

Work-Efficient Parallel Union-Find

Natcha Simsiri¹, Kanat Tangwongsan^{2*}, Srikanta Tirthapura³, and Kun-Lung Wu⁴

¹ *College of Information and Computer Sciences, University of Massachusetts-Amherst*

² *Computer Science Program, Mahidol University International College*

³ *Department of Electrical and Computer Engineering, Iowa State University*

⁴ *IBM T.J. Watson Research Center*

SUMMARY

The incremental graph connectivity (IGC) problem is to maintain a data structure that can quickly answer whether two given vertices in a graph are connected, while allowing more edges to be added to the graph. IGC is a fundamental problem and can be solved efficiently in the sequential setting using a solution to the classical union-find problem. However, sequential solutions are not sufficient to handle modern-day large, rapidly-changing graphs where edge updates arrive at a very high rate.

We present the first shared-memory parallel data structure for union-find (equivalently, IGC) that is both provably work-efficient (i.e. performs no more work than the best sequential counterpart) and has polylogarithmic parallel depth. We also present a simpler algorithm with slightly worse theoretical properties, but which is easier to implement and has good practical performance. Our experiments on large graph streams with various degree distributions show that it has good practical performance, capable of processing hundreds of millions of edges per second using a 20-core machine.

Copyright © 2016 John Wiley & Sons, Ltd.

Received ...

1. INTRODUCTION

The classical *Union-Find* problem (see, e.g., [5]) seeks to maintain a collection of disjoint sets, supporting the following two operations:

- (1) **union**(u, v): given elements u and v , combine the sets containing u and v into one set and return (a handle to) the combined set; and
- (2) **find**(v): given an element v , return (a handle to) the set containing v . If u and v belong to the same set, it is guaranteed that $\text{find}(u) = \text{find}(v)$.

This problem has many applications, including incremental graph connectivity on undirected graphs. Whereas the basic graph connectivity question asks whether there is a path between two

*Correspondence to: 999 Phutthamonthon 4 Road, Salaya, Nakhonpathom, Thailand 73170. Email: kanat.tan@mahidol.edu

given vertices, the *incremental graph connectivity* (IGC) problem seeks an answer to the graph connectivity question as edges are added over time. IGC has a well-known solution that uses Union-Find. In this case, each set in the union-find data structure is one connected component. Therefore, answering whether u and v are connected amounts to checking whether $\text{find}(u) = \text{find}(v)$. Furthermore, adding an edge (w, x) to the graph amounts to invoking $\text{union}(w, x)$ on the union-find data structure.

With a sharp increase in the amount of linked data, IGC now often has to be solved at a much larger scale than before. Every minute, a staggering amount of linked data is being generated from social media interactions, the Internet of Things (IoT) devices, among others—and timely insights from them are much sought after. These data are usually cast as a stream of edges with the goal of maintaining certain local and global properties on the accumulated data. On the one hand, modern streaming systems (e.g., IBM Streams [10] and Spark Streaming [23]) provide a software platform for using parallelism to achieve high-throughput processing. On the other hand, scalable, parallel, and dynamic algorithms are still scarce but needed to effectively utilize such platforms. This paper makes a step forward in designing parallel algorithms for IGC, where the input is a graph stream consisting of a sequence of edges, and the queries are connectivity queries. Our main ingredient is a novel parallel algorithm for Union-Find.

A data stream that has a high rate of flow is still inherently sequential in the sense that each data item, often called a tuple, arrives one after another. To facilitate parallel processing, many streaming systems (e.g., Apache Spark Streaming [23]) adopt a “discretized stream” input model: a stream is divided into a sequence of *minibatches*, so that each minibatch is simply an array of data tuples and can be processed as such. Using this model, we design parallel algorithms to efficiently process a minibatch of edges and a minibatch of queries. More precisely, for a fixed vertex set V , a graph stream \mathcal{A} is a sequence of minibatches A_1, A_2, \dots , where each minibatch A_i is a set of edges on V . The graph at the end of observing A_t , denoted by G_t , is $G_t = (V, \cup_{i=1}^t A_i)$ containing all the edges up to t . The minibatches A_i s may have different sizes. Equivalently, each edge (u, v) can be viewed as a single $\text{union}(u, v)$ operation that merges the sets containing u and v , and a minibatch A_i is a set of union ’s that have to be applied in parallel.

With this setup, the data structure that we will design has to support two operations. We use union-find and connectivity terminologies interchangeably:

- (1) The **Bulk-Union** operation takes as input a minibatch of edges A_i and adds them to the graph—this involves multiple union operations, one per edge.
- (2) The **Bulk-Same-Set** operation takes a minibatch of vertex-pairs $\{(u_i, v_i)\}_{i=1}^k$ and returns for each query, whether the two vertices are connected on the edges observed so far in the stream—this conceptually involves multiple find operations, two per vertex-pair.

On this data structure, the **Bulk-Union** and **Bulk-Same-Set** operations are each invoked with a (potentially large) minibatch of input, each processed using a parallel computation. But a bulk operation, say a **Bulk-Union**, must complete before the next operation, say a **Bulk-Same-Set**, can begin. By not allowing unions and finds in the same batch, the model provides clean semantics for what constitutes the graph being queried.

Contributions

The main contribution of this paper is the first shared-memory parallel algorithm for union-find, and hence, for IGC, that is provably work-efficient and has polylogarithmic parallel depth. The detailed contributions are as follows:

- *Simple Parallel Algorithm.* We first present a simple parallel algorithm for Union-Find (Section 4). This is easy to implement and has good theoretical properties. On a graph with n vertices, a Bulk-Union operation with b edges, and a Bulk-Same-Set involving b queries each require $O(b \log n)$ work and $O(\text{polylog}(n))$ parallel depth. The data structure consumes $O(n)$ memory.
- *Work-Efficient Parallel Algorithm.* We then present an improved parallel algorithm for Union-Find (Section 5), with total work $O((m + q)\alpha(m + q, n))$, where m is the total number of union operations across all minibatches and q is the total number of find operations across all minibatches, and α is an inverse Ackermann's function (see Section 2). Equivalently, this is a parallel algorithm for IGC with total work $O((m + q)\alpha(m + q, n))$, where m is the total number of edges across all minibatches, q is the total number of connectivity queries across all minibatches. This bound on work matches the work of the theoretically time-optimal sequential counterpart i.e., it is *work-efficient*. Further, processing a minibatch takes $O(\text{polylog}(m, n))$ parallel depth. Hence, the sequential bottleneck in the runtime of the parallel algorithm is very small, and the algorithm is capable of using almost a linear number of processors efficiently. We are not aware of a prior parallel algorithm for Union-Find with such provable properties on work and depth.
- *Implementation and Evaluation.* We implemented and benchmarked a variation of our simple parallel algorithm on a shared-memory machine. Our experimental results show that the algorithm achieves good speedups in practice and is able to efficiently use available parallelism. On a 20-core machine, it can process hundreds of millions of edges per second, and realize a speedup of 8–11x over its single threaded performance. Further analysis shows good scalability properties as the number of threads is varied. We describe this in Section 6.

We assume the concurrent-read, concurrent-write (CRCW) model of shared memory. Our algorithm designed for the CRCW model can work in other shared memory models such as exclusive-read, exclusive-write (EREW) PRAM, however with a depth, as well as work, that is a logarithmic factor worse. Hence, it will not be work-efficient in the EREW model. One of the bottlenecks in making it work-efficient in the EREW model is that our algorithm uses parallel integer sort, for which there is a work-efficient algorithm in the CRCW model, but not in the EREW model.

We note that this paper extended an earlier paper that appeared in Euro-Par 2016 [21].

2. RELATED WORK

Let n denote the number of vertices, m the number of operations, and α an inverse Ackermann's function (very slow-growing, practically a constant independent of n). In the sequential setting, the basic data structure for incremental connectivity is the well-studied union-find data structure [5]. Tarjan [22] presented an analysis that showed an amortized time of $O(\alpha(m, n))$ per find, which

has been shown to be optimal (see Seidel and Sharir [18] for an alternate analysis). Such solutions, however, often are unable to exploit parallelism.

Recent work on streaming graph algorithms mostly focuses on minimizing the memory requirement, with little attention paid to the use of parallelism. This type of work has largely focused on the “semi-streaming model” [8], which allows $O(n \cdot \text{polylog}(n))$ space usage. In this case, the union-find data structure [22] solves incremental connectivity in $O(n)$ space and a total time nearly linear in m . McGregor’s recent survey paper [15] discusses variants of this problem.

Under the condition that only $o(n)$ of workspace is allowed (sublinear memory), there are interesting tradeoffs for multi-pass algorithms. With $O(s)$ workspace, an algorithm can compute the connected components of a graph in $\Omega(n/s)$ passes [8]. Demetrescu et al. [7] consider the W-stream model, which allows the processing of streams in multiple passes in a pipelined manner: the output of the i -th pass is given as input to the $(i + 1)$ -th pass. They demonstrate a tradeoff between the number of passes and the memory needed. With $O(s)$ space, their algorithm computes connected components in $O((n \log n)/s)$ passes. Demetrescu et al. [6] present a simulation of a PRAM algorithm on the W-Stream model, allowing existing PRAM algorithms to run *sequentially* in the W-Stream model.

McColl et al. [14] present a parallel algorithm for maintaining connected components in a fully dynamic graph. As part of a bigger project (STINGER), their work focuses on engineering algorithms that work well on real-world graphs and gives no theoretical analysis of the parallel complexity. Our work gives a provable bound on the cost of the parallel computation. In particular, we achieve work equal to the best sequential algorithm while keeping depth polylogarithmic. However, compared with our work, their algorithm addresses a more general setting that allows for arbitrary edge deletions and additions. Manne and Patwary [13] present a parallel Union-Find algorithm. However, their algorithm is for distributed memory computers while ours is for shared-memory machines. Patwary et al. [16] present a shared-memory parallel algorithm for computing a spanning forest that uses the Union-Find data structure, but this does not provide theoretical guarantees of work-efficiency.

Berry et al. [2] present methods for maintaining connected components in a parallel graph stream model, called X-Stream, which periodically ages out edges. Their algorithm is essentially an “unrolling” of the algorithm of [7], and edges are passed from one processor to another until the connected components are found by the last processor in the sequence. When compared with our work, the input model and notions of correctness differ. Our work views the input stream as a sequence of batches, each a set of edges or a set of queries, which are unordered within the set. Their algorithm strictly respects the sequential ordering of edges and queries. Further, they age out edges (we do not). Also, they do not give provable parallel complexity bounds.

There are multiple parallel (batch) algorithms for graph connectivity including [20, 9] that are work-efficient (linear in the number of edges) and that have polylogarithmic depth. However, these algorithms assume that edges are all known in advance. The work of these parallel algorithms is of the same order as that of the best sequential algorithm for connected components of a static graph. However, these algorithms are not efficient for a dynamic graph. If the graph changes due to the addition of a few edges, then the entire algorithm will have to be re-run to update the connected components in the graph. Hence, they would not be suitable for parallel IGC.

Shiloach and Vishkin (as presented by Jájá [11]) describe a batch parallel algorithm for graph connectivity based on Union-Find that runs in $O(m \log n)$ work and polylogarithmic depth. Like ours, their algorithm also relies on a union-find-like forest and linking together trees (“grafting” in their terminology); however, theirs relies on pointer jumping to keep the tree shallow to the point of keeping stars. As far as we know, there is no easy way to port their algorithm to our setting, let alone perform parallel path compression. Our algorithm is simpler as it sidesteps concurrent grafting.

Prior work on wait-free implementations of the union-find data structure [1] focuses on the asynchronous model, where the goal is to be correct under all possible interleavings of operations; unlike us, they do not focus on bulk processing of edges and the parallel depth of the computation. There is also a long line of work on sequential algorithms for maintaining graph connectivity on an evolving graph. See the recent work by [12] that addresses this problem in the general dynamic case and the references therein.

3. PRELIMINARIES AND NOTATION

Throughout the paper, let $[n]$ denote the set $\{0, 1, \dots, n\}$. A sequence is written as $X = \langle x_1, x_2, \dots, x_{|X|} \rangle$, where $|X|$ denotes the length of the sequence. For a sequence X , the i -th element is denoted by X_i or $X[i]$. Following the set-builder notation, we denote by $\langle f(x) : \Phi(x) \rangle$ a sequence generated (logically) by taking all elements that satisfy $\Phi(x)$, preserving their original ordering, and transform them by applying f . For example, if T is a sequence of numbers, the notation $\langle 1 + f(x) : x \in T \text{ and } x \text{ odd} \rangle$ means a sequence created by taking each element x from T that are odd and map x to $1 + f(x)$, retaining their original ordering. Furthermore, we write $S \oplus T$ to mean the concatenation of S and T .

We design algorithms in the work-depth model assuming an underlying CRCW PRAM machine model. The *work* of a parallel algorithm is the total operation count across all processors, and the *depth* (also called parallel time or span) is the length of the longest chain of dependencies within a parallel computation. The gold standard for a parallel algorithm in this model is to perform the same amount of work as the best sequential counterpart (work-efficient) and to have polylogarithmic depth.

We use standard parallel operations such as filter, prefix sum, map (applying a constant-cost function), and pack, all of which can be implemented with $O(n)$ work and $O(\log^2(n))$ depth on an input sequence of length n . Given a sequence of n numbers, there is a duplicate removal algorithm `removeDup` that runs in $O(n)$ work and $O(\log^2 n)$ depth [11]. We also use the following result on sorting integer keys in a small range faster than a typical comparison-based algorithm:

Theorem 1 (Parallel Integer Sort [17])

There is an algorithm, `intSort`, that takes a sequence of integers a_1, a_2, \dots, a_n , where $a_i \in [0, c \cdot n]$, $c = O(1)$, and produces a sorted sequence in $O(n)$ work and $\text{polylog}(n)$ depth.

Parallel Connectivity: For a graph $G = (V, E)$, a connected component algorithm (CC) computes a sequence of connected components $\langle C_i \rangle_{i=1}^k$, where each C_i is a list of vertices in the component. There are algorithms for CC that have $O(|V| + |E|)$ work and $O(\text{polylog}(|V|, |E|))$ depth (e.g., [9, 20]), with Gazit’s algorithm [9] requiring $O(\log |V|)$ depth.

4. SIMPLE BULK-PARALLEL DATA STRUCTURE

This section describes a simple bulk-parallel data structure for Union-Find and IGC. The data structure is conceptually simple but will be instructive for the theoretical improvements presented in the section that follows. As before, n is the number of vertices in the graph stream, or equivalently the total number of elements across all disjoint sets. The main result for this section is as follows:

Theorem 2

There is a bulk-parallel data structure for Union-Find and IGC, given by Algorithms **Simple-Bulk-Union** and **Simple-Bulk-Same-Set**, where

- (1) The total memory consumption is $O(n)$ words.
- (2) A minibatch of b edges is processed by **Simple-Bulk-Union** in $O(b \log n)$ work and $O(\log \max(b, n))$ parallel depth.
- (3) A minibatch of q connectivity queries is answered by **Simple-Bulk-Same-Set** in $O(\log n)$ parallel depth and $O(q \log n)$ total work.

In a nutshell, we show how to bootstrap a standard union-find structure to take advantage of parallelism while preserving the height of the union-find forest to be at most $O(\log n)$. For concreteness, we will work with union by size, though other variants (e.g., union by rank) will also work. The crux here is to handle concurrent union operations efficiently.

Sequential Union-Find Implementation: We review a basic union-find implementation that uses union by size. Conceptually, the union-find data structure maintains a union-find forest with one tree for each set in the partition. In this view, $\text{find}(u)$ returns the vertex that is the root of the tree containing u and $\text{union}(u, v)$ joins together the roots of the tree containing u and the tree containing v , by pointing the root of one tree to the root of the other. The trees in a union-find forest are typically represented by remembering each node's parent, in an array parent of length n , where $\text{parent}[u]$ is the tree's parent of u or $\text{parent}[u] = u$ if it is the root of its component.

The running time of the union and find operations depends on the maximum height of a tree in the union-find forest. To keep the height small, at most $O(\log n)$, a simple strategy, known as *union by size*, is for union to always link the tree with fewer vertices into the tree with more vertices. The data structure also keeps an array for the sizes of the trees. The following results are standard (see [18], for example):

Lemma 3 (Sequential Union-Find)

On a graph with vertices $[n]$, a sequential union-find data structure implementing the union-by-size strategy consumes $O(n)$ space and has the following characteristics:

- Every union-find tree has height $O(\log n)$ and each find takes $O(\log n)$ sequential time.
- Given two distinct roots u and v , the operation $\text{union}(u, v)$ implementing union by size takes x sequential time.

Our data structure maintains an instance of this union-find data structure, called U . Notice that the find operation is read-only. Unlike the more sophisticated variants, this version of union-find does not perform path compression.

4.1. Answering Connectivity Queries in Parallel

Connectivity queries can be easily answered in parallel, using read-only finds on U . To answer whether u and v are connected, we compute $U.\text{find}(v)$ and $U.\text{find}(u)$, and report if the results are equal. To answer multiple queries in parallel, we note that because the finds are read-only, we can answer all queries simultaneously independently of each other. We present Simple-Bulk-Same-Set in Algorithm 1.

Algorithm 1: Simple-Bulk-Same-Set($U, \langle (u_i, v_i) \rangle_{i=1}^q$).

Input: U is the union find structure, and (u_i, v_i) is a pair of vertices, for $i = 1, \dots, q$.

Output: For each i , whether or not u_i and v_i are in the same set (i.e., connected in the graph).

```

1: for  $i = 1, 2, \dots, q$  do in parallel
2:    $a_i \leftarrow (U.\text{find}(u_i) == U.\text{find}(v_i))$ 
3: return  $\langle a_1, a_2, \dots, a_q \rangle$ 

```

Correctness follows directly from the correctness of the base union-find structure. The parallel complexity is simply that of applying q operations of $U.\text{find}$ in parallel:

Lemma 4

The parallel depth of Simple-Bulk-Same-Set is $O(\log n)$, and the work is $O(q \log n)$, where q is the number of queries input to the algorithm.

4.2. Adding a Minibatch of Edges

How can one incorporate (in parallel) a minibatch of edges A into an existing union-find structure? Sequentially, this is simple: invoke **union** on the endpoints of every edge of A . But it is dangerous to blindly apply the **union** operations in parallel since **union** updates the forest, potentially leading to inconsistencies in the structure.

However, it is safe to run multiple unions in parallel as long as they operate on different trees. Because there may be many union operations involving the same tree, this idea alone is not sufficient—running these sequentially will result in a large parallel depth. For instance, consider adding the edges of a star graph (with a very high degree) to an empty graph. Because all the edges share a common endpoint (the center), this vertex is involved in every **union**, and hence no two operations can proceed in parallel.

To tackle this problem, our algorithm transforms the minibatch of edges A into a structure that can be connected up easily in parallel. For illustration, we revisit the example when the minibatch is itself a star graph. Suppose there are seven edges within the minibatch: $(v_1, v_2), (v_1, v_3), (v_1, v_4), \dots, (v_1, v_8)$. By examining the minibatch, we find that all of v_1, \dots, v_8 will belong to the same component. We now apply these connections to the graph.

In terms of connectivity, it does not matter whether we apply the actual edges that arrived, or a different, but equivalent set of edges; it only matters that the relevant vertices are connected up. To connect up these vertices, our algorithm schedules the unions in only three parallel rounds as follows. The notation $X \parallel Y$ indicates that X and Y are run in parallel:

```

1: union( $v_1, v_2$ )  $\parallel$  union( $v_3, v_4$ )  $\parallel$  union( $v_5, v_6$ )  $\parallel$  union( $v_7, v_8$ )
2: union( $v_1, v_3$ )  $\parallel$  union( $v_5, v_7$ )
3: union( $v_1, v_5$ )

```

As we will soon see, such a schedule can be constructed for a component of any size, provided that no two vertices in the component are connected previously. The resulting parallel depth is logarithmic in the size of the minibatch.

Algorithm 2: Simple-Bulk-Union(U, A)

Input: U : the union find structure, A : a set of edges to add to the graph.
 \triangleright Relabel each (u, v) with the roots of u and v
1: $A' \leftarrow \langle (p_u, p_v) : (u, v) \in A \text{ where } p_u = U.\text{find}(u) \text{ and } p_v = U.\text{find}(v) \rangle$
 \triangleright Remove self-loops
2: $A'' \leftarrow \langle (u, v) : (u, v) \in A' \text{ where } u \neq v \rangle$
3: $\mathcal{C} \leftarrow \text{CC}(A'')$
4: **foreach** $C \in \mathcal{C}$ **do** in parallel
5: | Parallel-Join(U, C)

To add a minibatch of edges, our Simple-Bulk-Union algorithm, presented in Algorithm 2, proceeds in three steps:

\triangleright **Step 1:** Relabel edges as links between existing components. An edge $\{u, v\} \in A$ does not simply join vertices u and v . Due to potential existing connections in G , it joins together C_u and C_v , the component containing u and the component containing v , respectively. In our representation, the identifier of the component containing u is $U.\text{find}(u)$, so $C_u = U.\text{find}(u)$ and similarly $C_v = U.\text{find}(v)$. Lines 1–2 in Algorithm 2 create A'' by relabeling each endpoint of an edge with the identifier of its component, and dropping edges that are within the same component.

\triangleright **Step 2:** Discover new connections arising from A . After the relabeling step, we are implicitly working with the graph $\tilde{H} = (V_{\tilde{H}}, A'')$, where $V_{\tilde{H}}$ is the set of all connected components of G that pertain to A (i.e., all the roots in the union-find forest reachable from vertices incident on A) and A'' is the connections between them. In other words, \tilde{H} is a graph on “supernodes” and the connections between them using the edges of A . In this view, a connected component on \tilde{H} represents a group of existing components of G that have just become connected as a result of incorporating A . While never materializing the vertex set $V_{\tilde{H}}$, Line 3 in Algorithm 2 computes \mathcal{C} , the set of connected components of \tilde{H} , using a linear-work parallel algorithm for connected components, CC (see Section 3).

\triangleright **Step 3:** Commit new connections to U . With the preparation done so far, the final step only has to make sure that the pieces of each connected component in \mathcal{C} are linked together in U . Lines 4–5 of Algorithm 2 go over the components of \mathcal{C} in parallel, seeking help from Parallel-Join, the real workhorse that links together the pieces.

Connecting a Set of Components within U : Let $v_1, v_2, \dots, v_k \in [n]$ be distinct tree roots from the union-find forest U that form a component in \mathcal{C} , and need to be connected together. Algorithm Parallel-Join connects them up in $O(\log k)$ iterations using a divide-and-conquer approach. Given a sequence of tree roots, the algorithm splits the sequence in half and recursively connects the roots in the first half, in parallel with connecting the roots in the second half. Since components in the first half and the second half have no common vertices, handling them in parallel will not cause a conflict. Once both calls return with their respective new roots, they are unioned together.

Correctness of Parallel-Join is immediate since the order that the union calls are made does not matter, and we know that different union calls that proceed in parallel always work on separate sets of tree roots, posing no conflicts. This means that the work and depth of Parallel-Join follow

Algorithm 3: Parallel-Join(U, C)

Input: U : the union-find structure, C : a seq. of tree roots
Output: The root of the tree after all of C are connected

```

1: if  $|C| == 1$  then
2:   return  $C[1]$ 
3: else
4:    $\ell \leftarrow \lfloor |C|/2 \rfloor$ 
5:    $u \leftarrow \text{Parallel-Join}(U, C[1, 2, \dots, \ell])$  in parallel with
      $v \leftarrow \text{Parallel-Join}(U, C[\ell + 1, \ell + 2, \dots, |C|])$ 
6:   return  $U.\text{union}(u, v)$ 

```

the recurrences

$$W(k) = 2W(k/2) + O(1)$$

$$D(k) = D(k/2) + O(1),$$

which solve to $W(k) = O(k)$ and $D(k) = O(\log k)$. Hence:

Lemma 5

Given k distinct roots of U , Algorithm Parallel-Join runs in $O(k)$ work and $O(\log k)$ depth.

Lemma 6 (Correctness of Simple-Bulk-Union)

Let U be the shared-memory union-find data structure after calling Simple-Bulk-Union on every minibatch of edges whose union equals a graph G , then for any $u, v \in V$, $U.\text{find}(u) = U.\text{find}(v)$ if and only if u and v are connected in G .

Proof

We will prove this lemma inductively on the minibatches. The lemma is vacuously true initially when there are no edges. Now consider a minibatch of edges A . Let G_1 be the set of edges that arrived prior to A and U_1 the state of the union-find structure corresponding to G_1 (i.e., prior to A). Further, let $G_2 = G_1 \cup A$ and let U_2 be the state of the union-find structure after Simple-Bulk-Union(U_1, A). We will assume inductively that U_1 is correct with respect to G_1 and show that U_2 is correct with respect to G_2 .

Let $x \neq y$ be a pair of vertices in V . We consider the following two cases:

Case I: x and y are not connected in G_2 . We know that x and y are not connected in G_1 , either. Therefore, if $r_x = U_1.\text{find}(x)$ and $r_y = U_1.\text{find}(y)$, then $r_x \neq r_y$ by our inductive assumption. We also know that r_x and r_y will belong to different components of \mathcal{C} , the connected components generated in Simple-Bulk-Union. Hence, they will continue to be in separate components after Parallel-Join— $U_2.\text{find}(x) = U_2.\text{find}(r_x) \neq U_2.\text{find}(r_y) = U_2.\text{find}(y)$.

Case II: x and y are connected in G_2 . We know that there must be a path $x = v_1, v_2, \dots, v_t = y$ in G_2 . We will show that $U_2.\text{find}(v_1) = U_2.\text{find}(v_2) = \dots = U_2.\text{find}(v_t)$, leading to the conclusion that $U_2.\text{find}(x) = U_2.\text{find}(y)$. To this end, consider any pair v_i and v_{i+1} , $1 \leq i \leq (t-1)$. Let $r_i = U_1.\text{find}(v_i)$ and $r_{i+1} = U_1.\text{find}(v_{i+1})$ denote the roots of the trees that contain v_i and v_{i+1} , respectively, in U_1 .

If $r_i = r_{i+1}$, then it will remain true that $U_2.\text{find}(v_i) = U_2.\text{find}(r_i) = U_2.\text{find}(r_{i+1}) = U_2.\text{find}(v_{i+1})$. If, however, $r_i \neq r_{i+1}$, then v_i and v_{i+1} are not connected in G_1 , so $\{v_i, v_{i+1}\} \in A$. Thus, in Steps 1 and 2 of Simple-Bulk-Union, the edge $\{r_i, r_{i+1}\}$ is present in A'' (note

this edge is not a self-loop and is not eliminated in Step 2). Therefore, in Step 3, when the connected components of A'' are computed, r_i and r_{i+1} are in the same component of \mathcal{C} . Consequently, **Parallel-Join** will eventually join the trees of r_i and r_{i+1} into the same component in U_2 . Hence, $U_2.\text{find}(r_i) = U_2.\text{find}(r_{i+1})$. Since $U_2.\text{find}(v_i) = U_2.\text{find}(r_i)$ and $U_2.\text{find}(v_{i+1}) = U_2.\text{find}(r_{i+1})$, we have $U_2.\text{find}(v_i) = U_2.\text{find}(v_{i+1})$. Proceeding thus, we have $U_2.\text{find}(x) = U_2.\text{find}(y)$ in Case II. □

Lemma 7 (Complexity of Simple-Bulk-Union)

On input a minibatch A with b edges/union's, **Simple-Bulk-Union** takes $O(b \log n)$ work and $O(\log \max(b, n))$ depth.

Proof

There are three parts to the work and depth of **Simple-Bulk-Union**. First is the generation of A' and A'' . For each $(u, v) \in A$, we invoke $U.\text{find}$ on u and v , requiring $O(\log n)$ work and depth per edge. Since the edges are processed in parallel, this leads to $O(b \log n)$ work and $O(\log n)$ depth. Then, A'' is derived from A' through a parallel filtering algorithm, using $O(|A'|) = O(b)$ work and $O(\log b)$ depth. The second part is the computation of connected components of A'' which can be done in $O(|A''|) = O(b)$ work and $O(\log \max(b, n))$ depth using the algorithm of Gazit [9]. The third part is **Parallel-Join**. As the number of components cannot exceed b , and using Lemma 5, we have that the total work in **Parallel-Join** is $O(b)$ and depth is $O(\log n)$. Adding the three parts, we arrive at the lemma. □

5. WORK-EFFICIENT PARALLEL ALGORITHM

Whereas the best sequential data structures (e.g., [22]) require $O((m + q)\alpha(m + q, n))$ work to process m edges and q queries, our basic data structure from the previous section needs up to $O((m + q) \log n)$ work for the same input stream. This section describes improvements that make it match the best sequential work bound while preserving the polylogarithmic depth guarantee. The main result for this section is as follows:

Theorem 8

There is a bulk-parallel data structure for Union-Find and IGC with the following properties:

- (1) The total memory consumption is $O(n)$ words.
- (2) The depth of **Bulk-Union** and **Bulk-Same-Set** is $O(\text{polylog}(n))$ each.
- (3) Over the lifetime of the data structure, the total work for processing m edge updates (across all **Bulk-Union**) and q queries is $O((m + q)\alpha(m + q, n))$.

Overview: All sequential data structures with a $O((m + q)\alpha(n))$ bound use a technique called path compression, which shortens the path that **find** traverses on to reach the root, making subsequent operations cheaper. Our goal in this section is to enable path compression during parallel execution. We present a new parallel **find** procedure called **Bulk-Find**, which answers a set of **find** queries in parallel and performs path compression.

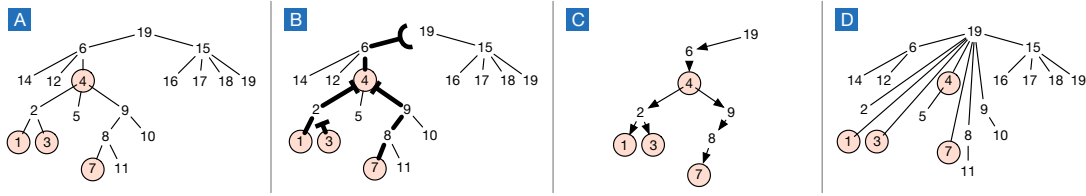


Figure 1. A: An example union-find tree with sample queries circled; B: Bolded edges are paths, together with their stopping points, that result from the traversal in Phase I; C: The traversal graph R_U recorded as a result of Phase I; and D: The union-find tree after Phase II, which updates all traversed nodes to point to their roots.

Algorithm 4: Bulk-Find(U, S)—find the root in U for each $s \in S$ with path compression.

Input: U is the union find structure. For $i = 1, \dots, |S|$, $S[i]$ is a vertex in the graph

Output: A response array res of length $|S|$ where $res[i]$ is the root of the tree of the vertex $S[i]$ in the input.

▷ **Phase I:** Find the roots for all queries

1: $R_0 \leftarrow \langle (S[k], \text{null}) : k = 0, 1, 2, \dots, |S| - 1 \rangle$

2: $F_0 \leftarrow \text{mkFrontier}(R_0, \emptyset)$, $roots \leftarrow \emptyset$, $visited \leftarrow \emptyset$, $i \leftarrow 0$

3: **while** $R_i \neq \emptyset$ **do**

4: $visited \leftarrow visited \cup F_i$

5: $R_{i+1} \leftarrow \langle (\text{parent}[v], v) : v \in F_i \text{ and } \text{parent}[v] \neq v \rangle$

6: $roots \leftarrow roots \cup \{v : v \in F_i \text{ where } \text{parent}[v] = v\}$

7: $F_{i+1} \leftarrow \text{mkFrontier}(R_{i+1}, visited)$, $i \leftarrow i + 1$

▷ Set up response distribution

8: Create an instance of RD with $R_U = R_0 \oplus R_1 \oplus \dots \oplus R_i$

▷ **Phase II:** Distribute the answers and shorten the paths

9: $D_0 \leftarrow \{(r, r) : r \in roots\}$, $i \leftarrow 0$

10: **while** $D_i \neq \emptyset$ **do**

11: For each $(v, r) \in D_i$, in parallel, $\text{parent}[v] \leftarrow r$

12: $D_{i+1} \leftarrow \bigcup_{(v,r) \in D_i} \{(u, r) : u \in RD.\text{allFrom}(v) \text{ and } u \neq \text{null}\}$. That is, create D_{i+1} by expanding every $(v, r) \in D_i$ as the entries of $RD.\text{allFrom}(v)$ excluding **null**, each inheriting r .

13: $i \leftarrow i + 1$

14: For $i = 0, 1, 2, \dots, |S| - 1$, in parallel, make $res[i] \leftarrow \text{parent}[S[i]]$

15: **return** res

def mkFrontier($R, visited$):

// nodes to go to next

1: $req \leftarrow \langle v : (v, -) \in R \wedge \text{not } visited[v] \rangle$

2: **return** removeDup(req)

To understand the benefits of path compression, consider a concrete example in Figure 1A, which shows a union-find tree T that is typical in a union-find forest. The root of T is $r = 19$. Suppose we need to support find 's from $u = 1$ and $v = 7$. When all is done, both $\text{find}(u)$ and $\text{find}(v)$ should return r . Notice that in this example, the paths to the root $u \rightsquigarrow r$ and $v \rightsquigarrow r$ meet at common vertex $w = 4$. That is, the two paths are identical from w onward to r . If find 's were done sequentially, say $\text{find}(u)$ before $\text{find}(v)$, then $\text{find}(u)$ —with path compression—would update all nodes on the $u \rightsquigarrow r$ path to point to r . This means that when $\text{find}(v)$ traverses the tree, the path to the root is significantly shorter: at this point, the next hop after w is already r .

The kind of sharing and shortcutting illustrated, however, is not possible when the find operations are run independently in parallel. Each find , unaware of the others, will proceed all the way to the root, missing out on possible sharing.

We fix this problem by organizing the parallel computation so that the work on different “flows” of `finds` is carefully coordinated. Algorithm 4 shows an algorithm `Bulk-Find`, which works in two phases, separating actions that only read from the tree from actions that only write to it:

▷ *Phase I*: Find the roots for all queries, coalescing flows as soon as they meet up. This phase should be thought of as running breadth-first search (BFS), starting from all the query nodes S at once. As with normal BFS, if multiple flows meet up, only one will move on. Also, if a flow encounters a node that has been traversed before, that flow no longer needs to go on. We also need to record the paths traversed so we can distribute responses to the requesting nodes.

▷ *Phase II*: Distribute the answers and shorten the paths. Using the transcript from Phase I, Phase II makes sure that all nodes traversed will point to the corresponding root—and answers delivered to all the `finds`. This phase, too, should be thought of as running breadth-first search (BFS) backwards from all the roots reached in Phase I. This BFS reverses the steps taken in Phase I using the trails recorded. There is a technical challenge in implementing this. Back in Phase I, to minimize the cost of recording these trails, the trails are kept as a list of directed edges (marked by their two endpoints) traversed. However, for the reverse traversal in Phase II to be efficient, it needs a means to quickly look up all the neighbors of a vertex (i.e., at every node, we must be able to find every flow that arrived at this node back in Phase I). For this, we design a data structure that takes advantage of hashing and integer sorting (Theorem 1) to keep the parallel complexity low. We discuss our solution to this problem in the section that follows (Lemma 9).

Example: We illustrate how the `Bulk-Find` algorithm works using the union-find forest from Figure 1A. The queries to the `Bulk-Find` are nodes that are circled. The paths traversed in Phase I are shown in panel B. If a flow is terminated, the last edge traversed on that flow is rendered as $\text{---}\perp$.

Notice that as soon as flows meet up, only one of them will carry on. In general, if multiple flows meet up at a point, only one will go on. Notice also that both the flow $1 \rightarrow 2 \rightarrow 4$ and the flow $7 \rightarrow 8 \rightarrow 9 \rightarrow 4$ are stopped at 4 because 4 is a source itself, which was started at the same time as 1 and 7. At the finish of Phase I, the graph (in fact a tree) given by R_\cup is shown in panel C. Finally, in Phase II, this graph is traversed and all nodes visited are updated to point to their corresponding root (as shown in panel D).

5.1. Response Distributor

Consider a sequence $R_\cup = \langle (from_i, to_i) \rangle_{i=1}^\lambda$. We need a data structure RD such that after some preprocessing of R_\cup , can efficiently answer the query $RD.allFrom(f)$ which returns a sequence containing all to_i where $from_i = f$.

To meet the overall running time bound, the preprocessing step cannot take more than $O(\lambda)$ work and $O(\text{polylog}(\lambda))$ depth. As far as we know, we cannot afford to generate, say, a sequence of sequences RD where $RD[f]$ is a sequence containing all to_i such that $from_i = f$. Instead, we propose a data structure with the following properties:

Lemma 9 (Response Distributor)

There is a data structure, response distributor (RD), that from input $R_\cup = \langle (from_i, to_i) \rangle_{i=1}^\lambda$ can be constructed in $O(\lambda)$ work and $O(\text{polylog}(n))$ depth. Each `allFrom` query can be answered in

$O(\log \lambda)$ depth. Furthermore, if \mathbb{F} is the set of unique $from_i$ (i.e., $\mathbb{F} = \{from_i : i = 1, \dots, \lambda\}$), then

$$\mathbf{E} \left[\sum_{f \in \mathbb{F}} \text{Work}(RD.\text{allFrom}(f)) \right] = O(\lambda).$$

Proof

Let h be a hash function from the domain of $from_i$'s (a subset of $[n]$) to $[\rho]$, where $\rho = 3\lambda$. To construct an RD , we proceed as follows. Compute the hash for each $from_i$ using $h(\cdot)$ and sort the ordered pairs $(from_i, to_i)$ by their hash values. Call this sorted array A . After sorting, we know that pairs with the same hash value are stored consecutively in A . Now create an array o of length $\rho + 1$ so that o_i marks the beginning of pairs whose hash value is i . If none of them hash to i , then $o_i = o_{i+1}$. These steps can be done using `intSort` and standard techniques in $O(\lambda)$ work and $O(\text{polylog}(\lambda))$ depth because the hash values range within $O(\lambda)$.

To support `allFrom(f)`, we compute $\kappa = h(f)$ and look in A between o_κ and $o_{\kappa+1} - 1$, selecting only pairs whose $from$ matches f . This requires at most $O(\log |o_{\kappa+1} - o_\kappa|) = O(\log \lambda)$ depth. The more involved question is, how much work is needed to support `allFrom` over all? To answer this, consider all the pairs in R_\cup with $from_i = f$. Let n_f denote the number of such pairs. These n_f pairs will be gone through by queries looking for f and other entries that happen to hash to the same value as f does. The exact number of times these pairs are gone through is $\beta_f := \#\{s \in \mathbb{F} : h(f) = h(s)\}$. Hence, across all queries $f \in \mathbb{F}$, the total work is $\sum_{f \in \mathbb{F}} n_f \beta_f$. But $\mathbf{E}[\beta_f] \leq 1 + \frac{|\mathbb{F}|}{\rho}$, so

$$\sum_{f \in \mathbb{F}} \mathbf{E}[n_f \beta_f] \leq \left(1 + \frac{|\mathbb{F}|}{\rho}\right) \sum_{f \in \mathbb{F}} n_f \leq \left(1 + \frac{\lambda}{3\lambda}\right) \lambda \leq 2\lambda$$

because $|\mathbb{F}| \leq \lambda$ and $\sum_{f \in \mathbb{F}} n_f = \lambda$, completing the proof. \square

With this lemma, the cost of **Bulk-Find** can be stated as follows.

Lemma 10

Bulk-Find(U, S) does $O(|R_\cup|)$ work and has $O(\text{polylog}(n))$ depth.

Proof

The R_i 's, F_i 's, and D_i 's can be maintained directly as arrays. The *roots* and *visited* sets can be maintained as bit flags on top of the vertices of U as all we need are setting the bits (adding/removing elements) and reading their values (membership testing). There are two phases in this algorithm. In Phase I, the cost of adding F_i to *visited* in iteration i is bounded by $|R_i|$. Using standard parallel operations [11], the work of the other steps is clearly bounded by $|R_{i+1}|$, including `mkFrontier` because `removeDup` does work linear in the input, which is bounded by $|R_{i+1}|$. Thus, the work of Phase I is at most $O(\sum_i |R_i|) = O(|R_\cup|)$. In terms of depth, because the union-find tree has depth at most $O(\log n)$, the **while** loop can proceed for at most $O(\log n)$ times. Each iteration involves standard operations with depth at most $O(\log^2 n)$, so the depth of Phase I is at most $O(\log^3 n)$.

In Phase II, the dominant cost comes from expanding D_i into D_{i+1} by calling `RD.allFrom`. By Lemma 9, across all iterations, the work caused by `RD.allFrom`, run on each vertex once, is expected $O(|R_\cup|)$, and the depth is $O(\text{polylog}(|R_\cup|)) \leq O(\text{polylog}(n))$. Overall, the algorithm requires $O(|R_\cup|)$ work and $O(\text{polylog}(n))$ depth. \square

5.2. Bulk-Find's Cost Equivalence to Serial find

In analyzing the work bound of the improved data structure, we will show that what Bulk-Find does is equivalent to some sequential execution of the standard `find` and requires the same amount of work, up to constants.

To gather intuition, we will manually derive such a sequence for the sample queries $S = \{1, 3, 4, 7\}$ used in Figure 1. The query of 4 went all the way to the root without merging with another flow. But the queries of 1 and 7 were stopped at 4 and in this sense, depended upon the response from the query of 4. By the same reasoning, because the query of 3 merged with the query of 1 (with 1 proceeding on), the query of 3 depended on the response from the query of 1. Note that in this view, although the query of 3 technically waited for the response at 2, it was the query of 1 that brought the response, so it depended on 1. To derive a sequence execution, we need to respect the “depended on” relation: if a depended on b , then a will be invoked after b . As an example, one sequential execution order that respects these dependencies is `find(4)`, `find(7)`, `find(1)`, `find(3)`.

We can check that by applying finds in this order, the paths traversed are exactly what the parallel execution does as `U.find` performs full path compression.

We formalize this idea in the following lemma:

Lemma 11

For a sequence of queries S with which `Bulk-Find(U, S)` is invoked, there is a sequence S' that is a permutation of S such that applying `U.find` to S' serially in that order yields the same union-find forest as Bulk-Find's and incurs the same traversal cost of $O(|R_\cup|)$, where R_\cup is as defined in the Bulk-Find algorithm.

Proof

For this analysis, we will associate every $(parent, child) \in R_\cup$ with a query $q \in S$. Logically, every query $q \in S$ starts a flow at q ascending up the tree. If there are multiple flows reaching the same node, `removeDup` inside `mkFrontier` decides which flow to go on. From this view, for any nonroot node u appearing in R_\cup , there is *exactly* one query flow from this node that proceeds up the tree. We will denote this flow by `own(u)`. All the auxiliary graphs mentioned in this proof are only for analysis purposes; they are never constructed in the execution of the algorithm.

If a query flow is stopped partway without reaching the corresponding root, the reason is either it merged in with another flow (via `mkFrontier`) or it recognized another flow that traversed the same path before (via `visited`). For every query q that is stopped partway, let $r(q)$ be the furthest point in the tree it has advanced to, i.e., $r(q)$ is the endpoint of the maximal path in R_\cup for the query flow q .

In this set up, a query flow whose furthest point is u will depend on the response from the query `own(u)`. Therefore, we form a dependency graph G_{dep} (“ u depends on v ”) as follows. The vertices are all the vertices from S . For every query flow q that is stopped partway, there is an arc `own($r(q)$)` \rightarrow q .

Let S' be a topologically-ordered sequence of G_{dep} . Multiple copies of the same query vertex can simply be placed next to each other. If we apply `U.find` serially on S' , then all queries that a query vertex q depends on in G_{dep} will have been called prior to `U.find(q)`. Because of full path compression, this means that `U.find(q)` will follow $u \rightsquigarrow r(q) \rightarrow t$ (note: $r(q) \rightarrow t$ is one step), where t is the root of the tree. Hence, every `U.find(q)` traverses the same number of edges as $u \rightsquigarrow r(q)$ plus 1. As every R_\cup edge is part of a query flow, we conclude that the work of running `U.find` on S' in that order is $O(|R_\cup|)$. \square

Finally, to obtain the bounds in Theorem 8, we modify `Simple-Bulk-Same-Set` and `Simple-Bulk-Union` (in the relabeling step) to use `Bulk-Find` on all query pairs. The depth clearly remains $O(\text{polylog}(n))$ per bulk operation. Aggregating the cost of `Bulk-Find` across calls from `Bulk-Union` and `Bulk-Same-Set`, we know from Lemma 11 that there is a sequential order that has the same work. Therefore, the total work is bounded by $O((m + q)\alpha(m + q, n))$.

6. IMPLEMENTATION AND EVALUATION

This section discusses an implementation of the proposed data structure and its empirical performance, evaluated using a number of large graphs of varying degree distributions.

6.1. Implementation

With an eye towards a simple implementation that delivers good practical performance, we set out to implement the simple bulk-parallel data structure from Section 4. The underlying union-find data structure U maintains two arrays of length n —`parent` and `sizes`—one storing a parent pointer for each vertex, and the other tracking the sizes of the trees. The `find` and `union` operations follow a standard textbook implementation. On top of these operations, we implemented `Simple-Bulk-Same-Set` and `Simple-Bulk-Union` as described earlier in the paper. We use standard sequence manipulation operations (e.g., `filter`, `prefix sum`, `pack`, `remove duplicate`) from the PBBS library [19]. There are two modifications that we made to improve the performance:

Path Compression: We wanted some benefits of path compression but without the full complexity of the work-efficient parallel algorithm from Section 5, to keep the code simple. We settled with the following pragmatic solution: The `find` operations inside `Simple-Bulk-Same-Set` and `Simple-Bulk-Union` still run independently in parallel. But after finding the root, each operation traverses the tree one more time to update all the nodes on the path to point to the root. This leads to shorter paths for later bulk operations with clear performance benefits. However, for large minibatches, the approach may still perform significantly more work than the work-efficient solution because the path compression from one `find` operation may not benefit other `find` operations within the same minibatch.

Connected Components: The algorithm as described uses as a subroutine a linear-work parallel algorithm to find connected components. These linear work algorithms expect a graph representation that gives quick random access to the neighbors of a vertex. We found the processing cost to meet this requirement to be very high and instead implemented the algorithm for connectivity described in Blleloch et al. [3]. Although this has worse theoretical guarantees, it can work with a sequence of edges directly and delivers good real-world performance.

6.2. Experimental Setup

Environment: We performed experiments on an Amazon EC2 instance with 20 cores (allowing for 40 threads via hyperthreading) of 2.4 GHz Intel Xeon E5-2676 v3 processors, running Linux 3.11.0-19 (Ubuntu 14.04.3). We believe this represents a baseline configuration of midrange workstations available in a modern cluster. All programs were compiled with Clang version 3.4 using the flag

Graph	#Vertices	#Edges	Notes
3Dgrid	99.9M	300M	3-d mesh
random	100M	500M	5 randomly-chosen neighbors per node
local5	100M	500M	small separators, avg. degree 5
local16	100M	1.6B	small separators, avg. degree 16
rMat5	134M	500M	power-law graph using rMat [4]
rMat16	134M	1.6B	a denser rMat graph

Table I. Characteristics of the graph streams used in our experiments, showing for every dataset, the total number of nodes (n), the total number of edges (m), and a brief description.

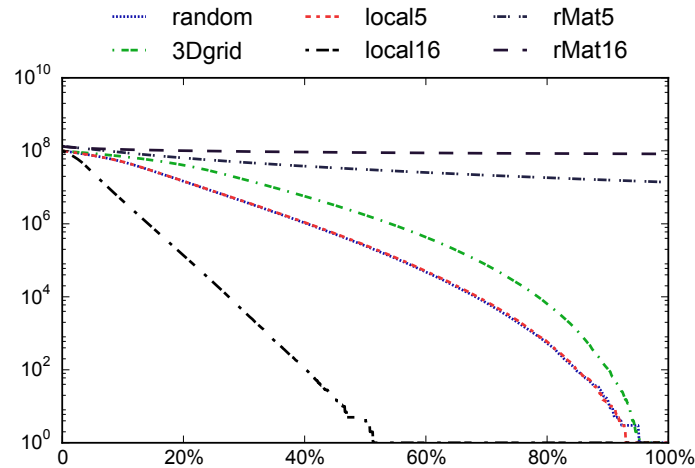


Figure 2. The numbers of connected components for each graph dataset at different percentages of the total graph stream processed.

-03. This version of Clang has the Intel Cilk runtime, which uses a work-stealing scheduler known to impose only a small overhead on both parallel and sequential code. We report wall-clock time measured using `std::chrono::high_resolution_clock`.

For robustness, we perform *three* trials and report the median running time. Although there is randomness involved in the connected component (CC) algorithm, we found no significant fluctuations in the running time across runs.

Datasets: Our experiments aim to study the behavior of the algorithm on a variety of graph streams. To this end, we use a collection of synthetic graph streams created using well-accepted generators. We include both power-law-type graphs and more regular graphs in the experiments. These are graphs commonly used in dynamic/streaming graph experiments (e.g., [14]). A summary of these datasets appear in Table I.

The graph streams in our experiments differ substantially in how quickly they become connected. This input characteristic influences the data structure's performance. Figure 2 shows for each graph stream, the number of connected components at different points in the stream. The local16 graph becomes fully connected right around the midpoint of the stream. Both rMat5 and rMat16 continue to have tens of millions of components after consuming the whole stream. Note that in this figure, random and local5 are almost visually indistinguishable until the very end.

<i>Graph</i>	UF	UF	Bulk-Parallel Using Batch Size			
	(no p.c.)	(p.c.)	500K	1M	5M	10M
random	44.63	18.42	65.43	66.57	75.20	77.89
3Dgrid	30.26	14.37	61.10	62.00	71.74	75.07
local5	44.94	18.51	65.84	66.77	75.33	78.23
local16	154.40	46.12	114.34	108.92	114.80	117.55
rMat5	33.39	18.47	66.98	68.48	74.97	78.69
rMat16	81.74	35.29	83.27	76.64	76.03	77.62

Table II. Running times (**in seconds**) on 1 thread of the baseline union-find implementation (UF) with and without path compression and the bulk-parallel version as the batch size is varied.

Baseline: Because most prior algorithms in this space either focus on parallel graphs or streaming graphs, not parallel streaming graphs, we directly compare our algorithms with sequential union find implementations (denoted UF), using both the union by size variant and a time-optimal path compression variant. We note that the algorithm of McColl et al. [14] that works in the parallel dynamic setting is not directly comparable to ours. Their algorithm focuses on supporting arbitrary insertion and deletion of edges, whereas ours is designed to take advantage of the insert-only setting.

6.3. Results

How does the bulk-parallel data structure perform on a multicore machine? To this end, we investigate the parallel overhead, speedup, and scalability trend:

- *Parallel Overhead:* How much slower is the parallel implementation running on *one* core compared to a sequential union-find? When running on one core, the parallel implementation should not take much longer than the sequential implementation, showing empirical “work efficiency.”
- *Parallel Speedup:* How much faster is the parallel version on many cores compared to the parallel version on one core? A good parallel algorithm should achieve a good speedup.
- *Scalability Trend:* How does the algorithm scale with more cores? Complementing the previous analysis, this helps predict performance on a machine with more cores.

What is the parallel overhead of our implementation? Table II shows the timings for the baseline sequential implementation of union-find UF with and without path compression and the bulk-parallel implementation *on a single thread* for four different batch sizes, 500K, 1M, 5M, and 10M. Notice that the sequential UF implementations do not depend on the batch size. To measure overhead, we first compare our implementation to union find *without* path compression: our implementation is between 1.01x and 2.5x slower except on local16, in which the bulk parallel achieves some speedups even on one thread. This is mainly because the number of connected components in local16 drops quickly to 1 as soon as midstream (Figure 2). With only 1 connected component, there is little work for bulk-parallel after that. Compared to union find with path compression, our implementation, which does pragmatic path compression, shows nontrivial—but still acceptable—overhead, as to be expected because our solution does not fully benefit finds within the same minibatch.

How much does parallelization help increase processing throughput? Table III shows the average throughputs (million edges/second) of Bulk-Union for different batch sizes. Here T_1 denotes the throughput on 1 thread and T_{20c} the throughput using 40 (hyper-) threads (the machine has only

Graph	Using $b = 500K$			Using $b = 1M$			Using $b = 5M$			Using $b = 10M$		
	T_1	T_{20c}	T_{20c}/T_1	T_1	T_{20c}	T_{20c}/T_1	T_1	T_{20c}	T_{20c}/T_1	T_1	T_{20c}	T_{20c}/T_1
random	7.64	36.87	4.8x	7.51	46.02	6.1x	6.65	60.66	9.1x	6.42	63.90	10.0x
3Dgrid	4.91	27.97	5.7x	4.83	34.97	7.2x	4.18	44.27	10.6x	3.99	45.24	11.3x
local5	7.59	38.41	5.1x	7.49	48.32	6.5x	6.64	64.61	9.7x	6.39	64.09	10.0x
local16	13.99	78.83	5.6x	14.69	95.57	6.5x	13.94	122.69	8.8x	13.61	122.03	9.0x
rMat5	7.47	26.08	3.5x	7.30	34.19	4.7x	6.67	49.92	7.5x	6.35	50.37	7.9x
rMat16	19.21	54.94	2.9x	20.88	78.10	3.7x	21.05	143.63	6.8x	20.61	167.68	8.1x

Table III. Average throughput (in million edges/second) and speedup of Bulk-Union for different batch sizes b , where T_1 is throughput on 1 thread and T_{20c} is the throughput on 20 cores.

20 physical cores). We also show the speedup factor, as given by T_{20c}/T_1 . We observe consistent speedup on all six datasets under all four batch sizes. Across all datasets, the general trend is that the larger the batch size, the higher the speedup factor. This is to be expected since a larger batch size means more work per core to fully utilize the cores, and less synchronization overhead.

How does the implementation scale with more cores? Figure 3 shows the average throughput (edges/sec) as the number of threads increases from 1 to 20 threads, then to 40 (denoted as 20c). Three different batch sizes were used for the experiments: 1M, 5M and 10M. The top chart show the results on the random dataset, the middle chart on the local16 dataset, and the bottom chart on the rMat16 dataset. Overall, the speedup is significant: with 40 (hyper-) threads, a batch size of 10M edges yields speedups between 8–11x. Smaller batch sizes yield somewhat lower speedups. In general, it is clear that *as the number of threads increases, the average throughput increases for all batch sizes*. An exception to this trend is the rMat16 dataset (bottom) using 1M batch size: the throughput drops slightly when the number of threads increases beyond 20. This is because rMat16 is sparsely connected, as is evident from the high number of components after consuming the whole stream (Figure 2), which in fact has the same order of magnitude as the number of components at the beginning. As a result, in a given minibatch, a large fraction of the edges are links between existing connected components and are discarded quickly at the start of the minibatch, resulting in relatively small work done per minibatch. Hence, at this batch size, there is not enough work to realize additional benefits beyond 20 threads.

7. CONCLUSION

We have presented a shared-memory parallel data structure for the union-find problem (hence, incremental graph connectivity) in the minibatch arrival model. Our algorithm has polylogarithmic parallel depth and its total work across all processors is of the same order as the work due to the best sequential algorithm for the problem. We also presented a simpler parallel algorithm that is easier to implement and has good practical performance.

This still leaves open several natural problems. We list some of them here. (1) In case all edge updates are in a single minibatch, the total work of our algorithm is (in a theoretical sense), superlinear in the number of edges in the graph. Whereas, the optimal batch algorithm for graph

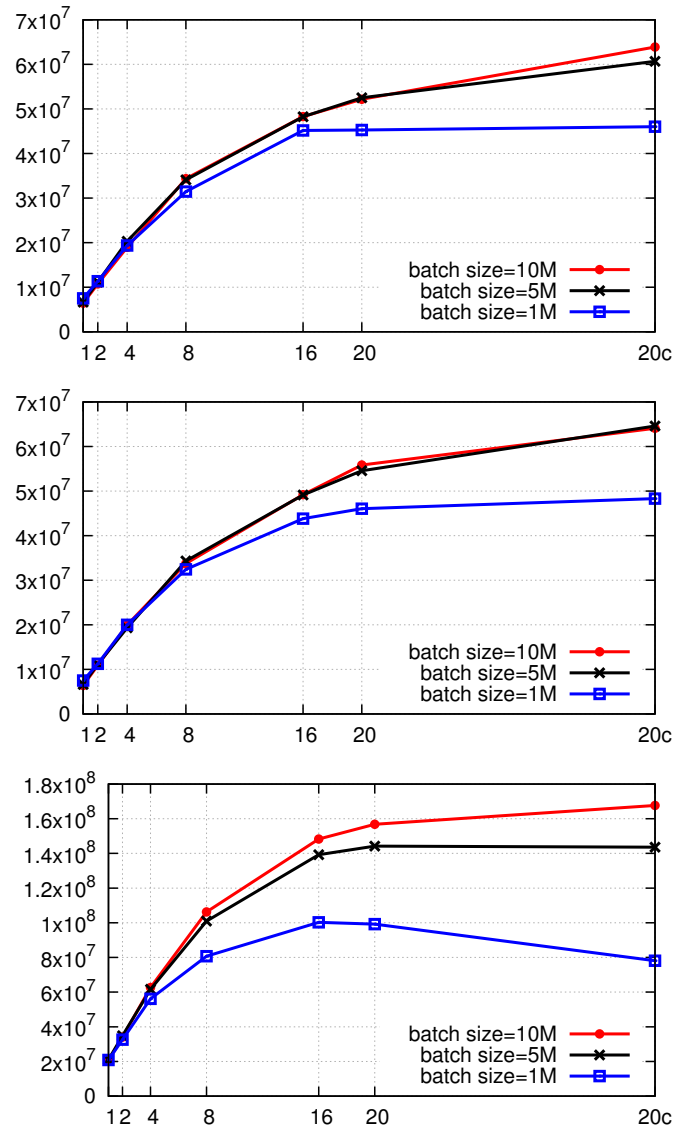


Figure 3. Average throughput (**edges per second**) as the number of threads is varied from 1 to 40 (denoted by 20c as they run on 20 cores with hyperthreading). The graph streams shown are (top) random, (middle) local16, and (bottom) rMat16.

connectivity, based on a depth-first search, has work linear in the number of edges. Is it possible to have an incremental algorithm whose work is linear in the case of very large batches, such as the above, and falls back to the union-find type algorithms for smaller minibatches? Note that for all practical purposes, the work of our algorithm is linear in the number of edges, due to very slow growth of the inverse Ackerman's function. (2) Can these results on parallel algorithms be extended to the fully dynamic case when there are both edge arrivals as well as deletions?

ACKNOWLEDGMENT

This material is based upon work partially supported by the National Science Foundation under grant number IIS 1527541.

REFERENCES

- [1] Richard J. Anderson and Heather Woll. Wait-free parallel algorithms for the union-find problem. In *Proceedings of the 23rd Annual ACM Symposium on Theory of Computing, May 5-8, 1991, New Orleans, Louisiana, USA*, pages 370–380, 1991.
- [2] Jonathan Berry, Matthew Oster, Cynthia A. Phillips, Steven Plimpton, and Timothy M. Shead. Maintaining connected components for infinite graph streams. In *Proc. 2nd International Workshop on Big Data, Streams and Heterogeneous Source Mining: Algorithms, Systems, Programming Models and Applications (BigMine)*, pages 95–102, 2013.
- [3] Guy E. Blelloch, Jeremy T. Fineman, Phillip B. Gibbons, and Julian Shun. Internally deterministic parallel algorithms can be fast. In *Proceedings of the 17th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming, PPOPP 2012, New Orleans, LA, USA, February 25-29, 2012*, pages 181–192, 2012.
- [4] Deepayan Chakrabarti, Yiping Zhan, and Christos Faloutsos. R-MAT: A recursive model for graph mining. In *Proceedings of the Fourth SIAM International Conference on Data Mining, Lake Buena Vista, Florida, USA, April 22-24, 2004*, pages 442–446, 2004.
- [5] Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. *Introduction to Algorithms* (3. ed.). MIT Press, 2009.
- [6] Camil Demetrescu, Bruno Escoffier, Gabriel Moruz, and Andrea Ribichini. Adapting parallel algorithms to the W-stream model, with applications to graph problems. *Theor. Comput. Sci.*, 411(44-46):3994–4004, 2010.
- [7] Camil Demetrescu, Irene Finocchi, and Andrea Ribichini. Trading off space for passes in graph streaming problems. *ACM Transactions on Algorithms*, 6(1), 2009.
- [8] Joan Feigenbaum, Sampath Kannan, Andrew McGregor, Siddharth Suri, and Jian Zhang. On graph problems in a semi-streaming model. *Theor. Comput. Sci.*, 348(2-3):207–216, 2005.
- [9] Hillel Gazit. An optimal randomized parallel algorithm for finding connected components in a graph. *SIAM J. Comput.*, 20(6):1046–1067, 1991.
- [10] IBM Corporation. Infosphere streams. <http://www-03.ibm.com/software/products/en/infosphere-streams>. Accessed Jan 2014.
- [11] Joseph JáJá. *An Introduction to Parallel Algorithms*. Addison-Wesley, 1992.
- [12] Bruce M. Kapron, Valerie King, and Ben Mountjoy. Dynamic graph connectivity in polylogarithmic worst case time. In *Proceedings of the Twenty-Fourth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2013, New Orleans, Louisiana, USA, January 6-8, 2013*, pages 1131–1142, 2013.
- [13] Fredrik Manne and Md. Mostofa Ali Patwary. A scalable parallel union-find algorithm for distributed memory computers. In *Parallel Processing and Applied Mathematics, 8th International Conference, PPAM 2009, Wroclaw, Poland, September 13-16, 2009. Revised Selected Papers, Part I*, pages 186–195, 2009.
- [14] Robert McColl, Oded Green, and David A. Bader. A new parallel algorithm for connected components in dynamic graphs. In *20th Annual International Conference on High Performance Computing, HiPC 2013, Bengaluru (Bangalore), Karnataka, India, December 18-21, 2013*, pages 246–255, 2013.
- [15] Andrew McGregor. Graph stream algorithms: a survey. *SIGMOD Record*, 43(1):9–20, 2014.
- [16] Md. Mostofa Ali Patwary, Peder Refsnes, and Fredrik Manne. Multi-core spanning forest algorithms using the disjoint-set data structure. In *26th IEEE International Parallel and Distributed Processing Symposium, IPDPS*, pages 827–835, 2012.
- [17] Sanguthevar Rajasekaran and John H. Reif. Optimal and sublogarithmic time randomized parallel sorting algorithms. *SIAM J. Comput.*, 18(3):594–607, 1989.
- [18] Raimund Seidel and Micha Sharir. Top-down analysis of path compression. *SIAM J. Comput.*, 34(3):515–525, 2005.
- [19] Julian Shun, Guy E. Blelloch, Jeremy T. Fineman, Phillip B. Gibbons, Aapo Kyrola, Harsha Vardhan Simhadri, and Kanat Tangwongsan. Brief announcement: the problem based benchmark suite. In *24th ACM Symposium on Parallelism in Algorithms and Architectures, SPAA '12, Pittsburgh, PA, USA, June 25-27, 2012*, pages 68–70, 2012.

- [20] Julian Shun, Laxman Dhulipala, and Guy E. Blelloch. A simple and practical linear-work parallel algorithm for connectivity. In *26th ACM Symposium on Parallelism in Algorithms and Architectures, SPAA '14, Prague, Czech Republic*, pages 143–153, 2014.
- [21] Natcha Simsiri, Kanat Tangwongsan, Srikanta Tirthapura, and Kun-Lung Wu. Work-efficient parallel union-find with applications to incremental graph connectivity. In *Proceedings of the 22nd International European Conference on Parallel and Distributed Computing*, pages 561–573, 2016.
- [22] Robert Endre Tarjan. Efficiency of a good but not linear set union algorithm. *Journal of the ACM*, 22(2):215–225, 1975.
- [23] M. Zaharia, T. Das, H. Li, T. Hunter, S. Shenker, and I. Stoica. Discretized streams: Fault-tolerant streaming computation at scale. In *Proc. ACM Symposium on Operating Systems Principles (SOSP)*, pages 423–438, 2013.