness of the CC scores, exploits replicated and pipelined parallelism, and takes the hierarchical architecture of modern clusters into account.

The system is fully scalable as each of its components can be made to use an arbitrary number of nodes. Also, we showed that we can easily use alternative implementation of the BFS computations to allow the use of novel algorithmic techniques or hardware. Using STREAMER on a 64 nodes, 8 cores/node cluster, we reached almost linear speedup in the experiments and the performance are orders of magnitude higher than the non-incremental computation. Maintaining the closeness centrality of large and complex graph in real-time is now within our grasp.

As a future work, approximate closeness centrality and different path-based centrality measures can be computed in a STREAMER like framework. For the approximate closeness centrality computation, we just need to change the way closeness centrality score is computed. In the current STREAMER, we gather the centrality score of each vertex while computing SSSP from it. For the approximate closeness centrality, we need to scatter the centrality scores to each traversed vertex while applying the SSSP. For this purpose, we just need to modify StreamingMaster component. For other path-based centrality measures, general layout will remain same, but some components might need to be modified. For example, betweenness centrality computation might require a different procedure in *ComputeCC*.

## Acknowledgments

This work was supported in parts by the NSF Grant OCI-0904809, and the Defense Threat Reduction Agency Grant HDTRA1-14-C-0007.

## References

- [1] V. Agarwal, F. Petrini, D. Pasetto, D.A. Bader, Scalable graph exploration on multicore processors, in: SuperComputing (SC), 2010, pp. 1–11.
- [2] M. Belgin, G. Back, C.J. Ribbens, Pattern-based sparse matrix representation for memory-efficient SMVM kernels, in: Proceedings of the 23rd International ACM Conference on International Conference on Supercomputing, ICS '09, 2009, pp. 100–109.
- [3] M.D. Beynon, T. Kurç, Ü.V. Çatalyürek, C. Chang, A. Sussman, J. Saltz, Distributed processing of very large datasets with DataCutter, *Parallel Comput.* 27 (11) (2001) 1457–1478.
- [4] U. Brandes, A faster algorithm for betweenness centrality, *J. Math. Sociol.* 25 (2) (2001) 163–177.
- [5] A. Buluç, J.R. Gilbert, The combinatorial BLAS: design, implementation, and applications, *Int. J. High Perform. Comput. Appl. (IJHPCA)* (2011).
- [6] A. Buluç, S. Williams, L. Oliker, J. Demmel, Reduced-bandwidth multithreaded algorithms for sparse matrix-vector multiplication, in: In Proc IPDPS, 2011.
- [7] S.Y. Chan, I.X.Y. Leung, P. Liò, Fast centrality approximation in modular networks, in: Proc. of the 1st ACM International Workshop on Complex Networks Meet Information and Knowledge Management (CNIKM), 2009.
- [8] Ö. Şimşek, A.G. Barto, Skill characterization based on betweenness, in: Proc. of Advances in Neural Information Processing Systems (NIPS), 2008.
- [9] D. Eppstein, J. Wang, Fast approximation of centrality, in: Proceedings of the Twelfth Annual ACM-SIAM Symposium on Discrete Algorithms (SODA), 2001.
- [10] O. Green, R. McColl, D.A. Bader, A fast algorithm for streaming betweenness centrality, in: Proc. of SocialCom, 2012.
- [11] T.D.R. Hartley, Ü. V. Çatalyürek, A. Ruiz, F. Igual, R. Mayo, M. Ujaldon, Biomedical image analysis on a cooperative cluster of GPUs and multicores, in: Proc. of the 22nd Annual International Conference on Supercomputing, ICS 2008, 2008, pp. 15–25.
- [12] T.D.R. Hartley, A.R. Fasih, C.A. Berdanier, F. Özgüner, Ü. V. Çatalyürek, Investigating the use of GPU-accelerated nodes for SAR image formation, in: Proc. of the IEEE International Conference on Cluster Computing, Workshop on Parallel Programming on Accelerator Clusters (PPAC), 2009.
- [13] T.D.R. Hartley, E. Saule, Ü.V. Çatalyürek, Improving performance of adaptive component-based dataflow middleware, *Parallel Comput.* 38 (6–7) (2012) 289–309.
- [14] J. Hopcroft, R. Tarjan, Algorithm 447: efficient algorithms for graph manipulation, *Commun. ACM* 16 (6) (1973) 372–378.
- [15] Y. Jia, V. Lu, J. Hoberock, M. Garland, J.C. Hart, Edge vs. node parallelism for graph centrality metrics, in: GPU Computing Gems: Jade Edition, Morgan Kaufmann, 2011.
- [16] S. Jin, Z. Huang, Y. Chen, D.G. Chavarría-Miranda, J. Feo, P.C. Wong, A novel application of parallel betweenness centrality to power grid contingency analysis, in: Proc. of IPDPS, 2010.
- [17] S. Kintali, Betweenness centrality: algorithms and lower bounds, CoRR, abs/0809.1906, 2008.
- [18] V. Krebs, Mapping networks of terrorist cells, *Connections* 24 (2002).
- [19] M.-J. Lee, J. Lee, J.Y. Park, R.H. Choi, C.-W. Chung, QUBE: a Quick algorithm for Updating BEtweenness centrality, in: Proc. of World Wide Web Conference (WWW), 2012.
- [20] R. Lichtenwalter, N.V. Chawla, Disnet: a framework for distributed graph computation, in: Proc. of ASONAM, 2011.
- [21] X. Liu, M. Smelyanskiy, E. Chow, P. Dubey, Efficient sparse matrix-vector multiplication on x86-based many-core processors, in: Proceedings of the 27th International ACM Conference on International Conference on Supercomputing, ICS '13, 2013.
- [22] K. Madduri, D. Ediger, K. Jiang, D.A. Bader, D.G. Chavarría-Miranda, A faster parallel algorithm and efficient multithreaded implementations for evaluating betweenness centrality on massive datasets. In 23rd International Parallel and Distributed Processing Symposium Workshops, Workshop on Multithreaded Architectures and Applications (MTAAP), 2009.
- [23] E.L. Merrer, G. Trédan, Centralities: capturing the fuzzy notion of importance in social graphs, in: Proc. of the Second ACM EuroSys Workshop on Social Network Systems (SNS), 2009.
- [24] K. Okamoto, W. Chen, X.-Y. Li, Ranking of closeness centrality for large-scale social networks, in: Proc. of Frontiers in Algorithmics Second Annual International Workshop (FAW), 2008.
- [25] P. Pande, D.A. Bader, Computing betweenness centrality for small world networks on a GPU, In 15th Annual High Performance Embedded Computing Workshop (HPEC), 2011.
- [26] S. Porta, V. Latora, F. Wang, E. Strano, A. Cardillo, S. Scellato, V. Iacoviello, R. Messori, Street centrality and densities of retail and services in Bologna, Italy, *Environ. Plan. B: Plan. Des.* 36 (3) (2009) 450–465.
- [27] A.E. Saryüce, K. Kaya, E. Saule, Ü. V. Çatalyürek, Incremental algorithms for closeness centrality, in: Proc of IEEE Int'l Conference on BigData, October 2013.
- [28] A.E. Saryüce, E. Saule, K. Kaya, Ü. V. Çatalyürek, Shattering and compressing networks for betweenness centrality, in: SIAM International Conference on Data Mining, (SDM), May 2013.
- [29] A.E. Saryüce, E. Saule, K. Kaya, Ü. V. Çatalyürek, STREAMER: a distributed framework for incremental closeness centrality computation, in: Proc. of IEEE Cluster 2013, September 2013.

- [30] A.E. Saryüce, E. Saule, K. Kaya, Ü. V. Çatalyürek, Hardware/software vectorization for closeness centrality on multi-/many-core architectures, in: 28th International Parallel and Distributed Processing Symposium Workshops, Workshop on Multithreaded Architectures and Applications (MTAAP), May 2014.
- [31] A.E. Saryüce, E. Saule, K. Kaya, Ü. V. Çatalyürek, Regularizing graph centrality computations, *J. Parallel Distrib. Comput.*, in press.
- [32] Z. Shi, B. Zhang, [Fast network centrality analysis using GPUs](#), *BMC Bioinf.* 12 (2011) 149.
- [33] R. Vuduc, J. Demmel, K. Yelick, OSKI: a library of automatically tuned sparse matrix kernels, in: Proc. SciDAC 2005, J. of Physics: Conference Series, 2005.