

# GRAPHZEPPELIN: Storage-Friendly Sketching for Connected Components on Dynamic Graph Streams.

David Tench  
Stony Brook University  
dtench@cs.stonybrook.edu

Michael A. Bender  
Stony Brook University

Martin Farach-Colton  
Rutgers University

Evan West  
Stony Brook University

Abiyaz Chowdhury  
Stony Brook University

Tyler Seip  
MongoDB

Victor Zhang  
Rutgers University

J. Ahmed Dellas  
Rutgers University

Kenny Zhang  
Stony Brook University

## ABSTRACT

Finding the connected components of a graph is a fundamental problem with uses throughout computer science and engineering. The task of computing connected components becomes more difficult when graphs are very large, or when they are dynamic, meaning the edge set changes over time subject to a stream of edge insertions and deletions. A natural approach to computing the connected components problem on a large, dynamic graph stream is to buy enough RAM to store the entire graph. However, the requirement that the graph fit in RAM is an inherent limitation of this approach and is prohibitive for very large graphs. Thus, there is an unmet need for systems that can process dense dynamic graphs, especially when those graphs are larger than available RAM.

We present a new high-performance streaming graph-processing system for computing the connected components of a graph. This system, which we call GRAPHZEPPELIN, uses new linear sketching data structures (CUBEskETCH) to solve the streaming connected components problem and as a result requires space asymptotically smaller than the space required for a lossless representation of the graph. GRAPHZEPPELIN is optimized for massive dense graphs: GRAPHZEPPELIN can process millions of edge updates (both insertions and deletions) per second, even when the underlying graph is far too large to fit in available RAM. As a result GRAPHZEPPELIN vastly increases the scale of graphs that can be processed.

## 1 INTRODUCTION

Finding the connected components of a graph is a fundamental problem with uses throughout computer science and engineering. A recent survey by Sahu et al. [62] of industrial uses of algorithms reports that, for both practitioners and academic researchers, connected components was the most frequently performed computation from a list of 13 fundamental graph problems that includes shortest paths, triangle counting, and minimum spanning trees. It has applications in scientific computing [58, 65], flow simulation [66], metagenome assembly [25, 54], identifying protein families [49, 72], analyzing cell networks [4], pattern recognition [28, 35], graph partitioning [43, 44], random walks [33], social network community detection [39], graph compression [34, 42], medical imaging [30], and object recognition [29]. It is a starting point for strictly harder problems such as edge/vertex connectivity, shortest paths, and  $k$ -cores. It is used as a subroutine for pathfinding algorithms

like Dijkstra and  $A^*$ , some minimum spanning tree algorithms, and for many approaches to clustering [21–23, 57, 70, 71].

The task of computing connected components becomes more difficult when graphs are very large, or when they are *dynamic*, meaning the edge set changes over time subject to a stream of edge insertions and deletions. Applications on large graphs include metagenome assembly tasks that may include hundreds of millions of genes with complex relations [25], and large-scale clustering which is a common machine learning challenge [23]. Applications using dynamic graphs include identifying objects from a video feed rather than a static image [36], or tracking communities in social networks that change as users add or delete friends [9, 10]. And of course graphs can be both large and dynamic. Indeed, Sahu et al.’s [62] survey reports that a majority of industry respondents work with large graphs ( $> 1$  million nodes or  $> 1$  billion edges) and a majority work with graphs that change over time.

A natural approach to computing the connected components problem on a large, dynamic graph stream is to buy enough RAM to store the entire graph. Indeed, dynamic graph stream processing systems such as Aspen and Terrace [19, 56] can efficiently query the connected components of a large graph subject to a stream of edge insertions and deletions. However, the requirement that the graph fit in RAM is an inherent limitation of this approach and is prohibitive for most large graphs. For example, a graph with ten million nodes and an average degree of 1 million, using 2B to encode an edge, would require 10TB of memory. We show in Section 6 that the Aspen and Terrace graph representations are significantly larger than this lower bound.

Indeed, in public graph-data-set repositories, most graphs are smaller than typical single-machine RAM sizes. As Figure 1 illustrates, nearly all graphs in Network Repository [60] can be stored as an adjacency list in less than 16GB. This fixed memory budget furthermore implies that graphs with large numbers of vertices must be sparse. Similarly, the Stanford SNAP graph repository and the SuiteSparse repository have few graphs larger than 16GB, and graphs with many nodes are always extremely sparse.

Large, dense graphs, we argue, are absent from graph repositories not because they are unworthy of study, but because there are few tools to analyze them. To illustrate: dense graphs do appear in Network Repository [60], but these graphs are never larger than a few GB; moreover, as the graphs’ vertex count increases, the maximum density decreases such that the densest graphs never require



**Figure 1: Published graphs have few nodes or are sparse. Each point represents a graph data set from NetworkRepository. Any point below the dark line indicates a graph that can be represented as an adjacency list in 16GB of RAM.**

more than 10 GB. A compelling explanation for the absence of large, dense graphs is selection bias: interesting dense graphs exist at all scales, but large, dense graphs are discarded as computationally infeasible and consequently are rarely published or analyzed.

Thus, there is an unmet need for systems that can process dense graphs, especially when those graphs are larger than available RAM. Existing systems are not designed for large, dense, dynamic graph streams, and instead optimize for other use cases. Aspen and Terrace are optimized for large, sparse, dynamic graphs that completely fit in RAM, but their performance degrades significantly on dense graphs larger than RAM. There is a deep literature on parallel systems for connected components computation in multicore [26], GPU [5], and distributed settings [12, 37] but these focus on static graphs which fit in RAM. Many external memory [11] and semi-external memory [1] systems focus on graphs that are too large for RAM and must be stored on disk, but none of these systems focus on graphs whose edges can be deleted dynamically.

In this paper, we explore the general problem of connected components on large, dense, dynamic graphs. We introduce GRAPHZEPPELIN, which computes the connected components of graph streams using space asymptotically smaller than an explicit representation of the graph. GRAPHZEPPELIN uses a new  $\ell_0$ -sketching data structure that outperforms the state of the art on graph sketching workloads. Additionally, GRAPHZEPPELIN employs node-based buffering strategies which improve I/O efficiency. These techniques allow GRAPHZEPPELIN to scale better than existing systems in several settings. First, for in-RAM computation, GRAPHZEPPELIN’s small size means it can process larger, denser graphs than Aspen or Terrace given the same amount of RAM. Moreover, even if the input graph fits in RAM on all systems, GRAPHZEPPELIN is 28% faster than Aspen and 26 times faster than Terrace on large dense graphs. Finally, GRAPHZEPPELIN scales to SSD at the cost of a minor decrease to ingestion rate, and is more than two orders of magnitude faster than Aspen and Terrace which suffer significant performance degradation when scaling out of RAM.

GRAPHZEPPELIN employs a new sketch algorithm, overcoming a computational bottleneck of existing linear sketching techniques in the semi-streaming graph algorithms literature [18]. The asymptotically best existing streaming connected components algorithm

is Ahn et al.’s STREAMINGCC [2, 52], which has asymptotically low space and update time complexity. STREAMINGCC relies on  $\ell_0$ -sampling, which it uses to sample edges across arbitrary graph cuts. However, the best known  $\ell_0$ -sampling algorithm suffers from high constant and polylogarithmic factors in its space and update time, as we show in Section 3. This overhead makes any implementation of the STREAMINGCC data structure infeasibly slow and large. GRAPHZEPPELIN employs what we call CUBESKETCH, a specialized  $\ell_0$ -sampling algorithm for sampling edges across graph cuts, to solve the connected components problem. For large graphs CUBESKETCH uses 8 times less space than the best general  $\ell_0$ -sampling algorithm and can process updates more than three orders of magnitude faster.

GRAPHZEPPELIN also uses new write-optimized data structures to overcome prohibitive resource requirements of existing semi-streaming algorithms. Streaming algorithms have had a significant impact in large part because they require a small (polylogarithmic) amount of RAM. In contrast, graph semi-streaming algorithms have higher RAM requirements: for most problems on a graph with  $V$  nodes, sublinear RAM is insufficient to even represent a solution so  $O(V \text{polylog}(V))$  RAM is typically assumed. With the large polylog factors, this is often more RAM than is feasible in practice; see Section 2. We propose a computational setting, the hybrid streaming model, that enjoys the memory advantage of the streaming model while allowing enough space in external memory to compute on dynamic graph streams. In this model there is still  $O(V \text{polylog}(V))$  space available, but only  $O(\text{polylog}(V))$  of this space is RAM and the rest is disk which may only be accessed in  $O(\text{polylog}(V))$ -size blocks. The dual challenges in this model are to design algorithms that use small total space but also that have fast update and query time complexity. While existing graph semi-streaming algorithms satisfy the small space requirement for this setting, their heavy reliance on hashing and random access patterns make them slow on disk. We show that GRAPHZEPPELIN is simultaneously a space-optimal in-RAM semi-streaming algorithm, and an I/O efficient external memory algorithm for the connected components problem. We also validate its performance experimentally, showing that GRAPHZEPPELIN can operate on modern consumer solid-state disk, increasing the scale of dynamic graph streams that it can process while incurring a 23% cost to stream ingestion rate.

**Results.** In this paper we establish the following:

- **GRAPHZEPPELIN:** We present a new high-performance streaming graph-processing system for computing the connected components of a graph. This system, which we call GRAPHZEPPELIN, uses new linear sketching data structures (CUBESKETCH, described below) to solve the streaming connected components problem and as a result requires space asymptotically smaller than the space required for a lossless representation of the graph. GRAPHZEPPELIN is optimized for massive dense graphs: GRAPHZEPPELIN can process millions of edge updates (both insertions and deletions) per second, even when the underlying graph is far too large to fit in available RAM. As a result GRAPHZEPPELIN vastly increases the scale of graphs that can be processed.
- **CUBESKETCH:  $\ell_0$ -sampling optimized for graph connectivity sketching.** We give a new  $\ell_0$ -sampling algorithm,

CUBESKETCH, for vectors of integers mod 2. Given a vector of length  $n$ , CUBESKETCH uses  $O(\log^2(n))$  words of space and  $O(\log^2(n))$  time per update, which is a factor of  $O(\log(n))$  faster than the best existing  $\ell_0$ -sampler for general vectors [18].

CUBESKETCH is a key subroutine in GRAPHZEPPELIN, where it is used to sample graph edges across arbitrary cuts as part of connected components computation. Here it is used to sketch vectors of length  $\binom{V}{2} = O(V^2)$ , where  $V$  denotes the number of nodes in the graph. We show experimentally that CUBESKETCH is more than 3 orders of magnitude faster than the state-of-the-art  $\ell_0$  sampling algorithm on graph streaming workloads.

In addition to the  $O(\log(V))$ -factor speedup, several non-asymptotic factors contribute to this performance improvement as well. First, the existing algorithm's update cost is dominated by  $O(\log^3(V))$  division operations, while CUBESKETCH's update cost is dominated by  $O(\log^2(V))$  bitwise XOR operations, which are much faster. In addition, the general algorithm performs 128-bit arithmetic operations (including division) when processing graphs with more than  $10^5$  nodes, whereas CUBESKETCH can use standard 64-bit operations to achieve the same polynomially small error probability. Finally, both algorithms match the asymptotic space lower bound but CUBESKETCH uses roughly 8 times less space than the general algorithm.

- **Asymptotic guarantees of GRAPHZEPPELIN: space-optimality, I/O efficiency,  $O(\log^3(V))$  time per update.** GRAPHZEPPELIN's core algorithm matches the  $O(V \log^3(V))$  space lower bound for the streaming connected components problem, and its per-update time cost of  $O(\log^3(V))$  is  $O(\log(V))$  times faster than the best existing algorithm [2]. Additionally, GRAPHZEPPELIN can efficiently ingest stream updates even when its sketch data structure is too large to fit in RAM: its I/O complexity is  $\text{sort}(\text{length of stream}) + O(V/B \log^3(V) + V \log^*(V))$  and for realistic block sizes it is an I/O-optimal external-memory connected components algorithm [16]. Given a fixed amount of RAM and disk, GRAPHZEPPELIN is capable of efficiently computing the connected components of larger graphs than existing algorithms in the streaming or external memory models.
- **Empirical achievements of GRAPHZEPPELIN: better scaling for in-memory, out-of-core, and parallel computation, and undetectable failure probability.** GRAPHZEPPELIN's CUBESKETCH-based design increases the size of input graphs that can be processed, scales well to persistent memory, and facilitates parallelism in stream ingestion. As a result, GRAPHZEPPELIN can ingest 2-4 million edge updates per second on a single scientific workstation (see Section 6), both when its data structures reside completely in RAM and also when they reside in fast persistent memory. As a result of these advantages, GRAPHZEPPELIN is faster and more scalable than the state of the art on large, dense graphs.
- **GRAPHZEPPELIN handles larger graphs for in-RAM computation.** GRAPHZEPPELIN's space-efficient CUBESKETCH allows it to process graph streams larger than can be stored explicitly in a fixed amount of RAM and give it an asymptotic  $O(V/\log^3(V))$  space advantage over state-of-art systems

on dense graphs. Given the polylogarithmic factors and constants, we need to determine the actual crossover point where GRAPHZEPPELIN processes graphs more compactly than Aspen and Terrace. We show empirically that this crossover point occurs when the space budget is between 32 and 64 gigabytes. That is, for dense graphs on several hundred thousand nodes, GRAPHZEPPELIN is 25% more compact than Aspen and several times more compact than Terrace, and this scaling advantage only increases for larger space budgets or input sizes. Additionally, for suitably large and dense graph streams GRAPHZEPPELIN ingests updates roughly 20 times faster than Terrace and 33% faster than Aspen, even when each system's data structures fit entirely within RAM.

- **GRAPHZEPPELIN can use persistent memory to handle even larger graphs.** GRAPHZEPPELIN's node-based work buffering strategy facilitates out-of-core computation, allowing GRAPHZEPPELIN to use SSD to increase the scale of graph streams it can process while incurring a small cost to performance. We show experimentally that GRAPHZEPPELIN ingests updates more than two orders of magnitude faster than Aspen and Terrace when all systems swap to disk, and that using SSD slows GRAPHZEPPELIN stream ingestion by only 23%.
- **GRAPHZEPPELIN's stream ingestion is highly parallel.** GRAPHZEPPELIN employs a node-based work buffering strategy that facilitates parallelism and improves data locality. We show experimentally that GRAPHZEPPELIN's multithreaded stream ingestion system scales well with more threads: its ingestion rate is 25 times higher with 46 threads than an optimized single-thread implementation.
- **GRAPHZEPPELIN's theoretical failure probability is undetectable in practice.** GRAPHZEPPELIN and similar graph sketching approaches achieve their remarkable space efficiency at the cost of a random chance of failure. We show empirically that GRAPHZEPPELIN's observed failure rate is even lower than the proved (polynomially small) upper bound: in fact, for 4000 thousand trials on various real-world and synthetic graphs it never was observed to fail.

**Paper Organization.** In Section 2 we formally define the dynamic streaming connected components problem, as well as other useful concepts, and summarize Ahn et al.'s STREAMINGCC algorithm. In Section 3 we demonstrate experimentally that STREAMINGCC's core sketching data structure is too large and slow to be of use given the limitations of existing computing technology, and we present CUBESKETCH, a novel variant of this core sketching data structure and show how it alleviates the algorithm's limitations. In Section 4 we describe how this modified algorithm can be extended to efficiently make use of sequentially-accessible disk as working memory, and in doing so introduce the hybrid streaming model. In Section 5 we describe the design of GRAPHZEPPELIN in detail. In Section 6 we demonstrate experimentally that GRAPHZEPPELIN outperforms existing systems on large, dense graph streams for both RAM-only and out-of-core computation, that it is highly parallelizable, and that its algorithmic failure probability is essentially 0 in practice.

## 2 PRELIMINARIES

### 2.1 Graph Streaming & Hybrid Graph Streaming

In the *graph semi-streaming* model [24, 51] (sometimes just called the *graph streaming* model), an algorithm is presented with a *stream*  $S$  of updates (each an edge insertion or deletion) where the length of the stream is  $N$ . Stream  $S$  defines an input graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  with  $V = |\mathcal{V}|$  and  $E = |\mathcal{E}|$ . The challenge in this model is to compute (perhaps approximately) some property of  $\mathcal{G}$  given a single pass over  $S$  and at most  $O(V \text{polylog}(V))$  words of memory. Each update has the form  $((u, v), \Delta)$  where  $u, v \in \mathcal{E}$ ,  $u \neq v$  and  $\Delta \in \{-1, 1\}$  where 1 indicates an edge insertion and  $-1$  indicates an edge deletion. Let  $s_i$  denote the  $i$ th element of  $S$ , and let  $S_i$  denote the first  $i$  elements of  $S$ . Let  $\mathcal{E}_i$  be the edge set defined by  $S_i$ , i.e., those edges which have been inserted and not subsequently deleted by step  $i$ . The stream may only insert edge  $e$  at time  $i$  if  $e \notin \mathcal{E}_{i-1}$ , and may only delete edge  $e$  at time  $i$  if  $e \in \mathcal{E}_{i-1}$ .

In Section 4 we additionally use a new variant of the graph semi-streaming model, which we call the *hybrid graph streaming setting* (since it incorporates some components of the external memory model [67] into the semi-streaming model). In this setting, there is an additional constraint on the type of memory available for computation: only  $M = \Omega(\text{polylog}(V)) = o(V)$  RAM is available, and  $D = O(V \text{polylog}(V))$  disk space is available. A word in RAM is accessed at unit cost, and disk is accessed in blocks of  $B = o(M)$  words at a cost of  $B$  per access. Any semi-streaming algorithm can be run with this additional constraint, but may become much slower if the algorithm makes many random accesses to disk. The algorithmic challenge in the hybrid graph streaming setting is to minimize time complexity (of ingesting stream updates and returning solutions) in addition to satisfying the typical limited-space requirement of the data stream setting. In Section 4 we show how GRAPHZEPPELIN can be adapted to this model, resulting in an algorithm which is simultaneously a space-optimal single pass streaming algorithm with  $O(\log^3(V))$  update time and an I/O efficient external memory algorithm.

**Problem 1 (The streaming Connected Components problem.).** *Given a insert/delete edge stream of length  $N$  that defines a graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , return a insert-only edge stream that defines a spanning forest of  $\mathcal{G}$ .*

### 2.2 Prior Work in Streaming Connected Components

We summarize STREAMINGCC, Ahn et al.’s [2] semi-streaming algorithm for computing a spanning forest (and therefore the connected components) of a graph.

For each node  $v_i$  in  $G$ , define the *characteristic vector*  $a_i$  of  $v_i$  to be a 1-dimensional vector indexed by the set of possible edges in  $\mathcal{G}$ .  $a_i[(j, k)]$  is only nonzero when  $i = j$  or  $i = k$  and edge  $(j, k) \in \mathcal{E}$ . That is,  $a_i \in \{-1, 0, 1\}^{\binom{V}{2}}$  s.t. for all  $0 \leq j < k < \binom{V}{2}$ :

$$a_i[(j, k)] = \begin{cases} 1 & i = j \text{ and } (v_j, v_k) \in \mathcal{E} \\ -1 & i = k \text{ and } (v_j, v_k) \in \mathcal{E} \\ 0 & \text{otherwise} \end{cases}$$

Crucially, for any  $S \subset \mathcal{V}$ , the sum of the characteristic vectors of the nodes in  $S$  is a direct encoding of the edges across the cut  $(S, \mathcal{V} \setminus S)$ . That is, let  $x = \sum_{v \in S} a_v$  and then  $|x[(j, k)]| = 1$  iff  $(j, k) \in E(S, \mathcal{V} \setminus S)$ .

Using these vectors, we immediately have a (very inefficient) algorithm for computing the connected components from a stream: Initialize  $a_i = \{0\}^{\binom{V}{2}}$  for all  $i$ . For each stream update  $s = ((u, v), \Delta)$ , set  $a_u[u, v] += \Delta$  and  $a_v[u, v] += -\Delta$ .

After the stream, run Boruvka’s algorithm [53] for finding a spanning forest as follows. For the first round of the algorithm, from each  $a_i$  arbitrarily choose one nonzero entry  $(w, y)$  (an edge in  $\mathcal{E}$  s.t.  $w = i$  or  $y = i$ ). Add  $e_i$  to the spanning forest. For each connected component  $C$  in the spanning forest, compute the characteristic vector of  $C$ :  $a_C = \sum_{v \in C} a_v$ . Proceed similarly for the remaining rounds of Boruvka’s algorithm: in each round, choose one nonzero entry from the characteristic vector of each connected component and add the corresponding edges to the spanning forest. Sum the characteristic vectors of the component nodes of the connected components in the spanning forest, and continue until no new merges are possible. This will take at most  $\log(V)$  rounds.

The key idea to make this a small-space algorithm is to use “ $\ell_0$ -sampling” [18] to run this version of Boruvka’s algorithm by compressing each characteristic vector  $a_i$  into a data structure of size  $O(\log^2(V))$  that can return a nonzero entry of  $a_i$  with high probability.

**Definition 1.** *A sketch algorithm is a  $\delta\ell_0$ -sampler if it is*

- (1) **Sampleable:** *it can take as its input a stream of updates to the coordinates of a non-zero vector  $a$ , and output a non-zero coordinate  $(j, a[j])$  of  $a$ . We use  $S(a)$  to denote the sketch of vector  $a$ .*
- (2) **Linear:** *for any vectors  $a$  and  $b$ ,  $S(a) + S(b) = S(a + b)$  and this operation preserves sampleability, i.e.,  $S(a + b)$  can output a nonzero coordinate of vector  $a + b$ .*
- (3) **Low Failure Probability.** *the algorithm returns an incorrect or null answer with probability at most  $\delta$ .*

For all  $\ell_0$  samplers discussed in this paper,  $S(a)$  is a vector and adding two sketches is equivalent to adding their vectors element-wise.

**Lemma 1.** (Adapted from [18], Theorem 1): *Given a 2-wise independent hash family  $\mathcal{F}$  and an input vector of length  $n$ , there is an  $\ell_0$ -sampler using  $O(\log^2(n))$  space that succeeds with probability at least  $1 - 1/n^c$  for constant  $c$ .*

We denote a  $\ell_0$  sketch of a vector  $x$  as  $S(x)$ . Since the sketch is linear,  $S(x) + S(y) = S(x + y)$  for any vectors  $x$  and  $y$ . This allows us to process stream updates as follows: we maintain a running sum of the sketches of each stream update, which is equivalent to a sketch of the vector defined by the stream. That is, let  $a_i^t$  denote  $a_i$  after stream prefix  $S_t$ . For the  $j$ th stream update  $s_j = ((i, x), \Delta)$  we obtain  $S(a_i^j) = S(s_j) + S(a_i^{j-1})$ .

Linearity also allows us to emulate the merging step of Boruvka’s algorithm by summing the sketches of all nodes in each connected component. We require  $\log(V)$  independent  $\ell_0$  sketches for each  $v \in \mathcal{V}$ , one for each round<sup>1</sup>, so the size of the sketch data

<sup>1</sup>In the original paper the authors note that adaptivity concerns require the use of new sketches for each round of Boruvka’s algorithm.

structure for each node is  $O(\log^3(V))$ . We refer to the sketch data structure for each node as a **node sketch** and each of its  $\log(V)$   $\ell_0$ -subsketches as CUBESKETCHS. The total size of the entire data structure is  $O(V \log^3(V))$ . Recent work [52] has shown that this is optimal space.

### 3 $\ell_0$ -SAMPLING REVISITED

Existing  $\ell_0$ -sampling algorithms are asymptotically small and fast to update, but in practice high constant and logarithmic overheads in size and update time prevent these algorithms from being useful for a streaming connected components algorithm. We now review some details of the best known  $\ell_0$ -sampling algorithm and demonstrate experimentally that using it to emulate Boruvka’s algorithm for graph connectivity would be prohibitively slow and would require an enormous amount of space. Then we introduce an  $\ell_0$ -sketching algorithm which exploits the structure of the connected components problem to improve performance, and experimentally demonstrate that it is 16 times smaller and X orders of magnitude faster to update than the state of the art.

$\ell_0$ -sampling algorithms have two key components [18]:

- A **selection** process in which multiple subsets  $\mathcal{I}_{1,1}, \mathcal{I}_{1,2}, \dots, \mathcal{I}_{1,ck}, \mathcal{I}_{2,1}, \dots, \mathcal{I}_{k,ck}$  of indices of input vector  $x$  are selected where  $k = \log(|x|)$  and  $c$  is a constant.
- A **1-sparse recovery** process which, given index subset  $\mathcal{I}$ , returns a nonzero element of  $\mathcal{I}$  if  $\text{supp}(\mathcal{I}) = 1$ , returns NO ELEMENT if  $\text{supp}(\mathcal{I}) = 0$ , and with high probability returns FAIL if  $\text{supp}(\mathcal{I}) > 1$ .

If  $\text{supp}(x) > 0$ ,  $\exists \mathcal{I}_{i,j}$  s.t.  $\text{supp}(\mathcal{I}_{i,j}) > 0$  with probability at least  $1 - 1/|x|^c$ .

**Existing  $\ell_0$ -samplers are slow to update for graph streaming workloads.** Each index subset  $\mathcal{I}_{w,y}$  is represented by a 1-sparse recovery data structure  $r_{y,w}$ . When a stream update  $(i, \Delta)$  arrives, the  $\ell_0$ -sampler determines which of the index subsets contain  $i$  and apply the update to their 1-sparse recovery data structures. State-of-the-art  $\ell_0$ -samplers use a recovery process which requires a small constant number of machine words, but applying an update requires expensive arithmetic computation including modular exponentiation, necessitating  $O(\log(|x|))$  multiplication operations and  $O(\log(|x|))$  modulo operations (where the modulus is a large prime) for each index subset. As a result, in the worst case processing a stream update requires  $O(\log^3(|x|))$  arithmetic operations. Moreover, for sufficiently large vectors, this modular exponentiation must be done on integers larger than a 64-bit machine word, drastically increasing computation time in practice.

The “Standard  $\ell_0$ ” column of Figure 2 displays the single-threaded ingestion rate in updates per second of the state-of-the-art  $\ell_0$ -sampling algorithm for vectors of various sizes. These results were obtained on a Dell Precision 7820 with 24-core 2-way hyperthreaded Intel(R) Xeon(R) Gold 5220R CPU @ 2.20GHz and 64GB 4x16GB DDR4 2933MHz RDIMM ECC Memory. Note how ingestion rate decreases as vector length increases, and in particular there is a catastrophic slowdown at vector length  $10^{10}$ . This drastic decrease in ingestion rate is due to the need to perform modular exponentiation on integers larger than  $2^{64}$ , requiring the use of 128-bit integers which slow computation. When sketching characteristic vectors of

Vector Length	Standard $\ell_0$	CUBESKETCH	Speedup
$10^3$	174,000	7,160,000	41.1 x
$10^4$	96,200	6,590,000	68.5 x
$10^5$	45,300	6,020,000	133 x
$10^6$	25,900	5,040,000	195 x
$10^7$	18,400	4,410,000	240 x
$10^8$	14,400	4,090,000	284 x
$10^9$	11,600	3,670,000	316 x
$10^{10}$	1,270	3,290,000	2,590 x
$10^{11}$	913	3,150,000	3,450 x
$10^{12}$	806	2,890,000	3,590 x

**Figure 2: CUBESKETCH is faster than standard  $\ell_0$  sketching. Ingestion rates (in updates/second) are listed for both  $\ell_0$  sketching methods.**

Vector Length	Standard $\ell_0$	CUBESKETCH	Size Reduction
$10^3$	3.00KiB	0.75KiB	4.0 x
$10^4$	5.00KiB	1.33KiB	3.8 x
$10^5$	7.00KiB	1.80KiB	3.9 x
$10^6$	10.0KiB	2.57KiB	3.9 x
$10^7$	14.0KiB	3.63KiB	3.9 x
$10^8$	18.0KiB	4.38KiB	4.1 x
$10^9$	22.0KiB	5.55KiB	4.0 x
$10^{10}$	56.0KiB	6.98KiB	8.0 x
$10^{11}$	66.0KiB	7.31KiB	9.0 x
$10^{12}$	76.0KiB	9.72KiB	7.8 x

**Figure 3: CUBESKETCH is significantly smaller than standard  $\ell_0$  sketching. Sizes are listed for both  $\ell_0$  sketching methods.**

length  $O(V^2)$  for streaming connected components, 128-bit integers are required when  $V \geq 10^5$ .

When using  $\ell_0$ -sampling for Boruvka emulation, each stream update  $((u, v), \Delta)$  must be applied to the node sketches of  $u$  and  $v$ . For any node  $u$ , the node sketch of  $u$  is made up of  $\log(V)$   $\ell_0$ -sketches of  $a_u$  so processing a stream update requires  $2 \cdot \log(V) \cdot O(\log^3(|a_u|)) = 2 \cdot \log(V) \cdot O(\log^3(V))$  multiplication and modulo operations. For a graph with a million nodes, STREAMINGCC must apply each update to  $2 \log(V) \approx 40$  sketch vectors of length  $10^{12}$ , so it can process roughly  $800/40 = 20$  edge updates per second.

**Existing  $\ell_0$ -samplers are large for graph streaming workloads.** Each node sketch consists of  $\log(V)$  CUBESKETCHS. Each such CUBESKETCH is an  $\ell_0$ -sketch: a vector of  $c \log^2(V^2) = 4c \log^2(V)$  1-sparse recovery data structures. Each recovery data structure is composed of three integers so a node sketch consists of  $12c \log^3(V)$  integers. As noted above, 128-bit(16B) integers are necessary when  $V \geq 10^5$ , so for  $c = 2$  the size of a node sketch is  $384 \log^3(V)$ B. Since there is a node sketch for each node in the graph, the entire streaming data structure has size  $384V \log^3(V)$ B. When  $V = 1$  million, this data structure is roughly 3 TB in size.

**Using existing  $\ell_0$ -samplers offers no advantage on modern hardware.** The goal of a streaming connected components algorithm is to use smaller space than would be required to store the entire graph explicitly. As we demonstrate empirically in Section 6, the most space-efficient dynamic graph processing system, Aspen, requires roughly 4B of space for each edge in the graph. A straightforward back-of-the-envelope calculation reveals that even for dense graphs with average degree  $V/2$ , STREAMINGCC would use less space than Aspen only on very large inputs which require enormous RAM capacities and thousands of years of processing time:  $384V \log^3(V)B \geq 4B \cdot V^2/4$  only when  $V \geq 4 \cdot 10^6$ . Processing a four million-node graph using STREAMINGCC would require more than 16 TB of RAM and, at an ingestion rate of less than 20 edges per second, would take more than six thousand years to process the graphs' roughly four trillion edges. While STREAMINGCC's space complexity is much smaller than explicit graph representations like Aspen's asymptotically, in absolute terms it offers no advantage on modern hardware.

### 3.1 Improved $\ell_0$ -Sampler for Graph Connectivity

We present CUBESKETCH, an  $\ell_0$ -sampling algorithm which is smaller than the best existing  $\ell_0$ -sampling algorithm and is asymptotically faster to update. It exploits the structure of the connected components problem to achieve these gains and as a result can only be used to sample from vectors of the form  $x = \{-1, 0, 1\}^n$ . This allows it to sample nonzero entries from the characteristic vectors of each node in the graph, a crucial component of streaming Boruvka emulation.

Since CUBESKETCH's goal is to recover a nonzero entry from a  $\{-1, 0, 1\}$  vector, its 1-sparse recovery procedure can be simpler and more efficient than a general 1-sparse recovery procedure, where the nonzero entry can take on any integer value. This 1-sparse recovery data structure maintains 2 values:  $\alpha$ , which is used to recover a single nonzero entry, and  $\gamma$ , which is used to verify that the value stored in  $\alpha$  is correct.  $\alpha$  and  $\gamma$  are each  $O(\log(n))$  bits, and therefore require  $O(1)$  machine words. Given a sequence of updates  $(e_1, \Delta_1), (e_2, \Delta_2), \dots, (e_k, \Delta_p)$  to the data structure,

$$\alpha = \bigoplus_{i \in [p]} b(e_i) \quad (1)$$

$$\gamma = \bigoplus_{i \in [p]} h(b(e_i)) \quad (2)$$

where  $\oplus$  denotes bitwise XOR,  $b(e_i)$  denotes the binary representation of  $e_i$  and  $h$  is a hash function drawn from a 2-wise independent family of hash functions. We denote the 1-sparse recovery data structure of  $I_{w,y}$  as  $r_{w,y}$  and its internal values as  $\alpha_{w,y}$  and  $\gamma_{w,y}$ . One may attempt to recover a nonzero entry from  $I_{w,y}$  as follows:

$$\text{result} = \begin{cases} e_i & \text{if } \alpha = b(e_i) \text{ and } \gamma = h(b(e_i)) \\ \text{NO ELEMENT} & \text{if } \alpha = 0 \text{ and } \gamma = 0 \\ \text{FAIL} & \text{if } \gamma \neq h_1(\alpha) \end{cases}$$

CUBESKETCH uses the selection process described in [18]: Given update  $(i, \Delta)$ , if  $h_y(i) \equiv 0 \pmod{2^w}$  then  $i$  is in  $I_{w,y}$ . For each

such  $I_{w,y}$ , its 1-sparse recovery data structure  $r_{w,y}$  is updated by setting  $\alpha_{w,y} = \alpha_{w,y} \oplus b(i)$  and  $\gamma_{w,y} = \gamma_{w,y} \oplus h(b(i))$ .

A nonzero entry is recovered from CUBESKETCH by attempting to recover a nonzero entry from each  $r_{w,y}$  until one returns a value other than NO ELEMENT or FAIL. If no such  $r_{w,y}$  exists, the algorithm returns NULL.

**THEOREM 1.** *CUBESKETCH is an  $\ell_0$  sampler that, for vector  $x \in \{-1, 0, 1\}^n$ , has space complexity  $O(\log^2(n))$ , worst-case update complexity  $O(\log^2(n))$ , and failure probability  $O(1/n^c)$  for constant  $c$ .*

**PROOF.** The space and update time results follow by construction: each 1-sparse recovery data structure  $r_{w,y}$  requires a constant number of machine words, and  $w = O(y) = O(\log(n))$ . Applying an update to any 1-sparse recovery data structure  $r_{w,y}$  requires constant time, and in the worst case, an update will be applied to each of the  $O(\log^2(n))$  such structures.

**Lemma 2.** *CUBESKETCH's selection process succeeds with high probability. Equivalently, CUBESKETCH contains a subset  $I_{w,y}$  with a single nonzero entry, that is,  $\Pr[\exists w, y \text{ s.t. } \text{supp}(I_{w,y}) = 1] = 1 - 1/n^c$ .*

**PROOF.** Adapted from [18]. Choose  $i \in [c \log(n)]$  such that  $2^{i-2} \leq \|x\|_0 < 2^{i-1}$  where  $\|x\|_0$  denotes the  $\ell_0$  norm of  $x$ , i.e., the number of nonzero entries of  $x$ . Then  $\forall j \in [c \log(n)]$ ,

$$\begin{aligned} \Pr[\text{supp}(I_{i,j}) = 1] &= \sum_{k \in I_{i,j}} \frac{1}{2} \left(1 - \frac{1}{2^i}\right)^{\|x\|_0 - 1} \\ &> \frac{\|x\|_0}{2} \left(1 - \frac{\|x\|_0}{2^i}\right) > 1/8. \end{aligned}$$

Then  $\Pr[\text{supp}(I_{i,j} \neq 1) \forall j \in [c \log(n)]] < 1 - 1/8^{-c \log(n)} = 1 - 1/n^{-3c}$ .  $\square$

**Lemma 3.** *CUBESKETCH's 1-sparse recovery error checking succeeds with high probability. That is,  $\forall w, y$ , if  $\text{supp}(I_{w,y}) > 1$  then  $\Pr[\gamma_{w,y} = h(\alpha_{w,y})] = 1 - 1/n^c$ .*

**PROOF.** When  $I_{w,y}$  has a single nonzero entry, it always passes the error check. That is, if  $\text{supp}(I_{w,y}) = 1$ ,  $\alpha_{w,y} = b(e_i)$  where  $e_i$  is the single nonzero element of  $I_{w,y}$ , and  $\gamma_{w,y} = h(b(e_i))$ .

When  $I_{w,y}$  has a single nonzero entry, then it passes the error check only in the rare event of a hash collision: If  $\text{supp}(I_{w,y}) > 1$ , fix  $e_i \in I_{w,y}$ .  $\gamma_{w,y} = h(\alpha_{w,y})$  iff  $\bigoplus_{j \in I_{w,y} \setminus e_i} h(b(j)) \oplus h(b(e_i)) = h(\alpha_{w,y})$ . Since  $h$  is a 2-wise independent hash function,

$$\Pr \left[ h(b(e_i)) = \bigoplus_{j \in I_{w,y} \setminus e_i} h(b(j)) \right] = \frac{1}{2^{c \log(n)}} = \frac{1}{n^c}.$$

$\square$

Lemmas 2 and 3 imply that CUBESKETCH is sampleable with high probability (see Definition 1). CUBESKETCH may be added via elementwise  $\oplus$  (exclusive or). Linearity of CUBESKETCH follows from the observation that exclusive or is a linear operation.  $\square$

Figure 2 illustrates that CUBESKETCH is far faster than the standard  $\ell_0$ -sampling algorithm. In fact, when applied to graphs with at least  $10^5$  nodes, it is more than 3 orders of magnitude faster. This dramatic speedup is a result both of CUBESKETCH's asymptotically

lower update time complexity, and the fact that its update cost is dominated by bitwise exclusive OR operations, which are in practice much faster than the division operations standard  $\ell_0$ -sampling performs. Finally, standard  $\ell_0$  sampling is slowed significantly by the need to perform  $O(\log^2(V))$  modular exponentiation operations on 128-bit integers for each update when  $V \geq 10^5$ . CUBESKETCH does not require 128-bit operations until processing graphs with tens of billions of nodes.

Figure 3 shows that, for the same input vector length and failure probability, CUBESKETCH is four times smaller than standard  $\ell_0$  sampling for smaller vectors, and eight times smaller for larger vectors. This is a result of the fact that CUBESKETCH's 1-sparse recovery data structures use four times fewer machine words than those of standard  $\ell_0$  sampling, and the fact that CUBESKETCH does not need to use 128-bit integers for longer vectors.

#### 4 BUFFERING FOR I/O EFFICIENCY AND IMPROVED PARALLELISM

In the streaming connectivity problem, stream updates are *fine-grained*: each update represents the insertion or deletion of a single edge. Since streams are ordered arbitrarily, even a short sequence of stream updates can be highly non-local, inducing changes throughout the graph. As a result, STREAMINGCC and similar graph streaming algorithms do not have good data locality in the worst case. This lack of locality can cause many CPU cache misses and therefore reduce the ingestion rate, even when sketches are stored in RAM. The cache-miss cost can be high since ingesting each stream update  $(u, v, \Delta)$  requires modifying a logarithmic number of sketches, and can thus induce a poly-logarithmic number of cache misses. The consequences are even worse if sketches are stored on disk since each edge update requires loading a logarithmic number of sketches from disk, leading to the following observation.

**Observation 1.** *In the hybrid semi-streaming model with  $M = o(V \log^3(V))$  RAM and  $D = \Omega(V \log^3(V))$  disk, STREAMINGCC uses  $\Omega(1)$  I/Os per update and processing the entire stream of length  $N$  uses  $\Omega(N) = \Omega(E)$  I/Os.*

Any sketching algorithm that scales out of core suffers severe performance degradation unless it amortizes the per-update overhead of accessing disk. Such an amortization is not straightforward, since sketching inherently makes use of hashing and as a result induces many random accesses, which are slow on persistent storage. We now introduce a sketching algorithm for the streaming connected components problem that amortizes disk access costs, even on adversarial graph streams, and as a result is simultaneously a space-efficient graph semi-streaming algorithm and an I/O-efficient external-memory algorithm. We also note that the design facilitates parallelism, which we experimentally verify in Section 6.

##### 4.1 I/O-Efficient Stream Ingestion

We describe GRAPHZEPPELIN's I/O efficient stream ingestion procedure in the hybrid streaming model.

Arbitrarily partition the nodes of the graph into **node groups** of cardinality  $\max\{1, B/\log^3(V)\}$ . Let  $\mathcal{U} \subset \mathcal{V}$  denote a node group, and let  $\mathcal{S}(\mathcal{U})$  denote the node sketches associated with the nodes in  $\mathcal{U}$ . Store  $\mathcal{S}(\mathcal{U})$  contiguously on disk. This allows  $\mathcal{S}(\mathcal{U})$  to be

read into memory I/O efficiently: if node groups are of cardinality 1, then  $B$  is smaller than the size of a node sketch, and if each node group has cardinality  $B/\log^3(V) > 1$ , then the sketches for the group have total size  $O(B)$ .

Applying stream update  $((u, v), \Delta)$  to node sketches of  $u$  and  $v$  immediately upon arrival takes  $\Omega(1)$  I/Os since the corresponding sketches must be read from disk. To amortize the cost of fetching sketches, GRAPHZEPPELIN only fetches  $\mathcal{S}(\mathcal{U}_i)$  when it has collected  $\max\{B, \log^3(V)\}$  updates for  $\mathcal{U}_i$ . Since there may be  $O(V)$  node groups, collecting these updates for each node group cannot be done in RAM. Instead, we collect these updates I/O efficiently on disk using a **gutter tree**, a simplified version of a buffer tree [8] which uses  $O(V \log^3(V))$  space.

Like a buffer tree, a gutter tree consists of a tree whose vertices each have buffers of size  $O(M)$ . Each non-leaf vertex has  $O(M/B)$  children. We refer to a leaf vertex of the gutter tree as a **gutter**, because it fills with stream data but is periodically emptied by applying the contained stream data to sketches. Each leaf vertex in the gutter tree is associated with a node group  $\mathcal{U}$  and has size  $\max\{B, \log^3(V)\}$ , the same size as  $\mathcal{S}(\mathcal{U})$ . When a gutter for node group  $\mathcal{U}$  fills, GRAPHZEPPELIN reads  $\mathcal{S}(\mathcal{U})$  and the updates stored in the gutter into memory, applies the updates to  $\mathcal{S}(\mathcal{U})$ , and writes  $\mathcal{S}(\mathcal{U})$  back to disk. Since data does not persist in leaf vertices, no rebalancing is necessary.

**Lemma 4.** *GRAPHZEPPELIN's stream ingestion uses  $O(V \log^3(V))$  space and  $\text{sort}(N) = O(N/B \log_{M/B}(V/B))$  I/Os in the hybrid streaming setting.*

**PROOF.** GRAPHZEPPELIN's sketch data structures use  $O(V \log^3(V))$  space.

Each leaf in the gutter tree has a gutter of size  $\max\{B, \log^3(V)\}$ . This is one gutter for each node group and there are  $V/(\max\{1, B/\log^3(V)\})$  node groups so the total space for the leaves of the gutter tree is  $O(V \log^3(V))$ .

In the level above the leaves, there are  $V \log^3(V)/B \cdot B/M$  vertices each with size  $M$ , so the total space used at this level is  $O(V \log^3(V))$ . Each subsequent higher level of the tree uses  $O(M/B)$  space less than the level below it, so the total space used for the entire gutter tree is  $O(V \log^3(V))$ .

The I/O complexity of the gutter tree is equivalent to that of the buffer tree, except that leaf gutters are flushed by reading in the appropriate sketches from disk and applying the updates in the gutter to these sketches. Asymptotically this incurs no additional cost so the total I/O complexity for ingestion is  $\text{sort}(N)$ .  $\square$

##### 4.2 I/O-Efficient Connectivity Computation

**Lemma 5.** *Once all stream updates have been processed, GRAPHZEPPELIN computes connected components using  $O((V \log^3(V)/B) + (V \log^*(V)))$  I/Os in the hybrid streaming model.*

**PROOF.** Each round of Boruvka's algorithm has three phases. In the first, an edge is recovered from the sketch of each current connected component. In the second, for each edge its endpoints are merged in a disjoint set union data structure which keeps track of the current connected components. In the third phase, for each pair of connected components merged in phase 2, the corresponding



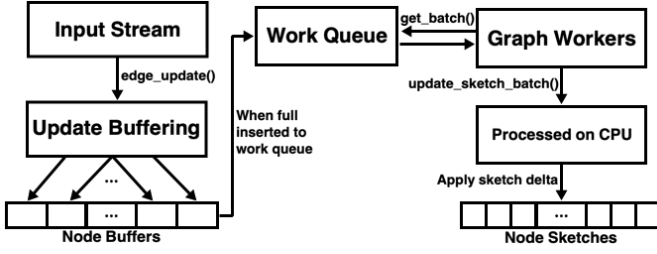


Figure 4: GRAPHZEPPELIN stream ingestion data flow during ingestion.

sketches are summed together. We analyze the I/O cost of each phase of a round separately.

In the first round, to query the sketches in the first phase, all of the sketches must be read into RAM which can be done with a single scan. This uses  $O(V \log^3(V)/B)$  I/Os.

The disjoint set union data structure has size  $(V)$  and must be stored on disk. In the second phase the cost of each DSU merge is  $\log^*(V)$  I/Os, so the total I/O cost is  $V \log^*(V)$ .

In the third phase, summing the sketches of the merged components together is I/O efficient if  $B = O(\log^3(V))$ , since the disk reads and writes necessary for summing sketches are the size of a block or larger. The cost for the third phase is  $(V \log^3(V)/B)$ .

If  $B = \omega(\log^3(V))$ , sketches are much smaller than the block size. Since the merges performed in each round of Boruvka are a function both of the input stream and of the randomness of the sketches, these merges induce random accesses to the sketches on disk and so summing the sketches for each merge takes  $O(1)$  I/Os. In total, the third phase takes  $O(V)$  I/Os in this case.  $\square$

**Corollary 1.** When  $B = o(\log^3(V))$  or  $M = O(V)$ , GRAPHZEPPELIN is an I/O optimal algorithm for the connected components problem; i.e., it uses  $\text{sort}(E) = O(E/B(\log_{M/B}(V/B)))$  I/Os.

In practice, for some graph streams  $M = O(VB)$  and  $D = O(V \log^3(V))$ . In this case, it is possible to do away with the upper levels of the gutter tree and write I/O efficiently to the gutters stored on disk. In Section 5 we describe how GRAPHZEPPELIN can perform stream ingestion using either a full gutter tree or just the leaf gutters, and evaluate the performance of both approaches in out-of-core and parallel settings in Section 6.

## 5 SYSTEM DESCRIPTION

The GRAPHZEPPELIN algorithm is split into two phases: a **stream ingestion** phase, in which edge updates are processed and stored using CUBESKETCH, and a **query-processing** phase, in which a spanning forest for the graph is recovered from these sketches. These components have been adapted to use SSD when the sketches are so large that they do not fit in RAM. They have also been adapted to run efficiently on a multi-core system.

GRAPHZEPPELIN’s user-facing API consists of `INITIALIZE()`, `EDGE_UPDATE()` for processing stream updates, and `LIST_SPANNING_FOREST()` to compute and return the connected components. When `INITIALIZE()` is called, GRAPHZEPPELIN allocates  $\log(V)$  CUBESKETCH data structures for each node in the graph. In

total GRAPHZEPPELIN allocates  $24V \cdot \log^3(V)$  bytes for sketches over all nodes. It additionally initializes a buffering data structure.

### 5.1 Stream Ingestion

Stream updates are immediately placed in a buffering system. Periodically, the buffering system produces a batch of updates that are all for the same graph node  $u$ . This batch is inserted into a work queue, which then hands the batch off to a **Graph Worker**, i.e., a thread for carrying out the sketch updates. GRAPHZEPPELIN runs many Graph Workers simultaneously to increase the stream ingestion rate. GRAPHZEPPELIN is designed so that stream ingestion is “embarrassingly parallel”. This is because each batch is only applied to a single node sketch, and also because updating each of the  $\log(V)$  CUBESKETCHES in a node sketch can be done in parallel. A high-level illustration of GRAPHZEPPELIN stream ingestion is shown in Figure 4.

Upon each call to `EDGE_UPDATE((e, v), Δ)`, GRAPHZEPPELIN inserts the updates  $(u, Δ)$  and  $(v, Δ)$  into the buffer data structure. This section describes the implementations of the next steps during GRAPHZEPPELIN’s stream ingestion phase.

**Buffering.** GRAPHZEPPELIN’s buffering process ingests updates from the stream and periodically outputs a batch of updates for a single node in the graph. GRAPHZEPPELIN implements two buffering data structures: a gutter tree, described in Section 4, and a simplified version of the gutter tree, which only includes the leaves. Depending upon available memory, GRAPHZEPPELIN uses only one of these two buffering structures. The gutter tree is designed to buffer updates I/O efficiently on SSD. The leaf-only version is fundamentally a special case gutter tree in which  $M > V \cdot B$  and is optimized for this case.

These buffering techniques confer several benefits. First, when GRAPHZEPPELIN’s sketches are so large that they do not fit in RAM and are stored on SSD, applying updates to a single node sketch in large batches amortizes the I/O cost of reading the node sketch into memory. Without buffering, each stream update would incur  $\Omega(1)$  I/Os in the worst case. Additionally, as we explain in Section 6.4, buffering facilitates parallelism.

**Gutter tree.** GRAPHZEPPELIN allocates 8MB for each non-leaf buffer in the gutter tree. The gutter tree writes updates to the disk in blocks of 16KB, and has a fan-out of  $\frac{8\text{MB}}{16\text{KB}} = 512$ . When  $V > 50000$ , the size of a sketch is greater than 100KB, much larger than this 16KB block. Therefore, the leaf nodes of the gutter tree accumulate updates for a single graph node. GRAPHZEPPELIN allocates space for each leaf gutter equal to the size of a node sketch.

When we initialize GRAPHZEPPELIN, we leverage the static structure of the gutter tree to pre-allocate its disk space. During stream ingestion, each edge update  $(u, v)$  is split into two updates  $(u, v)$  and  $(v, u)$  and inserted to the root buffer of the gutter tree. Another thread asynchronously flushes the contents of full buffers to the appropriate child using the `pwrite` system call. When a flush causes the buffer of a child node to fill, that child node is recursively flushed before the flush of the parent continues. When a leaf gutter is full we move the batch of updates to the work queue for processing by the Graph Workers.



*Leaf-only gutter tree.* For each graph node  $u$  we maintain a gutter that accumulates updates for  $u$ . When the system is initialized, we allocate the memory for each of these gutters. Each leaf gutter is  $1/4$  the size of a node sketch.

Upon call to `EDGE_UPDATE( $(e, v), \Delta$ )`, the edge  $e = (u, v)$  is directly inserted into the gutters for nodes  $u$  and  $v$ . As before, when the gutter for either node becomes full, it is flushed and the batch is inserted in the work queue.

It’s important to note that the leaf-only gutter data structure need not fit in RAM, because the GRAPHZEPPELIN design permits efficient swapping to SSD; see Section 6.

**Work queue.** The work queue functions as a simple solution for the producer-consumer problem, in which the thread filling buffers produces work and the Graph Workers consume it. Once a buffer is filled the `EDGE_UPDATE()` function inserts the batch of updates into the work queue. Once the batch reaches the front of the queue, it is removed from the queue for processing by a Graph Worker. Insertions to the queue are blocked while the queue is full, and Graph Workers in need of work are blocked while the queue is empty. The work queue can hold up to  $8g$  batches, where  $g$  is the number of Graph Workers.

**Sketch updates.** When a Graph Worker receives a batch of updates for node  $u$  from the work queue, it applies each update in the batch to each of the  $\log(V)$  CUBESKETCHES in the node sketch of  $u$ .

As described in Section 3, a CUBESKETCH is a vector of 1-sparse recovery data structures. Each 1-sparse recovery structure consists of a 64 bit  $\alpha$  value and a 32 bit  $\gamma$  value. Each CUBESKETCH stores a two dimensional array  $A$  of sparse recovery structures  $r_{w,y}$ , with dimensions  $\log(V^2) \times 1/2 \log(V^2)$ . To apply an update  $(i, \Delta)$  to a CUBESKETCH, the Graph Worker determines which sets  $I_{w,y}$  contain  $i$ , and modifies their sparse recovery structures’  $\alpha$  and  $\gamma$  values by setting  $\alpha := \alpha \oplus i$  and  $\gamma := \gamma \oplus h_y(i)$ . The hash values are calculated using xxHash [17].

Each CUBESKETCH data structure uses  $1/2 \log^2(V^2) = 2 \log^2(V)$  sparse recovery structures, each of which takes 12 bytes of space. In total, this is  $24 \log^2(V)$  bytes per CUBESKETCH, and  $24 \log^3(V)$  bytes per node sketch.

**Multithreading sketch updates.** Applying a batch to a node sketch is handled asynchronously by a Graph Worker, allowing what we call *batch-level parallelism*. We implement these workers using C++ STL threads that wait on the work queue.

We use OpenMP [55] to dispatch a group of threads to process the each CUBESKETCH update. We refer to this as *sketch-level parallelism*. OpenMP allows us to specify the number of threads to allocate to a task and handles work allocation transparently. When updating a node sketch, applying a batch to each CUBESKETCH is treated as one work unit and OpenMP allocates the  $\log(V)$  units between the apportioned threads.

Implementing both batch- and sketch-level parallelism gives us a natural way to tune GRAPHZEPPELIN’s performance. For instance, we can decide to configure more Graph Workers with fewer threads per group, or fewer Graph Workers with more threads per group. We experimentally determine a good configuration for our hardware and datasets (see Section 6.4).

A single work unit is never shared between threads in the same group. As a result, a CUBESKETCH is only modified by one thread in a group, so no locking is necessary at the sketch level. However, locking is necessary at the batch level because consecutive batch updates may be requested to the same node sketch, and thus multiple graph workers may seek to dispatch thread groups to the same sub-sketches. We minimize the size of this critical section by exploiting linearity of  $\ell_0$ -samplers. Rather than locking a node sketch  $S(x)$  for the entire batch operation, we apply the updates to an empty sketch  $S(x_0)$  and lock only to add  $S(x) = S(x) + S(x_0)$ .

## 5.2 Query Processing

When a connectivity query is issued, GRAPHZEPPELIN calls `LIST_SPANNING_FOREST()` which returns a spanning forest of the graph. The first step of post-processing is to flush the buffering data structure of any remaining updates, moving the batches to the work queue. We then wait for the Graph Workers to finish processing these batches. Finally, GRAPHZEPPELIN runs Boruvka’s algorithm to generate a spanning forest of the input graph.

## 6 EVALUATION

**Experimental setup.** We implemented GRAPHZEPPELIN as a C++11 executable compiled with g++ version 9.3 for Ubuntu. All experiments were run on a Dell Precision 7820 with 24-core 2-way hyperthreaded Intel(R) Xeon(R) Gold 5220R CPU @ 2.20GHz, 64GB 4x16GB DDR4 2933MHz RDIMM ECC Memory and two 1 TB Samsung 870 EVO SSDs. In some of our experiments we artificially limited RAM to force systems to page to disk using Linux Control Groups. We put a swap partition on one of the two SSDs, and the other SSD held the datasets.

### 6.1 Datasets

We used two types of data sets in this paper. First, we generated large, dense graphs using a Graph500 specification, and converted these to streams for our evaluation. We also evaluated correctness on graphs from the SNAP graph repository [41] and the Network Repository [60]. All data sets used are described in Table 5.

*Synthesizing Dense Graphs and Streams* We created undirected graphs using the Graph500 Kronecker generator. We produced five simple, undirected graphs. These graphs are dense: each has roughly one half of all possible edges. The Graph500 generator does not output simple graphs by default, so to produce our five simple graphs we pruned duplicate edges and self-loops [7].

We subsequently translated each of the 5 graphs into a random stream of edge insertions and deletions with the following guarantees: (i) an insertion always occurs before a deletion, (ii) an edge never receives two consecutive updates of the same type, (iii) we disconnect a small (fewer than 150) set of nodes from the rest of the graph, and (iv) by the end of the stream, exactly the input graph (with the exception of the edges removed to disconnect the vertices in (iii)) remains. Note that this mechanism deliberately adds edges not in the original graph, but they are always subsequently deleted. We implemented (iii) to guarantee some non-trivial connected components in each stream’s final graph.

Name	# of Nodes	# of Edges	# Stream Updates
<b>kron13</b>	$2^{13}$	$1.7 \times 10^7$	$1.8 \times 10^7$
<b>kron15</b>	$2^{15}$	$2.7 \times 10^8$	$2.8 \times 10^8$
<b>kron16</b>	$2^{16}$	$1.1 \times 10^9$	$1.1 \times 10^9$
<b>kron17</b>	$2^{17}$	$4.3 \times 10^9$	$4.5 \times 10^9$
<b>kron18</b>	$2^{18}$	$1.7 \times 10^{10}$	$1.8 \times 10^{10}$
<b>p2p-gnutella</b>	$6.3 \times 10^4$	$1.5 \times 10^5$	$2.9 \times 10^5$
<b>rec-amazon</b>	$9.2 \times 10^4$	$1.3 \times 10^5$	$2.5 \times 10^5$
<b>google-plus</b>	$1.1 \times 10^5$	$1.4 \times 10^7$	$2.7 \times 10^7$

**Figure 5: The dimensions of the data sets used. The first 5 rows detail the graphs and streams synthesized as outlined in Section 6.1. Each row lists the number of nodes for a given graph, the number of edges in the graph and the number of edge updates in the stream generated from that graph.**

*Publicly Available Datasets* We also used the following real-world data sets. **p2p-gnutella** is a graph representing the Gnutella peer-to-peer network [59]. **rec-amazon** is a co-purchase recommendation graph for products listed on Amazon [40], where each node represents a product and there is an edge between two nodes if their corresponding products are frequently purchased together. **google-plus** is a graph among users of the Google Plus social network [48] where edges represent follower relations. Each of these real-world graphs was converted to a stream using the process described above.

## 6.2 Speed Test & Memory Usage Experiment

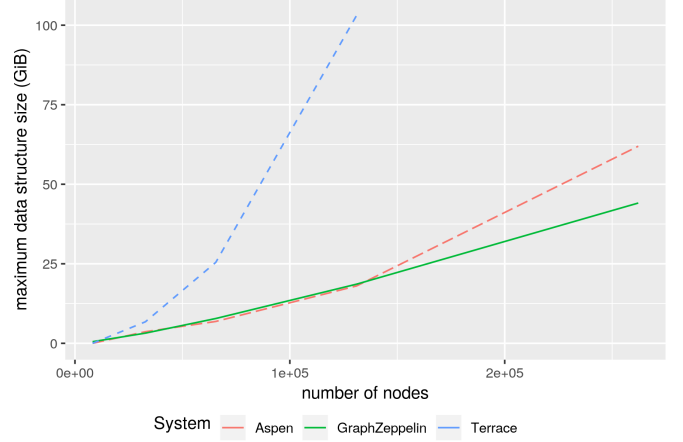
We now demonstrate that, given the same memory resources, GRAPHZEPPELIN can handle larger inputs than Aspen and Terrace on sufficiently large and dense graph streams. We also show that unlike these systems, GRAPHZEPPELIN does not suffer significant performance degradation once the scale of the streamed graphs exceeds available memory.

Both Aspen and Terrace are optimized for the *batch-parallel* model of dynamic graph processing. In this model, one begins with a non-empty graph and then applies updates to the graph in batches containing exclusively insertions or exclusively deletions. This contrasts with our streaming model, where we may only begin with an empty edge-set and construct the edges of the graph entirely from a stream of interspersed inserts and deletes. To avoid unfairly penalizing Aspen and Terrace, during their experiments we accumulate edge insertions into an array of insertions and edge deletions into an array of deletions. Whenever one of these arrays fills, we feed it into the appropriate batch update function provided by Aspen or Terrace. Note that Terrace does not currently support batch deletions, so we rely on its individual edge deletion functionality instead and do not maintain a deletions array. This naive batching approach does not guarantee correctness exactly the same set of connected components, but the point of this experiment is to make batched inputs so that we can fairly measure the performance of Aspen and Terrace.

We ran GRAPHZEPPELIN, Aspen, and Terrace on each Kronecker stream. We used a batch size of  $10^6$  for Aspen and Terrace because

Dataset	Aspen	Terrace	GRAPHZEPPELIN
kron13	0.000328	0.000519	0.52
kron15	3.40	6.30	3.20
kron16	6.40	23.7	7.70
kron17	16.8	96.0*	18.6
kron18	57.7	N/A	44.1

**(a) All number listed in GiB. The asterisk indicates that Terrace did not finish by the 24 hour termination mark, and so the reported number is memory consumption at the time of termination (counting swap space used in excess of the 64 GB of RAM available on our machine).**



**(b) GRAPHZEPPELIN is asymptotically more memory efficient than either Aspen or Terrace on large, dense graphs.**

**Figure 6: Memory consumption in the unrestricted RAM setting of Aspen, Terrace and the leaf-only version of GRAPHZEPPELIN.**

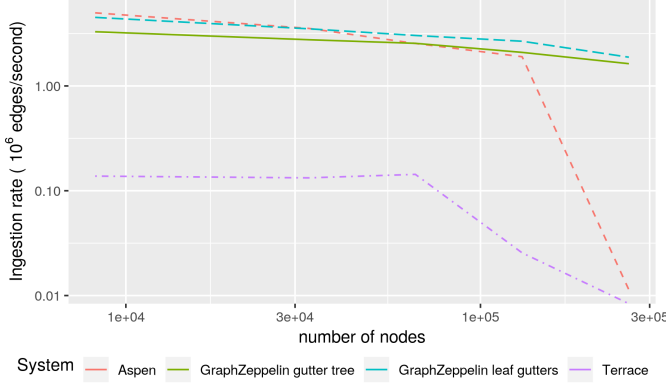
we found this to produce the highest ingestion rates for both systems. To record memory usage we logged the output of the Linux `top` command tracking each system every five seconds. To avoid penalizing Aspen and Terrace for input stream file read time, we dedicated a thread to buffering the input stream on all systems evaluated so that batches can be received at RAM latency. If any system did not complete a particular stream after 24 hours, then it was terminated.

*Memory profiling.* GRAPHZEPPELIN's space-efficient CUBESKETCHES give it an asymptotic  $O(V/\log^3(V))$  space advantage over Aspen and Terrace. Given the polylogarithmic factors and constants, this experiment determines the actual crossover point where GRAPHZEPPELIN processes graphs more compactly than Aspen and Terrace. As shown in Table 6a, GRAPHZEPPELIN is more compact than Terrace even on kron15, and is more compact than Aspen on kron18.

*I/O Performance and Ingestion Rate.* Unlike Aspen and Terrace, GRAPHZEPPELIN can leverage external memory effectively to maintain consistently high ingestion rates when its data structures are too large to fit in RAM. In Figure 7b we summarize the results of running Aspen, Terrace, and GRAPHZEPPELIN with only 16GB of RAM. The ingestion rates of both Aspen and Terrace plummet

Dataset	Aspen	Terrace	Gutter Tree	Leaf-Only Gutters
kron13	4.98	0.138	3.29	4.51
kron15	3.54	0.133	2.76	3.51
kron16	2.54	0.0143	2.55	3.03
kron17	1.90	0.0404*	2.09	2.67
kron18	0.0759*	N/A	1.63	1.88

(a) Ingestion rates in millions of updates per second are listed for Aspen, Terrace and both the gutter tree and leaf-only gutter GRAPHZEPPELIN implementations. Numbers with an asterisk indicate that the respective system did not finish.



(b) Plot of ingestion rates for Aspen, Terrace and GRAPHZEPPELIN.

Figure 7: GRAPHZEPPELIN eventually outperforms Aspen and Terrace in restricted RAM setting. Update throughput evaluated for all systems when restricted to 16 GiB of RAM.

once they exceed available RAM. Neither Aspen nor Terrace were able to finish their largest evaluated stream within 24 hours ( $2^{17}$  for Terrace and  $2^{18}$  for Aspen), and so the corresponding numbers in Table 7a are estimates based on the portion of these streams the systems were able to process within this time. In comparison, GRAPHZEPPELIN’s ingestion rate remains high when its memory consumption extends into secondary storage. GRAPHZEPPELIN finished the kron18 stream with an average ingestion rate of 1.88 million updates per second, lagging only 0.55 million updates per second behind its performance on the same stream given unrestricted access to RAM.

We also evaluated Aspen, Terrace and GRAPHZEPPELIN without the 16 GB memory restriction. For sufficiently large, dense graphs GRAPHZEPPELIN can process the same inputs faster than either Aspen or Terrace entirely within RAM (see Figure 8). GRAPHZEPPELIN exhibited far higher ingestion rates than Terrace on every Kroecker stream we evaluated. Aspen proved faster than GRAPHZEPPELIN on inputs up to and including the kron15 stream, but the GRAPHZEPPELIN ingestion rate overtook that of Aspen on kron16 and larger graphs.

### 6.3 Correctness Experiment

Our algorithm has a small but nonzero failure probability. Here we present a correctness experiment, where we give stream updates to

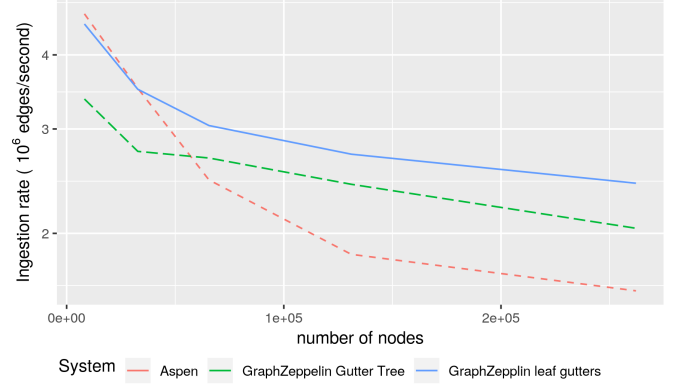


Figure 8: Ingestion rates for Aspen and GRAPHZEPPELIN without any RAM restriction. It is already evident from our results for ingestion rate with a 16 GB RAM restriction that GRAPHZEPPELIN is faster than Terrace on inputs where the data structures fit entirely within RAM, so Terrace data is omitted here.

Dataset	Aspen	Terrace	Gutter Tree	Leaf-Only Gutters
kron13	0.041	0.126	0.11	0.100
kron15	0.202	0.800	0.53	0.520
kron16	0.746	1.260	1.22	1.21
kron17	3.11	N/A	31.7	48.8
kron18	N/A	N/A	537	553

Figure 9: Time to compute connected components after stream ingestion for Aspen, Terrace and both the gutter tree and leaf-only GRAPHZEPPELIN implementations. Cells labeled N/A indicate connected components was never queried due to ingestion time exceeding the 24-hour termination mark.

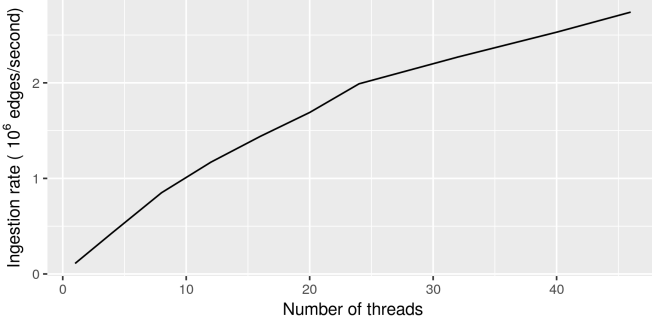
both our sketching algorithm and to an in-memory adjacency matrix stored as a bit vector. To perform a single correctness check, we pause the stream updates, and compare the connected components returned by our algorithm with the output of running Kruskal’s algorithm on the adjacency matrix. We ran 1000 correctness checks using the kron17 input stream, and none ever failed.

While our algorithm’s performance is optimized for dense graphs, it still succeeds with high probability for sparser graphs. To ensure this theoretical guarantee pans out in practice, we ran 1000 correctness checks on each real graph in Table 5. None of these correctness checks failed.

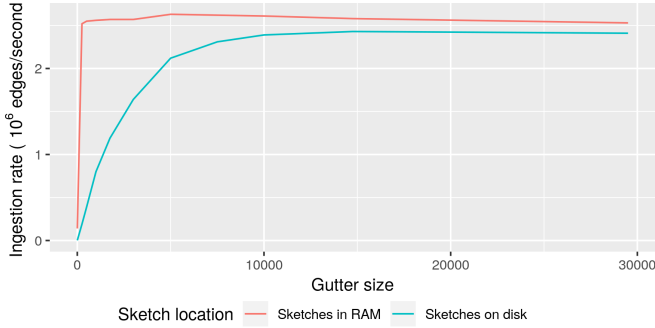
### 6.4 GRAPHZEPPELIN Multithreading Experiment

Due to the atomized nature of sketch updates, we expect stream ingestion to scale well on multi-core systems. In this section, we experimentally evaluate this claim by varying the number of threads used for processing updates and observe a significant speed-up.

Figure 10 shows the ingestion rate of GRAPHZEPPELIN as the number of threads processing the kron17 graph stream increases.



**Figure 10: GRAPHZEPPELIN processes sketch updates in parallel to increase ingestion rate. The insertion rate increases as threads are added resulting in an overall increase of 25 $\times$ .**



**Figure 11: GRAPHZEPPELIN gutter size vs ingestion speed. When node sketches are stored in RAM a small buffer size is sufficient for good performance. However, when sketches are on disk a much higher buffer size is required.**

The threads use a memory of 64GB so that the parallel performance can be measured without memory contention. To avoid external memory accesses, we use the leaf-only gutters for buffering. The per-thread increase in ingestion rate is significant; the ingestion rate for 46 threads is approximately 25 times higher than that of a single thread. Additionally, at 46 threads the marginal ingestion rate is still positive, suggesting that adding more threads would further increase performance.

We also experimentally determined that a group size of one gives the best performance with our combination of machine and inputs.

### 6.5 GRAPHZEPPELIN Buffer Size Experiment

Applying sketch updates is highly scaleable, only if updates are buffered and applied in batches. When node sketches reside partially on disk, directly applying each update to node sketches incurs  $\Omega(1)$  IOs per update. Additionally, cache contention and thread synchronization bottleneck the ingestion rate even when data is entirely within RAM. For these reasons we retain buffers of a constant factor  $c \leq 1$  of the node sketch size.

Figure 11 demonstrates that this bottleneck is significant both when node sketches are entirely in RAM and especially when more than half the node sketches are on disk. Both systems ingest the

kron17 graph stream using 46 Graph Workers and a group size of 1. When directly placing updates into the work queue, the performance of the in RAM system, 138 thousand updates per second, approximates the performance of a single Graph Worker that buffers updates. If node sketches are on disk, the ingestion rate is only two thousand insertions per second a further reduction of 65 $\times$ .

When the sketches and buffers fit within memory, performance increases rapidly indicating that  $c$  can be quite small while providing a high ingestion rate. However, once memory requirements exceed main memory,  $c$  must be larger to offset disk IOs. To achieve an ingestion rate within 5% of peak performance on kron17,  $c = 0.008$  was sufficient for entirely in RAM computation, while  $c = .25$  was required when node sketches partially reside on disk.

## 7 RELATED WORK

**Graph Streaming Systems.** Existing graph stream processing systems are designed primarily to handle updates in batches consisting entirely of insertions or entirely of deletions. Streaming systems that process updates in batches are generally divided into two categories. The first consists of those systems which finish ingestion prior to beginning queries and finish queries prior to accepting any additional edges [6, 13, 20, 50, 56, 63, 64]. The second allows updates to be applied asynchronously by periodically taking “snapshots” of the graph during ingestion to be used in conducting queries [15, 19, 31, 32, 47]. Terrace falls within the former category, while Aspen falls within the latter.

The batching employed in these systems serves to amortize the cost of applying updates, but also limits the granularity at which insertions and deletions may be interspersed during ingestion. In contrast, the model of GRAPHZEPPELIN allows for insertions and deletions to be arbitrarily interspersed during ingestion without sacrificing query correctness.

**External Memory Systems.** There is a rich literature of graph processing systems that exploit external memory to process static graphs. Such systems can be broadly partitioned into those designed to store the entire graph out-of-core [27, 38, 46, 74, 76], and semi-external memory systems that maintain only the vertex-set in RAM [3, 45, 61, 73, 75]. Some systems provide (at least theoretical) design extensions to handle queries on graphs with insert-only updates [14, 38, 68, 69, 74], but to the best of our knowledge GRAPHZEPPELIN is the first to leverage external-memory effectively in the streaming model of insertions *and* deletions.

## 8 CONCLUSION

GRAPHZEPPELIN computes the connected components of graph streams using space asymptotically smaller than an explicit representation of the graph. It is based on CUBESKETCH, a new  $\ell_0$ -sketching data structure that outperforms the state of the art on graph-streaming workloads. This new sketching technique allows GRAPHZEPPELIN to process larger, denser graphs than existing graph-streaming systems given a fixed RAM budget and to ingest these graph streams more quickly. GRAPHZEPPELIN employs work-buffering strategies, which allow sketches to be stored and on SSD. Even when GRAPHZEPPELIN’s sketch data structures are too large to fit in RAM, it can process graph streams at the cost of

a only small decrease in ingestion rate. Thus, GRAPHZEPPELIN is simultaneously a space-optimal graph semi-streaming algorithm and an I/O-efficient external-memory algorithm.

The asymptotically small space complexity of GRAPHZEPPELIN's linear sketching algorithm allows it to optimize for large, dense graph streams, unlike prior graph-processing systems, which often focus on sparse graphs. Thus, GRAPHZEPPELIN demonstrates that computational questions on graphs once thought intractably large and dense are now within reach.

GRAPHZEPPELIN illustrates that additional algorithmic improvements help make graph semi-streaming algorithms into a powerful engineering tool by reducing the update-time complexity and allowing sketches to be stored efficiently on SSD. These techniques may generalize to other graph-analytics problems.

## REFERENCES

- [1] J. Abello, A. L. Buchsbaum, and J. R. Westbrook. A functional approach to external graph algorithms. *Algorithmica*, 32(3):437–458, 2002.
- [2] K. J. Ahn, S. Guha, and A. McGregor. Analyzing graph structure via linear measurements. In *Proceedings of the twenty-third annual ACM-SIAM symposium on Discrete Algorithms*, pages 459–467. SIAM, 2012.
- [3] Z. Ai, M. Zhang, Y. Wu, X. Qian, K. Chen, and W. Zheng. Clip: A disk i/o focused parallel out-of-core graph processing system. *IEEE Transactions on Parallel and Distributed Systems*, 30(1):45–62, 2018.
- [4] R. Albert. Scale-free networks in cell biology. *Journal of cell science*, 118(21):4947–4957, 2005.
- [5] S. Allegretti, F. Bolelli, M. Cancilla, and C. Grana. Optimizing gpu-based connected components labeling algorithms. In *2018 IEEE International Conference on Image Processing, Applications and Systems (IPAS)*, pages 175–180, 2018.
- [6] K. Ammar, F. McSherry, S. Salihoglu, and M. Joglekar. Distributed evaluation of subgraph queries using worstcase optimal lowmemory dataflows. *arXiv preprint arXiv:1802.03760*, 2018.
- [7] J. Ang, B. W. Barrett, K. B. Wheeler, and R. C. Murphy. Introducing the graph 500. 2010.
- [8] L. Arge. The buffer tree: A new technique for optimal i/o-algorithms (extended abstract). In *University of Aarhus*, pages 334–345. Springer-Verlag, 1995.
- [9] T. Y. Berger-Wolf and J. Saia. A framework for analysis of dynamic social networks. In *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '06*, pages 523–528, New York, NY, USA, 2006. Association for Computing Machinery.
- [10] I. Bordino and D. Donato. Dynamic characterization of a large web graph.
- [11] G. S. Brodal, R. Fagerberg, D. Hammer, U. Meyer, M. Penschuck, and H. Tran. An experimental study of external memory algorithms for connected components. In *19th International Symposium on Experimental Algorithms (SEA 2021)*. Schloss Dagstuhl-Leibniz-Zentrum für Informatik, 2021.
- [12] L. Buš and P. Tvrđik. A parallel algorithm for connected components on distributed memory machines. In *European Parallel Virtual Machine/Message Passing Interface Users' Group Meeting*, pages 280–287. Springer, 2001.
- [13] F. Busato, O. Green, N. Bombieri, and D. A. Bader. Hornet: An efficient data structure for dynamic sparse graphs and matrices on gpus. In *2018 IEEE High Performance Extreme Computing Conference (HPEC)*, pages 1–7. IEEE, 2018.
- [14] J. Cheng, Q. Liu, Z. Li, W. Fan, J. C. Lui, and C. He. Venus: Vertex-centric streamlined graph computation on a single pc. In *2015 IEEE 31st International Conference on Data Engineering*, pages 1131–1142. IEEE, 2015.
- [15] R. Cheng, J. Hong, A. Kyrola, Y. Miao, X. Weng, M. Wu, F. Yang, L. Zhou, F. Zhao, and E. Chen. Kineograph: taking the pulse of a fast-changing and connected world. In *Proceedings of the 7th ACM european conference on Computer Systems*, pages 85–98, 2012.
- [16] Y.-J. Chiang, M. T. Goodrich, E. F. Grove, R. Tamassia, D. E. Vengroff, and J. S. Vitter. External-memory graph algorithms. In *Proc. 6th annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 139–149. Society for Industrial and Applied Mathematics, 1995.
- [17] Y. Collet. xxhash-extremely fast non-cryptographic hash algorithm. URL <https://github.com/Cyan4973/xxHash>, 2016.
- [18] G. Cormode and D. Firmani. A unifying framework for  $\ell_0$ -sampling algorithms. *Distributed and Parallel Databases*, 32, 09 2014.
- [19] L. Dhulipala, G. Blelloch, and J. Shun. Low-latency graph streaming using compressed purely-functional trees. *Proceedings of the ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*.
- [20] D. Ediger, R. McColl, J. Riedy, and D. A. Bader. Stinger: High performance data structure for streaming graphs. In *2012 IEEE Conference on High Performance Extreme Computing*, pages 1–5. IEEE, 2012.
- [21] M. Ester, H.-P. Kriegel, J. Sander, and X. Xu. A density-based algorithm for discovering clusters in large spatial databases with noise. pages 226–231. AAAI Press, 1996.
- [22] Y. Fang, R. Cheng, S. Luo, and J. Hu. Effective community search for large attributed graphs. *Proc. VLDB Endow.*, 9(12):1233–1244, Aug. 2016.
- [23] Y. Fang, R. Cheng, S. Luo, and J. Hu. Effective community search for large attributed graphs. *Proc. VLDB Endow.*, 9(12):1233–1244, Aug. 2016.
- [24] J. Feigenbaum, S. Kannan, A. McGregor, S. Suri, and J. Zhang. On graph problems in a semi-streaming model. *Theor. Comput. Sci.*, 348(2):207–216, Dec. 2005.
- [25] E. Georganas, R. Egan, S. Hofmeyr, E. Goltsman, B. Arndt, A. Tritt, A. Buluç, L. Olikek, and K. Yelick. Extreme scale de novo metagenome assembly. SC '18. IEEE Press, 2018.
- [26] J. Greiner. A comparison of parallel algorithms for connected components. In *Proceedings of the sixth annual ACM symposium on Parallel algorithms and architectures*, pages 16–25, 1994.
- [27] W.-S. Han, S. Lee, K. Park, J.-H. Lee, M.-S. Kim, J. Kim, and H. Yu. Turbograph: a fast parallel graph engine handling billion-scale graphs in a single pc. In *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 77–85, 2013.
- [28] L. He, Y. Chao, K. Suzuki, and K. Wu. Fast connected-component labeling. *Pattern recognition*, 42(9):1977–1987, 2009.
- [29] L. He, X. Ren, Q. Gao, X. Zhao, B. Yao, and Y. Chao. The connected-component labeling problem: A review of state-of-the-art algorithms. *Pattern Recognition*, 70:25–43, 2017.
- [30] M. M. Hossain, A. E. Hassanien, and M. Shoman. 3d brain tumor segmentation scheme using k-mean clustering and connected component labeling algorithms. In *2010 10th International Conference on Intelligent Systems Design and Applications*, pages 320–324, 2010.
- [31] A. Iyer, L. E. Li, and I. Stoica. Celliq: Real-time cellular network analytics at scale. In *12th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 15)*, pages 309–322, 2015.
- [32] A. P. Iyer, L. E. Li, T. Das, and I. Stoica. Time-evolving graph processing at scale. In *Proceedings of the fourth international workshop on graph data management experiences and systems*, pages 1–6, 2016.
- [33] J. Jung, K. Shin, L. Sael, and U. Kang. Random walk with restart on large graphs using block elimination. *ACM Transactions on Database Systems (TODS)*, 41(2):1–43, 2016.
- [34] U. Kang and C. Faloutsos. Beyond 'caveman communities': Hubs and spokes for graph compression and mining. pages 300–309, 12 2011.
- [35] U. Kang, M. McGlohon, L. Akoglu, and C. Faloutsos. Patterns on the connected components of terabyte-scale graphs. pages 875–880, 12 2010.
- [36] M. Korn, D. Sanders, and J. Pauli. Moving object detection by connected component labeling of point cloud registration outliers on the gpu. In *VISIGRAPP (6: VISAPP)*, pages 499–508, 2017.
- [37] A. Krishnamurthy, S. Lumetta, D. E. Culler, and K. Yelick. Connected components on distributed memory machines. *Third DIMACS Implementation Challenge*, 30:1–21, 1997.
- [38] A. Kyrola, G. Blelloch, and C. Guestrin. Grapchi: Large-scale graph computation on just a pc. In *10th USENIX Symposium on Operating Systems Design and Implementation*, pages 31–46. USENIX, 2012.
- [39] W. Lee, J. J. Lee, and J. Kim. Social network community detection using strongly connected components. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, pages 596–604. Springer, 2014.
- [40] J. Leskovec, L. A. Adamic, and B. A. Huberman. The dynamics of viral marketing. *ACM Transactions on the Web (TWEB)*, 1(1):5–es, 2007.
- [41] J. Leskovec and A. Krevl. SNAP Datasets: Stanford large network dataset collection. <http://snap.stanford.edu/data>, June 2014.
- [42] Y. Lim, U. Kang, and C. Faloutsos. Slashburn: Graph compression and mining beyond caveman communities. *IEEE Transactions on Knowledge and Data Engineering*, 26(12):3077–3089, 2014.
- [43] Y. Lim, W.-J. Lee, H.-J. Choi, and U. Kang. Discovering large subsets with high quality partitions in real world graphs. In *2015 International Conference on Big Data and Smart Computing (BIGCOMP)*, pages 186–193. IEEE, 2015.
- [44] Y. Lim, W.-J. Lee, H.-J. Choi, and U. Kang. Mtp: discovering high quality partitions in real world graphs. *World Wide Web*, 20(3):491–514, 2017.
- [45] H. Liu and H. H. Huang. Graphene: Fine-grained {IO} management for graph computing. In *15th {USENIX} Conference on File and Storage Technologies ({FAST} 17)*, pages 285–300, 2017.
- [46] S. Maass, C. Min, S. Kashyap, W. Kang, M. Kumar, and T. Kim. Mosaic: Processing a trillion-edge graph on a single machine. In *Proceedings of the Twelfth European Conference on Computer Systems*, pages 527–543, 2017.
- [47] P. Macko, V. J. Marathe, D. W. Margo, and M. I. Seltzer. Llama: Efficient graph analytics using large multiversioned arrays. In *2015 IEEE 31st International Conference on Data Engineering*, pages 363–374. IEEE, 2015.
- [48] J. J. McAuley and J. Leskovec. Learning to discover social circles in ego networks. In *NIPS*, volume 2012, pages 548–56. Citeseer, 2012.
- [49] D. Medini, A. Covacci, and C. Donati. Protein homology network families reveal step-wise diversification of type iii and type iv secretion systems. *PLoS*

- computational biology*, 2(12):e173, 2006.
- [50] D. G. Murray, F. McSherry, M. Isard, R. Isaacs, P. Barham, and M. Abadi. Incremental, iterative data processing with timely dataflow. *Communications of the ACM*, 59(10):75–83, 2016.
  - [51] S. Muthukrishnan. Data streams: Algorithms and applications. *Foundations and Trends® in Theoretical Computer Science*, 1(2):117–236, 2005.
  - [52] J. Nelson and H. Yu. Optimal lower bounds for distributed and streaming spanning forest computation. In *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 1844–1860. SIAM, 2019.
  - [53] J. Nešetřil, E. Milková, and H. Nešetřilová. Otakar borůvka on minimum spanning tree problem translation of both the 1926 papers, comments, history. *Discrete Math.*, 233(1–3):3–36, Apr. 2001.
  - [54] S. Nurk, D. Meleshko, A. Korobeynikov, and P. Pevzner. Metaspades: A new versatile metagenomic assembler. *Genome Research*, 27:gr.213959.116, 03 2017.
  - [55] OpenMP Architecture Review Board. *OpenMP Application Programming Interface*, 5th edition, Nov. 2018.
  - [56] P. Pandey, B. Wheatman, H. Xu, and A. Buluc. Terrace: A hierarchical graph container for skewed dynamic graphs. In *Proceedings of the 2021 International Conference on Management of Data*, pages 1372–1385, 2021.
  - [57] M. M. A. Patwary, D. Palsetia, A. Agrawal, W.-k. Liao, F. Manne, and A. Choudhary. A new scalable parallel dbscan algorithm using the disjoint-set data structure. In *SC '12: Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis*, pages 1–11, 2012.
  - [58] A. Pothen and C.-J. Fan. Computing the block triangular form of a sparse matrix. 16(4):303–324, Dec. 1990.
  - [59] M. Ripeanu, I. Foster, and A. Iamnitchi. Mapping the gnutella network: Properties of large-scale peer-to-peer systems and implications for system design. *arXiv preprint cs/0209028*, 2002.
  - [60] R. A. Rossi and N. K. Ahmed. The network data repository with interactive graph analytics and visualization. In *AAAI*, 2015.
  - [61] A. Roy, I. Mihailovic, and W. Zwaenepoel. X-stream: Edge-centric graph processing using streaming partitions. In *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles*, pages 472–488, 2013.
  - [62] S. Sahu, A. Mhedhbi, S. Salihoglu, J. Lin, and M. T. Özsu. The ubiquity of large graphs and surprising challenges of graph processing. *Proc. VLDB Endow.*, 11(4):420–431, Dec. 2017.
  - [63] D. Sengupta and S. L. Song. Evograph: On-the-fly efficient mining of evolving graphs on gpu. In *International Supercomputing Conference*, pages 97–119. Springer, 2017.
  - [64] D. Sengupta, N. Sundaram, X. Zhu, T. L. Willke, J. Young, M. Wolf, and K. Schwan. Graphin: An online high performance incremental graph processing framework. In *European Conference on Parallel Processing*, pages 319–333. Springer, 2016.
  - [65] H. Thornquist, E. Keiter, R. Hoekstra, D. Day, and E. Boman. A parallel preconditioning strategy for efficient transistor-level circuit simulation. pages 410–417, 01 2009.
  - [66] S. van Dongen. Graph clustering by flow simulation. Technical report, 2000.
  - [67] J. S. Vitter. External memory algorithms and data structures: Dealing with massive data. *ACM Computing Surveys*, 33:2001, 2001.
  - [68] K. Vora. {LUMOS}: Dependency-driven disk-based graph processing. In *2019 {USENIX} Annual Technical Conference*, pages 429–442, 2019.
  - [69] K. Vora, G. Xu, and R. Gupta. Load the edges you need: A generic i/o optimization for disk-based graph processing. In *2016 {USENIX} Annual Technical Conference*, pages 507–522, 2016.
  - [70] D. Wen, L. Qin, Y. Zhang, L. Chang, and X. Lin. Efficient structural graph clustering: An index-based approach. *The VLDB Journal*, 28(3):377–399, June 2019.
  - [71] D. Wen, L. Qin, Y. Zhang, L. Chang, and X. Lin. Efficient structural graph clustering: An index-based approach. *The VLDB Journal*, 28(3):377–399, June 2019.
  - [72] M. Wu, X. li, C.-K. Kwok, and S.-K. Ng. A core-attachment based method to detect protein complexes in ppi networks. *BMC bioinformatics*, 10:169, 02 2009.
  - [73] P. Yuan, C. Xie, L. Liu, and H. Jin. Pathgraph: A path centric graph processing system. *IEEE Transactions on Parallel and Distributed Systems*, 27(10):2998–3012, 2016.
  - [74] M. Zhang, Y. Wu, Y. Zhuo, X. Qian, C. Huan, and K. Chen. Wonderland: A novel abstraction-based out-of-core graph processing system. *ACM SIGPLAN Notices*, 53(2):608–621, 2018.
  - [75] D. Zheng, D. Mhembere, R. Burns, J. Vogelstein, C. E. Priebe, and A. S. Szalay. Flashgraph: Processing billion-node graphs on an array of commodity ssds. In *13th {USENIX} conference on file and storage technologies ({FAST} 15)*, pages 45–58, 2015.
  - [76] X. Zhu, W. Han, and W. Chen. Gridgraph: Large-scale graph processing on a single machine using 2-level hierarchical partitioning. In *2015 {USENIX} Annual Technical Conference*.