



**HAL**  
open science

## Triangle Counting in Dynamic Graph Streams

Laurent Bulteau, Vincent Froese, Konstantin Kutzkov, Rasmus Pagh

► **To cite this version:**

Laurent Bulteau, Vincent Froese, Konstantin Kutzkov, Rasmus Pagh. Triangle Counting in Dynamic Graph Streams. *Algorithmica*, 2016, 76 (1), pp.259 - 278. 10.1007/s00453-015-0036-4 . hal-01494399

**HAL Id: hal-01494399**

**<https://hal.science/hal-01494399>**

Submitted on 23 Mar 2017

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Triangle counting in dynamic graph streams

Laurent Bulteau <sup>\*1</sup>, Vincent Froese<sup>†1</sup>, Konstantin Kutzkov<sup>‡2</sup>, and Rasmus Pagh<sup>§3</sup>

<sup>1</sup>Technische Universität Berlin, Germany

<sup>2</sup>NEC Laboratories Europe, Heidelberg, Germany

<sup>3</sup>IT University of Copenhagen, Denmark

## Abstract

Estimating the number of triangles in graph streams using a limited amount of memory has become a popular topic in the last decade. Different variations of the problem have been studied, depending on whether the graph edges are provided in an arbitrary order or as incidence lists. However, with a few exceptions, the algorithms have considered *insert-only* streams. We present a new algorithm estimating the number of triangles in *dynamic* graph streams where edges can be both inserted and deleted. We show that our algorithm achieves better time and space complexity than previous solutions for various graph classes, for example sparse graphs with a relatively small number of triangles. Also, for graphs with constant transitivity coefficient, a common situation in real graphs, this is the first algorithm achieving constant processing time per edge. The result is achieved by a novel approach combining sampling of vertex triples and sparsification of the input graph. In the course of the analysis of the algorithm we present a lower bound on the number of pairwise independent 2-paths in general graphs which might be of independent interest. At the end of the paper we discuss lower bounds on the space complexity of triangle counting algorithms that make no assumptions on the structure of the graph.

## 1 Introduction

Many relationships between real life objects can be abstractly represented as graphs. The discovery of certain structural properties in a graph, which abstractly describes a given real-life problem, can often provide important insights into the nature of the original problem. The number of triangles, and the closely related clustering and transitivity coefficients, have proved to be an important measure used in applications ranging from social network analysis and spam detection to motif detection in protein interaction networks. We refer to [27] for a detailed discussion on the applications of triangle counting.

The best known algorithm for triangle counting in the RAM model runs in time  $O(m^{\frac{2\omega}{\omega+1}})$  [4] where  $\omega$  is the matrix multiplication exponent, the best known bound is  $\omega = 2.3727$  [28]. However, this algorithm is mainly of theoretical importance since exact fast matrix multiplication algorithms do not admit an efficient implementation for input matrices of reasonable size.

The last decade has witnessed a rapid growth of available data. This has led to a shift in attitudes in algorithmic research and solutions storing the whole input in main memory are not any more considered a feasible choice for many real-life problems. Classical algorithms have been adjusted in order to cope with the new requirements and many new techniques have been developed. This has led to the *streaming* model of computation where only a single scan of the data is possible. Both efficient algorithms and impossibility

---

\*l.bulteau@gmail.com, Supported by the Alexander von Humboldt Foundation, Bonn, Germany.

†vincent.froese@tu-berlin.de, Supported by the DFG project DAMM (NI 369/13).

‡kutzkov@gmail.com, Work done while the author was at IT University of Copenhagen and supported by the Danish National Research Foundation under the Sapere Aude program.

§pagh@itu.dk, Supported by the Danish National Research Foundation under the Sapere Aude program.

results have shed light on the computational complexity of many problems in the model [19].

**Approximate triangle counting in streamed graphs.** For many applications one is satisfied with a good approximation of the number of triangles instead of their exact number, thus researchers have designed randomized approximation algorithms returning with high probability a precise estimate using only small amount of main memory. Two models of streamed graphs have been considered in the literature. In the *incidence list stream* model the edges incident to each vertex arrive consecutively and in the *adjacency stream* model edges arrive in arbitrary order. Also, a distinction has been made between algorithms using only a single pass over the input, and algorithms assuming that the input graph can be persistently stored on a secondary device and multiple passes are allowed. A simple approach for estimating the number of triangles in insert-only streams is to sample a certain number of 2-paths, then compute the ratio of 2-paths in the sample that are completed to triangles and multiply the obtained number with the total number of 2-paths in the graphs. For incidence list streams this is easy since we can assume that the stream consists of (an implicit representation of) all 2-paths [8]. For the more difficult model of adjacency streams where edges arrive in arbitrary order the approach was adjusted such that we sample a random 2-path [12, 22]. The one-pass algorithm with the best known space complexity and constant processing time per edge for adjacency streams is due to Pavan et al. [22], and when several passes are allowed – by Kolountzakis et al. [15]. For a more detailed overview of results and developed techniques we refer to [27].

Dynamic graph streams have a wider range of applications. Consider for example a social network like Facebook where one is allowed to befriend and “unfriend” other members, or join and leave groups of interest. Estimating the number of triangles in a network is a main building block in algorithms for the detection of emerging communities [7], and thus it is required that triangle counting algorithms can also handle edge deletions. The problem of designing triangle counting algorithms for dynamic streams matching the space and time complexity of algorithms for insert-only streams has been presented as an open question in the 2006 IITK Workshop on Algorithms for Data Streams [17]. The best known algorithms for insert-only streams work by sampling a non-empty subgraph on three vertices from the stream (e.g. an edge  $(u, v)$  and a vertex  $w$ ). Then one checks whether the arriving edges will complete the sampled subgraph to a triangle (we look for  $(u, w)$  and  $(v, w)$ ). The approach does not work for dynamic streams because an edge in the sampled subgraph might be deleted later. Proposed solutions [1, 18, 11] have explored different ideas. These approaches, however, only partially resolve the open problem from [17] because of high processing time per edge update, see Section 3 for more details.

**Our contribution.** In this work we propose a method to adjust sampling to work in dynamic streams and show that for graphs with constant transitivity coefficient, a ubiquitous assumption for real-life graphs, we can achieve constant processing time per edge. We also show that *some* assumption on the graph is needed to be able to estimate small triangle counts, by showing a lower bound on the space usage in terms of the number of edges and triangles, matching the upper bound of Manjunath et al. [18] for constant approximation factor.

At a very high level, the main technical contribution of the present work can be summarized as follows: For dynamic graph streams sampling-based approaches fail because we don’t know how many of the sampled subgraphs will survive after edges have been deleted. On the other hand, graph sparsification approaches [21, 26, 27] can handle edge deletions but the theoretical guarantees on the complexity of the algorithms depend on specific properties of the underlying graph, e.g., the maximum number of triangles an edge is part of. The main contribution in the present work is a novel technique for sampling 2-paths *after* the stream has been processed. It is based on the combination of standard 2-path sampling with graph sparsification. The main technical challenge is to show that sampling at random a 2-path in a sparsified graph is (almost) equivalent to sampling at random a 2-path in the original graph. In the course of the analysis, we also obtain combinatorial results about general graphs that might be of independent interest.

## 2 Preliminaries

**Notation.** A simple undirected graph without loops is denoted as  $G = (V, E)$  with  $V = \{1, 2, \dots, n\}$  being a set of vertices and  $E$  a set of edges. The edges are provided as a stream of insertions and deletions in arbitrary order. We assume the strict turnstile model where each edge can be deleted only after being inserted. We assume that  $n$  is known in advance<sup>1</sup> and that the number of edges cannot exceed  $m$ . For an edge connecting the vertices  $u$  and  $v$  we write  $(u, v)$  and  $u$  and  $v$  are the *endpoints* of the edge  $(u, v)$ . Vertex  $u$  is *neighbor* of  $v$  and vice versa and  $N(u)$  is the set of  $u$ 's neighbors. We say that edge  $(u, v)$  is *isolated* if  $|N(u)| = |N(v)| = 1$ . We consider only edges  $(u, v)$  with  $u < v$ . A *2-path centered at  $v$* ,  $(u, v, w)$ , consists of the edges  $(u, v)$  and  $(v, w)$ . A  $k$ -clique in  $G$  is a subgraph of  $G$  on  $k$  vertices  $v_1, \dots, v_k$  such that  $(v_i, v_j) \in E$  for all  $1 \leq i < j \leq k$ . A 3-clique on  $u, v, w$  is called a *triangle* on  $u, v, w$ , and is denoted as  $\langle u, v, w \rangle$ . We denote by  $P_2(v)$  the number of 2-paths centered at a vertex  $v$ , and  $P_2(G) = \sum_{v \in V} P_2(v)$  and  $T_3(G)$  the number of 2-paths and number of triangles in  $G$ , respectively. We will omit  $G$  when clear from the context.

We say that two 2-paths are *independent* if they have at most one common vertex. The *transitivity coefficient* of  $G$  is

$$\alpha(G) = \frac{3T_3}{\sum_{v \in V} \binom{d_v}{2}} = \frac{3T_3}{P_2},$$

i.e., the ratio of 2-paths in  $G$  contained in a triangle to all 2-paths in  $G$ . When clear from the context, we will omit  $G$ .

**Hashing.** A family  $\mathcal{F}$  of functions from  $U$  to a finite set  $S$  is  *$k$ -wise independent* if for a function  $f : U \rightarrow S$  chosen uniformly at random from  $\mathcal{F}$  it holds

$$\Pr[f(u_1) = c_1 \wedge f(u_2) = c_2 \wedge \dots \wedge f(u_k) = c_k] = 1/s^k$$

for  $s = |S|$ , distinct  $u_i \in U$  and any  $c_i \in S$  and  $k \in \mathbb{N}$ . We will call a function chosen uniformly at random from a  $k$ -wise independent family  *$k$ -wise independent function* and a function  $f : U \rightarrow S$  *fully random* if  $f$  is  $|U|$ -wise independent. We will say that a function  $f : U \rightarrow S$  *behaves like a fully random function* if for any set of input from  $U$ , with high probability  $f$  has the same probability distribution as a fully random function.

We will say that an algorithm returns an  $(\varepsilon, \delta)$ -*approximation* of some quantity  $q$  if it returns a value  $\tilde{q}$  such that  $(1 - \varepsilon)q \leq \tilde{q} \leq (1 + \varepsilon)q$  with probability at least  $1 - \delta$  for every  $0 < \varepsilon, \delta < 1$ .

**Probability inequalities.** In the analysis of the algorithm we use the following inequalities.

- *Chebyshev's inequality.* Let  $X$  be a random variable and  $\lambda > 0$ . Then

$$\Pr[|X - \mathbb{E}[X]| \geq \lambda \sigma(X)] \leq \frac{1}{\lambda^2}$$

- *Chernoff's inequality.* Let  $X_1, \dots, X_\ell$  be  $\ell$  independent identically distributed Bernoulli random variables and  $\mathbb{E}[X_i] = \mu$ . Then for any  $\varepsilon > 0$  we have

$$\Pr\left[\left|\frac{1}{\ell} \sum_{i=1}^{\ell} X_i - \mu\right| > \varepsilon \mu\right] \leq 2e^{-\varepsilon^2 \mu \ell / 2}$$

## 3 Results

The following theorem is our main result.

---

<sup>1</sup>More generally, our results hold when the  $n$  vertices come from some arbitrary universe  $U$  known in advance.

	Space	Update time
Ahn et al. [1]	$O(\frac{mn}{\varepsilon^2 T_3} \log \frac{1}{\delta})$	$O(n \log n)$
Manjunath et al. [18]	$O(\frac{m^3}{\varepsilon^2 T_3^2} \log \frac{1}{\delta})$	$O(\frac{m^3}{\varepsilon^2 T_3^2} \log \frac{1}{\delta})$
Section 5	$\Omega(\frac{m^3}{T_3^2})$	—
Section 4	$O(\frac{\sqrt{m}}{\varepsilon^3 \alpha} \log \frac{1}{\delta})$	$O(\frac{1}{\varepsilon^2 \alpha} \log \frac{1}{\delta})$

Table 1: Overview of time and space bounds.

	$n \log n$	$m^3/T_3^2$
$z < 1/2$	$T_3 = \omega(C^2/(n^{2z} \log n))$	$T_3 = o(Cn^{2-z})$
$1/2 < z < 1$	$T_3 = \omega(C^2/(n \log n))$	$T_3 = o(Cn^{3-3z})$
$z > 1$	$T_3 = \omega(C^2/(n \log n))$	$T_3 = o(C)$

Table 2: Comparison of the theoretical guarantees for the per edge processing time for varying  $z$ .

**Theorem 1** *Let  $G = (V, E)$  be a graph given as a stream of edge insertions and deletions with no isolated edges and vertices,  $V = \{1, 2, \dots, n\}$  and  $|E| \leq m$ . Let  $P_2, T_3$  and  $\alpha$  be the number of 2-paths, number of triangles and the transitivity coefficient of  $G$ , respectively. Let  $\varepsilon, \delta \in (0, 1)$  be user defined. Assuming fully random hash functions, there exists a one-pass algorithm running in expected space  $O(\frac{\sqrt{m}}{\varepsilon^3 \alpha} \log \frac{1}{\delta})$  and  $O(\frac{1}{\varepsilon^2 \alpha} \log \frac{1}{\delta})$  processing time per edge. After processing the stream, an  $(\varepsilon, \delta)$ -approximation of  $T_3$  can be computed in expected time  $O(\frac{\log n}{\varepsilon^2 \alpha} \log \frac{1}{\delta})$  and worst case time  $O(\frac{\log^2 n}{\varepsilon^2 \alpha} \log \frac{1}{\delta})$  with high probability.*

(For simplicity, we assume that there are no isolated edges in  $G$ . More generally, the result holds by replacing  $n$  with  $n_C$ , where  $n_C$  is the number of vertices in connected components with at least two edges. We assume  $m$  and  $n$  can be described in  $O(1)$  words.)

Before presenting the algorithm, let us compare the above to the bounds in [1, 18]. The algorithm in [1] estimates  $T_3$  by applying  $\ell_0$  sampling [10] to non-empty subgraphs on 3 vertices. There are  $O(mn)$  such subgraphs, thus  $O(\frac{mn}{\varepsilon^2 T_3} \log \frac{1}{\delta})$  samples are needed for an  $(\varepsilon, \delta)$ -approximation. However, each edge insertion or deletion results in the update of  $n - 2$  non-empty subgraphs on 3 vertices. Using the  $\ell_0$  sampling algorithm from [13], this results in processing time of  $O(n \log n)$  per edge. The algorithm by Manjunath et al. [18] (which builds upon the the work of Jowhari and Ghodsi [11]) estimates the number of triangles (and more generally of cycles of fixed length) in streamed graphs by computing complex valued sketches of the stream. Each of them yields an unbiased estimator of  $T_3$ . The average of  $O(\frac{m^3}{\varepsilon^2 T_3^2} \log \frac{1}{\delta})$  estimators is an  $(\varepsilon, \delta)$ -approximation of  $T_3$ . However, each new edge insertion or deletion has to update all estimators, resulting in update time of  $O(\frac{m^3}{\varepsilon^2 T_3^2} \log \frac{1}{\delta})$ . The algorithm was generalized to counting arbitrary subgraphs of fixed size in [14].

The time and space bounds are summarized in Table 1. Comparing our space complexity to the bounds in [1, 18], we see that for several graph classes our algorithm is more time and space efficient. (We ignore  $\varepsilon$  and  $\delta$  and logarithmic factors in  $n$  for the space complexity.) For  $d$ -regular graphs the processing time per edge is better than  $O(n \log n)$  for  $T_3 = \omega(d^2 / \log n)$ , and better than  $O(m^3/T_3^2)$  for  $T_3 = o(n^2 d)$ . Our space bound is better than  $O(mn/T_3)$  when  $d = o(n^{1/3})$ , and better than  $O(m^3/T_3^2)$  for  $T_3 = o(n^{3/2} d^{1/2})$ . Most real-life graphs exhibit a skewed degree distribution adhering to some form of power law, see for example [2]. Assume vertices are sorted according to their degree in decreasing order such that the  $i$ th vertex has degree  $C/i^z$  for some  $C \leq n$ , and constant  $z > 0$ , i.e., we have Zipfian distribution with parameter  $z$ . It holds  $\sum_{i=1}^n i^{-z} = O(n^{1-z})$  for  $z < 1$  and  $\sum_{i=1}^n i^{-z} = O(1)$  for  $z > 1$ . Table 2 summarizes for which values of  $T_3$  our algorithm achieves faster processing time than [1, 18], and Table 3 – for which values of  $C$  our algorithm is more space-efficient than [1], and for which values of  $T_3$  – more space-efficient than [18].

	$mn/T_3$	$m^3/T_3^2$
$z < 1/2$	$C = o(n^{1/3+z})$	$T_3 = o(C^{1/2} n^{\frac{3-z}{2}})$
$1/2 < z < 1$	$C = o(n^{1-z/3})$	$T_3 = o(C^{1/2} n^{\frac{5-5z}{2}})$
$z > 1$	$C = o(n^{2/3})$	$T_3 = o(C^{1/2})$

Table 3: Comparison of the theoretical guarantees for the space usage for varying  $z$ .

However, the above values are for arbitrary graphs adhering to a certain degree distribution. We consider the main advantage of the new algorithm to be that it achieves constant processing time per edge for graphs with constant transitivity coefficient. This is a common assumption for real-life networks, see for instance [3, 8]. Note that fast update is essential for real life applications. Consider for example the Facebook graph. In May 2011, for less than eight years existence, there were about 69 billion friendship links [6]. This means an average of above 300 new links per second, without counting deletions and peak hours.

In Section 6 we compare the theoretical guarantees for several real life graphs. While such a comparison is far from being a rigorous experimental evaluation, it clearly indicates that the processing time per edge in [1, 18] is prohibitively large and the assumption that the transitivity coefficient is constant is justified. Also, for graphs with a relatively small number of triangles our algorithm is much more space-efficient.

## 4 Our algorithm

The main idea behind our algorithm is to design a new sampling technique for dynamic graph streams. It exploits a combination of the algorithms by Buriol et al. [8] for the incidence stream model, and the Doulion algorithm [26] and its improvement [21]. Let us briefly describe the approaches.

**The Buriol et al. algorithm for incidence list streams.** Assume we know the total number of 2-paths in  $G$ . One chooses at random one of them, say  $(u, v, w)$ , and checks whether the edge  $(u, w)$  appears later in the stream. For a triangle  $\langle u, v, w \rangle$  the three 2-paths  $(u, v, w)$ ,  $(w, u, v)$ ,  $(v, w, u)$  appear in the incidence list stream, thus the probability that we sample a triangle is exactly  $\alpha$ . One chooses independently at random  $K$  2-paths and using standard techniques shows that for  $K = O(\frac{1}{\varepsilon^2 \alpha} \log \frac{1}{\delta})$  we compute an  $(\varepsilon, \delta)$ -approximation of  $\alpha(G)$ . One can get rid of the assumption that the number of 2-paths is known in advance by running  $O(\log n)$  copies of the algorithm in parallel, each guessing the right value. The reader is referred to the original work for more details. For incidence streams, the number of 2-paths in  $G$  can be computed exactly by updating a single counter, thus if  $\tilde{\alpha}$  is an  $(\varepsilon, \delta)$ -approximation of the transitivity coefficient then  $\tilde{T}_3 = \tilde{\alpha}P_2$  is an  $(\varepsilon, \delta)$ -approximation of  $T_3$ .

**Doulion and monochromatic sampling.** The Doulion algorithm [26] is a simple and intuitive sparsification approach. Each edge is sampled independently with probability  $p$  and added to a sparsified graph  $G_S$ . We expect  $pm$  edges to be sampled and a triangle survives in  $G_S$  with probability  $p^3$ , thus multiplying the number of triangles in  $G_S$  by  $1/p^3$  we obtain an estimate of  $T_3$ . The algorithm was improved in [21] by using *monochromatic sampling*. Instead of throwing a biased coin for each edge, we uniformly at random color each vertex with one of  $1/p$  colors. Then we keep an edge in the sparsified graph iff its endpoints have the same color. A triangle survives in  $G_S$  with probability  $p^2$ . It is shown that for a fully random coloring the variance of the estimator is better than in Doulion. However, in both algorithms it depends on the maximum number of triangles an edge is part of, and one might need constant sampling probability in order to obtain an  $(\varepsilon, \delta)$ -approximation on  $T_3$ . The algorithm can be applied to dynamic streams because one counts the number of triangles in the sparsified graph after all edges have been processed. However, it can be expensive to obtain an estimate since the exact number of triangles in  $G_S$  is required.

**Combining the above approaches.** The basic idea behind the new algorithm is to use the estima-

tor of Buriol et al. for the incidence stream model: (i) estimate the transitivity coefficient  $\alpha(G)$  by choosing a sufficiently large number of 2-paths at random and check which of them are part of a triangle, and (ii) estimate the number of 2-paths  $P_2$  in the graph. We first observe that estimating  $P_2$  in dynamic graph streams can be reduced to second moment estimation of streams of items in the turnstile model, see e.g. [25]. For (i), we will estimate  $\alpha(G)$  by adjusting the monochromatic sampling approach. Its main advantage compared to the sampling of edges separately is that if we have sampled the 2-path  $(u, v, w)$ , then we must also have sampled the edge  $(u, w)$ , if existent. So, the idea is to use monochromatic sampling and then in the sparsified graph to pick up at random a 2-path and check whether it is part of a triangle. Instead of random coloring of the vertices, we will use a suitably defined hash function and we will choose a sampling probability guaranteeing that for a graph with no isolated edges (or rather a small number of isolated edges) the sparsified graph will contain a sufficiently big number of 2-paths. A 2-path in the sparsified graph picked up at random, will then be used to estimate  $\alpha(G)$ . Thus, unlike in [8], we sample *after* the stream has been processed and this allows to handle edge deletions. The main technical obstacles are to analyze the required sampling probability  $p$  and to show that this sampling approach indeed provides an unbiased estimator of  $\alpha(G)$ . We will obtain bounds on  $p$  and show that even if the estimator might be biased, the bias can be made arbitrarily small and one can still achieve an  $(\varepsilon, \delta)$ -approximation of  $\alpha(G)$ . Also, we present an implementation for storing a sparsified graph  $G_S$  such that each edge is added or deleted in constant time and a random 2-path in  $G_S$ , if existent, can be picked up without explicitly considering all 2-paths in  $G_S$ .

## 4.1 Algorithm details

Pseudocode description of the algorithm is given in Figure 1. We assume that the graph is given as a stream  $\mathcal{S}$  of pairs  $((u, v), \$)$ , where  $(u, v) \in E$  and  $\$ \in \{+, -\}$  with the obvious meaning that the edge  $(u, v)$  is inserted or deleted from  $G$ . In ESTIMATENUMBEROFTWOPATHS each incoming pair  $((u, v), \$)$  is treated as the insertion, respectively deletion, of two items  $u$  and  $v$ , and these update a second moment estimator  $SME$ , working as a blackbox algorithm. We refer to the proof of Lemma 1 for more details. In SPARSIFYGRAPH we assume access to a fully random coloring hash function  $f : V \rightarrow C$ . Each edge  $(u, v)$  is inserted/deleted to/from a sparsified graph  $G_S$  iff  $f(u) = f(v)$ . At the end  $G_S$  consists of all monochromatic edges that have not been deleted. In ESTIMATENUMBEROFTRIANGLES we run in parallel the algorithm estimating  $P_2$  and  $K$  copies of SAMPLERANDOM2PATH. For each  $G_S^i$ ,  $1 \leq i \leq K$ , with at least  $s$  pairwise independent 2-paths we choose at random a 2-path and check whether it is a triangle. (Note that we require the existence of  $s$  pairwise independent 2-paths but we choose a 2-path at random from *all* 2-paths in  $G_S$ .) The ratio of triangles to all sampled 2-paths and the estimate of  $P_2$  are then used to estimate  $T_3$ . In the next section we obtain bounds on the user defined parameters  $C, K$  and  $s$ . In Lemma 7 we present an efficient implementation of  $G_S$  that guarantees constant time updates and allows the sampling of a random 2-path in expected time  $O(\log n)$  and worst case time  $O(\log^2 n)$  with high probability.

## 4.2 Theoretical analysis

We will prove the main result in several lemmas. The first lemma provides an estimate of  $P_2$  using an estimator for the second frequency moment of data streams [25].

**Lemma 1** *Let  $G$  be a graph with no isolated edges given as a stream of edge insertions and deletions. There exists an algorithm returning an  $(\varepsilon, \delta)$ -approximation of the number of 2-paths in  $G$  in one pass over the stream of edges which needs  $O(\frac{1}{\varepsilon^2} \log \frac{1}{\delta})$  space and  $O(\log \frac{1}{\delta})$  processing time per edge.*

**Proof:** We show that ESTIMATENUMBEROFTWOPATHS in Figure 1 returns an  $(\varepsilon, \delta)$ -approximation of the number of 2-paths in  $G$ . We reduce the problem of computing the number of 2-paths in dynamic graph streams to the problem of estimating the second frequency moment in streams of items in the turnstile model. By associating vertices with items and treating each incoming pair  $((u, v), \$)$  as the insertion, respectively deletion, of two new items  $u$  and  $v$ , it is a simple observation that the second moment of the so defined

#### ESTIMATENUMBEROFTWOPATHS

**Input:** stream of edge deletions and insertions  $\mathcal{S}$ , algorithm  $SME$  estimating the second moment items streams

```

1:  $m = 0$ 
2: for each  $((u, v), \$)$  in  $\mathcal{S}$  do
3:   if  $\$ = +$  then
4:      $m = m + 1$ 
5:      $SME.update(u, 1), SME.update(v, 1)$ 
6:   else
7:      $m = m - 1$ 
8:      $SME.update(u, -1), SME.update(v, -1)$ 
9: return  $SME.estimate/2 - m$ 

```

#### SPARSIFYGRAPH

**Input:** stream of edge deletions and insertions  $\mathcal{S}$ , coloring function  $f : V \rightarrow C$

```

1:  $G_S = \emptyset$ 
2: for each  $((u, v), \$) \in \mathcal{S}$  do
3:   if  $f(u) = f(v)$  then
4:     if  $\$ = +$  then
5:        $G_S = G_S \cup (u, v)$ .
6:     else
7:        $G_S = G_S \setminus (u, v)$ .
8: Return  $G_S$ .

```

#### SAMPLERANDOM2PATH

**Input:** sparsified graph  $G_S$

```

1: choose at random a 2-path  $(u, v, w)$  in  $G_S$ 
2: if the vertices  $\{u, v, w\}$  form a triangle then
3:   return 1
4: else
5:   return 0

```

#### ESTIMATENUMBEROFTRIANGLES

**Input:** streamed graph  $\mathcal{S}$ , set of  $K$  independent fully random coloring functions  $\mathcal{F}$ , algorithm  $SME$  estimating the second moment of streams of items, threshold  $s$

```

1: run in parallel  $ESTIMATENUMBEROFTWOPATHS(\mathcal{S}, SME)$  and let  $\tilde{P}_2$  be the returned estimate
2: run in parallel  $K$  copies of  $SPARSIFYGRAPH(\mathcal{S}, f_i), f_i \in \mathcal{F}$ 
3:  $\ell = 0$ 
4: for each  $G_S^i$  with at least  $s$  pairwise independent 2-paths do
5:    $X_+ = SAMPLERANDOM2PATH(G_S^i)$ 
6:    $\ell_+ = 1$ 
7:  $\tilde{\alpha} = X/\ell$ 
8: return  $\frac{\tilde{\alpha}\tilde{P}_2}{3}$ 

```

Figure 1: Estimating the number of 2-paths in  $G$ , the transitivity coefficient and the number of triangles.



stream of items corresponds to  $F_2 = \sum_{v \in V} d_v^2$ . For the number of 2-paths in  $G$  we have

$$\sum_{v \in V} \binom{d_v}{2} = \left( \sum_{v \in V} d_v^2 \right) / 2 - m.$$

Since  $G$  contains no isolated edges, for each edge  $(u, v)$  we can assume that at least one of its endpoints has degree more than one, w.l.o.g.  $d_u \geq 2$ . Thus,  $(u, v)$  must be part of at least one 2-path, namely a 2-path centered at  $u$ . Each edge can be assigned thus to at least one 2-path and this implies a lower bound on the number of 2-paths in  $G$ ,  $\sum_{v \in V} \binom{d_v}{2} \geq m/2$ , thus

$$F_2 = 2 \sum_{v \in V} \binom{d_v}{2} + 2m \geq 3m.$$

Assume now that we have computed an  $(1 \pm \varepsilon/2)$ -approximation of  $F_2$  and we have the exact value of  $m$ . Then the over- and underestimation returned by ESTIMATENUMBEROFTWOPATHS is bounded by  $(\varepsilon/2)F_2$ . For  $F_2 \geq 3m$  it holds

$$(\varepsilon/2)F_2 \leq \varepsilon(F_2 - m) = \varepsilon P_2.$$

As for the complexity of ESTIMATENUMBEROFTWOPATHS, we observe that recording the exact value of  $m$  requires constant processing time per edge and  $O(\log n)$  bits. For *SME*, we use the algorithm from [25] for the estimation of the second moment, which computes an  $(\varepsilon, \delta)$ -approximation of  $\sum_{v \in V} d_v^2$  with the claimed time and space complexity.  $\square$

The next two lemmas show a lower bound of  $\Omega(m)$  on the number of pairwise independent 2-paths in a graph without isolated edges. The result is needed in order to obtain bounds on the required sampling probability. We first show that in order to obtain a lower bound for general graphs it is sufficient to consider bipartite connected graphs. This is true because every connected graph contains a cut with at least  $m/2$  edges such that the bipartite graph between the two vertex subsets is connected. We prove this first.

**Lemma 2** *Let  $G = (V, E)$  be an arbitrary connected graph with  $|E| = m$ . There exists a bipartition  $(U, W)$  of  $V$  such that the bipartite graph  $B := (V, E')$  with  $E' := E \cap \{\{u, w\} \mid u \in U, w \in W\}$  is connected and  $|E'| \geq m/2$ .*

**Proof:** Let  $G = (V, E)$  be a connected graph. We first show that there exists  $(U, W)$  such that  $B$  contains at least  $m/2$  edges. To this end, initialize  $U := \emptyset$  and  $W := V$ . Now, as long as there either exists a vertex  $v \in W$  with  $|N_G(v) \cap W| > |N_G(v) \cap U|$  or a vertex  $v \in U$  with  $|N_G(v) \cap U| > |N_G(v) \cap W|$ , we exchange this vertex  $v$ , that is, in the first case, we add  $v$  to  $U$  and delete it from  $W$  and vice versa in the second case. Clearly, this increases the number of edges between  $U$  and  $W$ ; hence, this procedure terminates after a finite number of steps. Moreover, after termination, it holds  $|N_G(v) \cap W| \leq |N_G(v) \cap U|$  for each  $v \in W$  and  $|N_G(v) \cap U| \leq |N_G(v) \cap W|$  for each  $v \in U$ . Thus, there are at least  $m/2$  edges in  $B$ .

Assume now that  $B$  is not connected and let  $B_1, \dots, B_c$  denote the  $c \geq 2$  connected components of  $B$  and let  $U_i \subseteq U, W_i \subseteq W$  be the bipartition of  $B_i$ . We simply “swap” an arbitrary component, say  $B_1$ , that is, we delete the vertices  $U_i$  from  $U$  and add them to  $W$  and delete the vertices  $W_i$  from  $W$  adding them to  $U$ . Clearly, the edges between  $U_i$  and  $W_i$  remain in  $B$ . Additionally, the edges between  $U_i$  and  $U \setminus U_i$  as well as the edges between  $W_i$  and  $W \setminus W_i$  are added to  $B$ . Hence, the number of edges in  $B$  increased. Moreover, since  $G$  is connected, the component  $B_1$  is now connected to some other component in  $B$  and thus the number of connected components decreased. This proves the claim.  $\square$

We also need, for the case where there are few edges in the graph, the following lower bound.

**Lemma 3** *Let  $G = (V, E)$  be a connected graph. The number of independent 2-paths in  $G$  is at least  $|V|/2 - 1$ .*

**Proof:** This can easily be obtained by taking a rooted spanning tree  $T$  of  $G$ . Assuming  $|V| \geq 3$ , take any leaf  $u$  of  $T$  with maximum depth (distance to the root). Let  $v$  be the parent of  $u$ . If  $v$  has another child  $w$ , then  $w$  is a leaf of  $T$ . Select the 2-path  $\{\{u, v\}, \{v, w\}\}$ , remove vertices  $u$  and  $w$  from  $V$ , and start again. Otherwise,  $v$  has a parent  $w$ : Similarly, select the 2-path  $\{\{u, v\}, \{v, w\}\}$ , remove vertices  $u$  and  $v$  from  $V$ , and start again. The set of 2-paths thus selected is independent (each one uses 2 new vertices) and contains at least  $(m-1)/2 = (|V|-2)/2$  2-paths.  $\square$

We now prove the lower bound for bipartite connected graphs.

**Lemma 4** *Let  $G = (V, E)$  with  $|E| = m$  be a connected bipartite graph. The number of independent 2-paths in  $G$  is at least  $\lfloor m/9 \rfloor$ .*

**Proof:** We say that a set of independent 2-paths  $\mathcal{P}$  is *maximal* if for any 2-path  $P$  of  $G$  not yet in  $\mathcal{P}$ ,  $\mathcal{P} \cup \{P\}$  is not independent. We first prove the following property by induction on  $|E|$ :

*In any bipartite graph  $G = (V, E)$ , if  $\mathcal{P}$  is a maximal set of independent 2-paths, then*

$$|\mathcal{P}| \geq \frac{m}{6} - \frac{|V|}{4}.$$

Let  $F$  be the set of edges which are not used in any 2-path of  $\mathcal{P}$ . Note that  $|\mathcal{P}| = \frac{m-|F|}{2}$ . We consider the subgraph  $H := (V, F)$  of  $G$ . First, consider the easy case where all vertices have degree at most 1 in  $H$ . It follows that  $|F| \leq \frac{|V|}{2}$ . Thus,

$$|\mathcal{P}| \geq \frac{m - \frac{|V|}{2}}{2} \geq \frac{m}{6} - \frac{|V|}{4}.$$

Otherwise, pick a vertex  $u$  having maximum degree in  $H$ , write  $k$  for its degree ( $k \geq 2$ ), and  $N(u)$  for the set of its  $k$  neighbors in  $H$ . Consider any two distinct vertices  $v, v' \in N(u)$ . Then  $\{\{v, u\}, \{u, v'\}\}$  is a 2-path of  $G$ . Since  $\mathcal{P}$  is maximal, this 2-path must be in conflict with some selected 2-path  $P_{v, v'} \in \mathcal{P}$ . So,  $P_{v, v'}$  covers two vertices among  $\{v, u, v'\}$ . Since the graph is bipartite, and  $P_{v, v'}$  does not use the edges  $\{u, v\}, \{u, v'\} \in F$ , it follows that  $P_{v, v'}$  must cover both vertices  $v$  and  $v'$ . Let  $F(u) \subseteq F$  be the subset of edges of  $F$  which are incident to a vertex of  $N(u)$ . Note that  $|F(u)| \leq k^2$ . Indeed,  $N(u)$  contains  $k$  vertices, all of them having degree at most  $k$  in  $H$ . Let  $P(u) \subseteq E \setminus F$  be the subset of  $E$  containing all edges used by a 2-path  $P_{v, v'} \in \mathcal{P}$  for any pair  $v, v' \in N(u)$ , then  $|P(u)| = 2 \binom{k}{2} = k(k-1)$ . Indeed, there are  $\binom{k}{2}$  pairs  $v, v' \in N(u)$ , and each of them yields a unique 2-path  $P_{v, v'}$  containing two edges (two 2-paths cannot share any edge since they are independent). Now, we define

$$E' := E \setminus (P(u) \cup F(u)) \text{ and } \mathcal{P}' := \mathcal{P} \setminus \{P_{v, v'} \mid v \neq v' \in N(u)\},$$

and claim that  $\mathcal{P}'$  is a maximal set of 2-paths of the subgraph  $G' := (V, E')$ . Assume towards a contradiction that there is a 2-path  $P^*$  in  $G'$  which is not in conflict with any 2-path of  $\mathcal{P}'$ . Since  $G'$  is a subgraph of  $G$ , and  $\mathcal{P}$  is maximal, it follows that  $P^*$  is in conflict with some  $P_{v, v'} \in \mathcal{P} \setminus \mathcal{P}'$ . Thus,  $P^*$  contains at least one of the two vertices  $v, v'$ , and must contain an edge in  $F(u) \subseteq E \setminus E'$ ; a contradiction. Hence,  $\mathcal{P}'$  is maximal for the bipartite graph  $G'$  and, by induction, we have

$$|\mathcal{P}'| \geq \frac{|E'|}{6} - \frac{|V|}{4}.$$

Putting everything together, we have

$$|\mathcal{P}'| = |\mathcal{P}| - \binom{k}{2} \text{ and } |E'| \geq |E| - k(k-1) - k^2,$$

which yields

$$\begin{aligned}
|\mathcal{P}| &= |\mathcal{P}'| + \binom{k}{2} \geq \frac{|E'|}{6} - \frac{|V|}{4} + \binom{k}{2} \\
&\geq \frac{|E| - k(k-1) - k^2}{6} - \frac{|V|}{4} + \binom{k}{2} \\
&= \frac{|E|}{6} - \frac{|V|}{4} + \frac{2k(k-1) - k^2}{6} \\
&\geq \frac{|E|}{6} - \frac{|V|}{4} \quad (\text{since } 2k(k-1) \geq k^2 \text{ for all } k \geq 2).
\end{aligned}$$

This completes the proof of the induction property.

To prove the lemma, it remains to distinguish between the following two cases:

- If  $m = |E| \geq \frac{9|V|}{2}$ , then with the induction property above, any maximal set of independent 2-paths has size at least

$$\frac{m}{6} - \frac{|V|}{4} \geq \frac{m}{6} - \frac{1}{4} \cdot \frac{2m}{9} = \frac{3m - m}{18} \geq \left\lfloor \frac{m}{9} \right\rfloor.$$

- If  $m = |E| < \frac{9|V|}{2}$ , then with Lemma 3, there exists a set of independent 2-paths with size at least  $\frac{|V|}{2} - 1 \geq \lceil \frac{m}{9} - 1 \rceil = \lfloor \frac{m}{9} \rfloor$ .

□

Combining Lemmas 2 and 4 we obtain that any connected graph has at least  $\lfloor \frac{m}{18} \rfloor$  independent 2-paths. For dense graphs, Lemma 4 gives a lower bound of  $\frac{m}{12} - \frac{|V|}{4} \sim \frac{m}{12}$ . This yields the following result:

**Theorem 2** *Let  $G = (V, E)$  with  $|E| = m$  be a connected graph. The number of independent 2-paths in  $G$  is in  $\Omega(m)$ .*

Next we obtain bounds on the sampling probability such that we can guarantee there are sufficiently many pairwise independent 2-paths in  $G_S$ . As we show later, this is needed to guarantee that `SAMPLERANDOM2PATH` will return an almost unbiased estimator of the transitivity coefficient. The events for two 2-paths being monochromatic are independent, thus the next lemma follows from Theorem 2 and Chebyshev's inequality. Note that we still don't need the coloring function  $f$  to be fully random.

**Lemma 5** *Let  $f$  be 6-wise independent and  $p \geq \frac{5}{\varepsilon\sqrt{b}}$  where  $b$  is the number of independent 2-paths in  $G$  and  $\varepsilon \in (0, 1]$ . Then with probability at least  $3/4$  `SPARSIFYGRAPH` returns  $G_S$  such that there are at least  $18/\varepsilon^2$  independent 2-paths in  $G_S$ .*

**Proof:** Let  $D_2$  be a set of pairwise independent 2-paths in  $G$ ,  $|D_2| = b$ . Clearly, each 2-path in  $D_2$  is monochromatic with probability  $p^2$ . Consider two pairwise independent 2-paths  $(u_1, v_1, w_1)$  and  $(u_2, v_2, w_2)$ . If they do not share a vertex, then since  $f$  is 6-wise independent, the two of them are monochromatic with probability  $p^4$ . Otherwise, assume w.l.o.g.  $u_1 = u_2$ . Under the assumption that  $f(u_1) = c$  for some  $c \in C$ , we have that evaluating  $f$  on  $v_i, w_i$ ,  $i = 1, 2$ , is 5-wise independent. Thus,  $\Pr[f(v_i) = f(w_i) = c] = p^4$ . The events of sampling the two 2-paths are thus independent. We introduce an indicator random variable  $X_{(u,v,w)}$  for each 2-path  $(u, v, w) \in D_2$  denoting whether  $(u, v, w) \in G_S$  and set  $X = \sum_{(u,v,w) \in D_2} X_{(u,v,w)}$ . We have  $\mathbb{E}[X] \geq p^2 b = 25/\varepsilon^2$  and since the events that any two 2-paths in  $D_2$  are sampled are independent, we have  $V[X] \leq \mathbb{E}[X]$ . From Chebyshev's inequality with some algebra we then obtain that with probability at least  $3/4$  we have at least  $18/\varepsilon^2$  monochromatic pairwise independent 2-paths in  $G_S$ . □

**Lemma 6** *Assume we run `ESTIMATENUMBEROFTRIANGLES` with  $s = 18/\varepsilon^2$  and let  $X$  be the value returned by `SAMPLERANDOM2PATH`. Then  $(1 - \varepsilon)\alpha \leq \mathbb{E}[X] \leq (1 + \varepsilon)\alpha$ .*

**Proof:** We analyze how much differs the probability between 2-paths to be selected by SAMPLERANDOM2PATH. Consider a given 2-path  $(u, v, w)$ . It will be sampled if the following three events occur:

1.  $(u, v, w)$  is monochromatic, i.e., it is in the sparsified graph  $G_S$ .
2. There are  $i \geq 18/\varepsilon^2$  pairwise independent 2-paths in  $G_S$ .
3.  $(u, v, w)$  is selected by SAMPLERANDOM2PATH.

The first event occurs with probability  $p^2$ . Since  $f$  is fully random, the condition that  $(u, v, w)$  is monochromatic does not alter the probability for any 2-path independent from  $(u, v, w)$  to be also monochromatic. The probability to be in  $G_S$  changes only for 2-paths containing two vertices from  $\{u, v, w\}$ , which in turn changes the number of 2-paths in  $G_S$  and thus probability for  $(u, v, w)$  to be picked up by SAMPLERANDOM2PATH. In the following we denote by  $p_{G_S}$  the probability that a given 2-path is monochromatic and there are at least  $18/\varepsilon^2$  pairwise independent 2-paths in  $G_S$ , note that  $p_{G_S}$  is equal for all 2-paths.

Consider a fixed coloring to  $V \setminus \{u, v, w\}$ . We analyze the difference in the number of monochromatic 2-paths depending whether  $f(u) = f(v) = f(w)$  or not. There are two types of 2-paths that can become monochromatic conditioning on  $f(u) = f(v) = f(w)$ : either (i) 2-paths with two endpoints in  $\{u, v, w\}$  centered at some  $\{u, v, w\}$ , or (ii) 2-paths with two vertices in  $\{u, v, w\}$  centered at a vertex  $x \notin \{u, v, w\}$ . For the first case assume w.l.o.g. there is a 2-path  $(u, v, w) \in G_S$  centered at  $v$  and let  $d_v^{f(v)} = |\{z \in N(v) \setminus \{u, w\} : f(z) = f(v)\}|$ , i.e.,  $d_v^{f(v)}$  is the number of  $v$ 's neighbors different from  $u$  and  $w$ , having the same color as  $v$ . Thus, the number of monochromatic 2-paths centered at  $v$  varies by  $2d_v^{f(v)}$  conditioning on the assumption that  $f(u) = f(v) = f(w)$ . The same reasoning applies also to the 2-paths centered at  $u$  and  $w$ . For the second case consider the vertices  $u$  and  $v$ . Conditioning on  $f(u) = f(w)$ , we additionally add to  $G_S$  2-paths  $(u, x_i, w)$  for which  $f(x_i) = f(u) = f(w)$  and  $x_i \in N(u) \cap N(w)$ . The number of such 2-paths is at most  $\min(d_u^{f(u)}, d_w^{f(w)})$ . The same reasoning applies to any pair of vertices from  $\{u, v, w\}$ . Therefore, depending on whether  $f(u) = f(v) = f(w)$  or not, the number of monochromatic 2-paths centered at a vertex from  $\{u, v, w\}$  varies between

$$\sum_{y \in \{u, v, w\}} \binom{d_y^{f(y)}}{2} \quad \text{and} \quad \sum_{y \in \{u, v, w\}} \binom{d_y^{f(y)}}{2} + 3d_y^{f(y)}.$$

Set  $k = 18/\varepsilon^2$ . Consider now two different, but not necessarily independent, 2-paths  $(u_1, v_1, w_1), (u_2, v_2, w_2) \in G$ . We analyze the probability for each of them to be selected by SAMPLERANDOM2PATH. Let  $\mathcal{C}$  be a partial coloring to  $V \setminus \{u_j, v_j, w_j\}, j = 1, 2$ . If  $\mathcal{C}$  is completed to a coloring of all vertices such that both  $(u_1, v_1, w_1)$  and  $(u_2, v_2, w_2)$  are monochromatic, then clearly they are picked up with the same probability. Assume that with probability  $p_i$ ,  $i - 1$  2-paths are colored monochromatic by  $\mathcal{C}$  and consider extensions of  $\mathcal{C}$  that make exactly one of  $(u_1, v_1, w_1)$  and  $(u_2, v_2, w_2)$  monochromatic. Under the assumption there are  $i \geq k$  2-paths in  $G_S$  and following the above discussion about the number of 2-paths with at least two vertices from  $\{u_j, v_j, w_j\}$ , we see that the number of monochromatic 2-paths can vary between  $i$  and  $i + 3\sqrt{2i}$ . Thus, the probability for  $(u_1, v_1, w_1)$  and  $(u_2, v_2, w_2)$  to be sampled varies between

$$p_{G_S} \sum_{i \geq k} \frac{p_i}{i} \quad \text{and} \quad p_{G_S} \sum_{i \geq k} \frac{p_i}{i + 3\sqrt{2i}}.$$

We assume  $G_S$  contains at least  $k$  2-paths, thus  $\sum_{i \geq k} p_i = 1$  and there exists  $r \geq k, r \in \mathbb{R}$  such that  $\sum_{i \geq k} p_i i^{-1} = 1/r$ . Thus we bound

$$\sum_{i \geq k} \frac{p_i}{i + 3\sqrt{2i}} = \sum_{i \geq k} \frac{p_i}{i(1 + 3\sqrt{2/i})} \geq \frac{1}{1 + 3\sqrt{2/k}} \sum_{i \geq k} \frac{p_i}{i} = \frac{1}{r(1 + 3\sqrt{2/k})}.$$

Since the function  $f$  is fully random, each coloring is equally probable. The above reasoning applies to any pair of 2-paths in  $G$ , thus for any 2-path the probability to be sampled varies between

$$\frac{p_{G_S}}{r} \quad \text{and} \quad \frac{p_{G_S}}{(1 + \sqrt{18/k})r} = \frac{p_{G_S}}{(1 + \varepsilon)r}.$$

Assume first the extreme case that 2-paths which are not part of a triangle are sampled with probability  $\frac{1}{r}$  and 2-paths part of a triangle with probability  $\frac{1}{(1+\varepsilon)r}$ . We have  $X = \sum_{(u,v,w) \in P_2} I_{(u,v,w)}$ , where  $I_{(u,v,w)}$  is an indicator random variable denoting whether  $(u, v, w)$  is part of a triangle. Thus

$$\mathbb{E}[X] \geq \frac{p_{G_S} 3T_3}{(1 + \varepsilon)r p_{G_S} P_2} r = \frac{\alpha}{1 + \varepsilon} \geq (1 - \varepsilon)\alpha.$$

On the other extreme, assuming that we select 2-paths part of triangles with probability  $\frac{1}{r}$  and 2-paths not part of a triangle with probability  $\frac{1}{r(1+\varepsilon)}$ , using similar reasoning we obtain  $\mathbb{E}[X] \leq (1 + \varepsilon)\alpha$ .  $\square$

Applying a variation of rejection sampling, in the next lemma we show how to store a sparsified graph  $G_S$  such that we efficiently sample a 2-path uniformly at random and  $G_S$  is updated in constant time.

**Lemma 7** *Let  $G_S = (V, E_S)$  be a sparsified graph over  $m'$  monochromatic edges. There exists an implementation of  $G_S$  in space  $O(m')$  such that an edge can be inserted to or deleted from  $G_S$  in constant time with high probability. A random 2-path, if existent, can be selected from  $G_S$  in expected time  $O(\log n)$  and  $O(\log^2 n)$  time with high probability.*

**Proof:** We implement  $G_S$  as follows. Let  $\Delta \leq n$  be the maximum vertex degree in  $G_S$ . We maintain a hash table  $H_i$  for all vertices  $v \in G_S$  with  $P_2(v) \in \{2^i, 2^i + 1, \dots, 2^{i+1} - 1\}$ ,  $0 \leq i \leq \log \Delta$ , i.e., there are between  $2^i$  and  $2^{i+1} - 1$  2-paths centered at  $v$ . The hash table contains a vertex  $v$  as a key together with the set of its neighbors  $N(v)$ . We also store the vertices with only one neighbor in  $G_S$  in a hashtable  $H_0$ . (Note that if there are no vertices in a given  $H_i$ , we don't maintain  $H_i$ .)

In another hash table  $T$  we maintain for each vertex incident to at least one sampled edge, a link to the  $H_i$  it is contained in. Whenever a new edge  $(u, v)$  is inserted or deleted from  $G_S$ , we first look-up in  $T$  for the  $H_i$  containing  $u$  and  $v$  and then update the corresponding numbers of 2-paths centered at  $u$  and  $v$ . It may happen that we need to move  $u$  and/or  $v$  from a hashtable  $H_i$  to a hashtable  $H_{i\pm 1}$ . For each  $H_i$  we also maintain the total number of 2-paths centered at vertices  $v \in H_i$ , denote this number as  $P_2(H_i)$ . Implementing  $T$  and the  $H_i$  using the implementation from [5], each update takes constant time with high probability and the total space is  $O(m')$ .

We sample a 2-path from  $G_S$  at random as follows. We compute  $P_2(G_S) = \sum_i H_i$  and select an  $H_i$  where each  $H_i$  has a chance of being picked up of  $P_2(H_i)/P_2(G_S)$ . This is done by generating a random number  $r \in (0, 1]$  and then computing a prefix-sum of  $P_2(H_i)$ 's until the sum reaches  $r$  in time  $O(\log n)$ . Once we have chosen an  $H_i$ , we select a vertex  $v \in H_i$  at random as follows. Assume we maintain the set of vertices in each  $H_i$  in a dynamic dictionary  $V_i$  implemented as a hashtable using tabulation hashing. We assume that the longest chain in  $V_i$  is bounded by a  $\kappa = O(\log n / \log \log n)$  [24]. We select a chain in  $V_i$  at random and keep it with probability  $\ell/\kappa$ , where  $\ell \geq 0$  is the length of the selected chain, otherwise we reject it and repeat the step until a chain is kept. Then, we select at random one of the vertices on the chain, let this be  $v$ . We apply one more time rejection sampling in order to decide whether we keep  $v$  or not: Let  $q = 2^{i+1} - 1$ . We get the value  $d_v$  from  $H_i$  and keep  $v$  with probability  $\binom{d_v}{2}/q$  and reject it with probability  $1 - \binom{d_v}{2}/q$ . Once we keep a vertex  $v$ , we choose at random two of its neighbors in  $G_S$  which constitutes the sampled 2-path. The expected number of sampling a chain and a random vertex until we keep a vertex is  $O(\log n / \log \log n)$ , thus the expected number of trials is  $O(\log n)$  and with high probability we determine a 2-path in time  $O(\log^2 n)$ . It is easy to see that each 2-path is selected with equal probability  $p$  such that  $\frac{1}{2\kappa P_2(G_S)} \leq p \leq \frac{1}{P_2(G_S)}$ , thus we sample uniformly at random from the set  $P_2(G_S)$ .  $\square$

Now we have all components in order to prove the main result.

**Proof:** (of Theorem 1).

Assume ESTIMATENUMBEROFTRIANGLES runs  $K$  copies in parallel of SPARSIFYGRAPH with  $p = \frac{5}{\varepsilon\sqrt{b}}$  for  $b = \lfloor m/18 \rfloor$ . By Lemma 2, Lemma 4 and Lemma 5 with probability  $3/4$  we have a sparsified graph with at least  $s = 18/\varepsilon^2$  pairwise independent 2-paths. Thus, we expect to obtain from  $3K/4$  of them an indicator random variable. A standard application of Chernoff's inequality yields that with probability  $O(2^{-K/36})$  we will have  $\ell \geq K/2$  indicator random variables  $X_i$  denoting whether the sampled 2-path is part of a triangle. By Lemma 6 we have  $(1 - \varepsilon)\alpha \leq \mathbb{E}[X_i] \leq (1 + \varepsilon)\alpha$  and as an estimate of  $\alpha$  we return  $\sum_{i=1}^{\ell} X_i/\ell$ . Observe that  $(1 + \varepsilon/3)^2 \leq 1 + \varepsilon$ , respectively  $(1 - \varepsilon/3)^2 \geq 1 - \varepsilon$ . From the above discussion and applying Chernoff's inequality and the union bound, we see that for  $K = \frac{36}{\varepsilon^2\alpha} \log \frac{2}{\delta}$ , we obtain an  $(\varepsilon, \delta/2)$ -approximation of  $\alpha$ .

By Lemma 1 we can compute an  $(\varepsilon, \delta/2)$ -approximation of the number of 2-paths in space  $O(\frac{1}{\varepsilon^2} \log \frac{1}{\delta})$  and  $O(\log \frac{1}{\delta})$  per edge processing time. It is trivial to show that this implies an  $(3\varepsilon, \delta)$ -approximation of the number of triangles for  $\varepsilon < 1/3$ . Clearly, one can rescale  $\varepsilon$  in the above, i.e.  $\varepsilon = \varepsilon/3$ , such that ESTIMATENUMBEROFTRIANGLES returns an  $(\varepsilon, \delta)$ -approximation.

By Lemma 7, each sparsified graph with  $m'$  edges uses space  $O(m')$  and each update takes constant time with high probability, thus we obtain that each edge is processed with high probability in time  $O(K)$ . Each monochromatic edge and its color can be represented in  $O(\log n)$  bits.

By Lemma 7, in expected time  $O(\log n)$  and worst case time  $O(\log^2 n)$  with high probability we sample uniformly at random a 2-path from each  $G_S$  with at least  $18/\varepsilon^2$  pairwise independent 2-paths.  $\square$

## 5 Lower bound

Pavan et al. [22] show that every streaming algorithm approximating the number of triangles in adjacency streams (edges are inserted in arbitrary order) needs space  $\omega(1/\alpha(G))$ . For the general case, we show another lower bound on the memory needed which matches the upper bound by Manjunath et al. [18]. Our lower bound works for the *promise version* of this problem where the algorithm is required to distinguish between the case where there are no triangles and the case where there are at least  $T_3$  triangles, for a parameter  $T_3$ . The behavior if the number of triangles is between 1 and  $T_3 - 1$  is unspecified, i.e., the algorithm may return any result. Clearly this problem is solved as a special case by any streaming algorithm that is able to approximate the number of triangles, so our space lower bound will also apply to the setting of our upper bound.

**Theorem 3** *Let  $G = (V, E)$  be a graph over  $m$  vertices. Any one-pass streaming algorithm distinguishing between the cases where  $G$  has at least  $T_3$  triangles and  $G$  is triangle-free, needs  $\Omega(m^3/T_3^2)$  bits.*

**Proof:** We obtain the lower bound by a reduction from 1-way protocols for the INDEX problem in communication complexity. In this problem Alice is given a bit string  $x \in \{0, 1\}^a$  and needs to send a message to Bob who holds an index  $i$ , such that Bob is able to output  $x_i$  with probability at least  $2/3$  (where the probability is over random coin tosses made by Alice and Bob). It is known that this problem requires a message of size  $\Omega(a)$  [16].

From a streaming algorithm with space usage  $s$  that distinguishes the cases of 0 and  $\geq T_3$  triangles we obtain a communication protocol for the indexing problem with  $a = \Theta(m^3/T_3^2)$ . The reduction is as follows: Consider  $a$  vertex disjoint bicliques  $C_1, \dots, C_a$  with  $\Theta(T_3/m)$  vertices on each side, and form the stream with edge set  $\cup_{i,x_i=1} C_i$ . The number of edges in this stream is  $\Theta(a(T_3/m)^2)$ , which is  $\Theta(m)$ . Adjusting constants we can make the number of edges less than  $m/2$ . Alice can simulate the streaming algorithm on this input and send its state of  $s$  bits to Bob. In order to determine the value of  $x_i$  Bob now connects  $\Theta(m^2/T_3)$  new (i.e., previously isolated) vertices to all vertices in  $C_i$ . Adjusting constants this gives another  $m/2$  edges, and will create either no triangles (if  $C_i$  was not in the stream simulated by Alice) or  $\Theta((m^2/T_3)(T_3/m)^2) = \Theta(T_3)$  triangles (if  $C_i$  was in the stream). Thus, if the streaming algorithm succeeds with probability  $2/3$ , so does Bob. In conclusion we must have a space usage of  $s = \Omega(a)$  bits, which is  $\Omega(m^3/T_3^2)$  as desired.  $\square$

Dataset	$n$	$m$	$T_3$	$\alpha$	$m^3/T_3^2$	$mn/T_3$	$\sqrt{m}/\alpha$
Enron	36K	367K	727K	0.0853	93.5K	18K	<b>7.1K</b>
AS20000102	6.4K	13.2K	6.5K	0.0095	52K	12.9K	<b>12.1K</b>
Astro-Ph	18.7K	396.1K	1.35M	0.318	34K	5.5K	<b>2K</b>
Cond-Mat	23.1K	186.9K	173.3K	0.2643	215K	24.9K	<b>1.6K</b>
roadNet-CA	1.96M	5.53M	120.6K	0.0603	11.55M	89.8M	<b>39K</b>
DBLP	317K	1.05M	2.2M	0.3064	239K	148K	<b>3.4K</b>
Oregon	10.6K	22K	17.1K	0.0093	37.1K	<b>13.6K</b>	16K
Facebook	4K	88.2K	1.61M	0.2647	264	<b>218</b>	1.1K
LiveJournal	4M	34.6M	177.8M	0.1154	1.3M	778K	<b>51K</b>
Youtube	1.1M	2.9M	3M	0.0062	2.7M	1.06M	<b>275K</b>
Amazon	334K	928K	667K	0.2052	1.79M	464K	<b>4.7K</b>
Skitter	1.7M	11.1M	28.7M	0.0053	1.66M	657K	<b>629K</b>

Table 4: Comparison of the theoretical guarantees for several real-life graphs,  $b = \max(n, P_2/n)$ . The lowest space usage in each row is given in bold font.

## 6 Comparison on the theoretical guarantee for real graphs

In Table 4 we compare the theoretical guarantee on the complexity of our algorithm to the ones shown in [1, 18] for real graphs. We used the publicly available information from the Stanford Large Network Dataset Collection<sup>2</sup>. (Note that there the transitivity coefficient is called *fraction of closed triangles*.) As can be seen, the transitivity coefficient  $\alpha$  appears to be constant. The update time of  $O(n \log n)$  from [1] is always impractical, and for [18] it is almost always prohibitively large, the only exception being the Facebook graph for which the  $m^3/T_3^2$  is small. (Note that this is a graph collected from users who used the Social Circles application and does not reflect the structure of the whole social network graph.) The space usage of our algorithm is also much better for several graphs which are not very dense with respect to the number triangles.

## References

- [1] K. J. Ahn, S. Guha, A. McGregor. Graph sketches: sparsification, spanners, and subgraphs. *PODS 2012*: 5–14
- [2] W. Aiello, F. R. K. Chung, L. Lu. A random graph model for massive graphs. *STOC 2000*: 171–180
- [3] R. Albert, A.-L. Barabási. Statistical mechanics of complex networks. *Rev. Mod. Phys.* 74, 47–97 (2002)
- [4] N. Alon, R. Yuster, U. Zwick. Finding and Counting Given Length Cycles. *Algorithmica* 17(3): 209–223 (1997)
- [5] Y. Arbitman, M. Naor, G. Segev. Backyard Cuckoo Hashing: Constant Worst-Case Operations with a Succinct Representation. *FOCS 2010*: 787–796
- [6] L. Backstrom, P. Boldi, M. Rosa, J. Ugander, S. Vigna. Four degrees of separation. *WebSci 2012*: 33–42
- [7] J.W. Berry, B. Hendrickson, R. LaViolette, C.A. Phillips. Tolerating the Community Detection Resolution Limit with Edge Weighting. *Phys. Rev. E*, 83(5)
- [8] L. S. Buriol, G. Frahling, S. Leonardi, A. Marchetti-Spaccamela, C. Sohler. Counting triangles in data streams. *PODS 2006*: 253–262

<sup>2</sup><http://snap.stanford.edu/data/index.html>

- [9] L. Carter, M. N. Wegman. Universal Classes of Hash Functions. *J. Comput. Syst. Sci.* 18(2): 143–154 (1979)
- [10] G. Frahling, P. Indyk, C. Sohler. Sampling in dynamic data streams and applications. *Symposium on Computational Geometry 2005*: 142–149
- [11] H. Jowhari, M. Ghodsi. New Streaming Algorithms for Counting Triangles in Graphs. *COCOON 2005*: 710–716
- [12] M. Jha, C. Seshadhri, A. Pinar. A space efficient streaming algorithm for triangle counting using the birthday paradox. *KDD 2013*: 589–597
- [13] H. Jowhari, M. Saglam, G. Tardos. Tight bounds for Lp samplers, finding duplicates in streams, and related problems. *PODS 2011*: 49–58
- [14] D. M. Kane, K. Mehlhorn, T. Sauerwald, H. Sun. Counting Arbitrary Subgraphs in Data Streams. *ICALP (2) 2012*: 598–609
- [15] M. N. Kolountzakis, G. L. Miller, R. Peng, C. E. Tsourakakis. Efficient Triangle Counting in Large Graphs via Degree-based Vertex Partitioning. *Internet Mathematics* 8(1-2), 161–185 (2012)
- [16] I. Kremer, N. Nisan, D. Ron. On Randomized One-Round Communication Complexity. *Computational Complexity* 8(1): 21–49 (1999) 1995
- [17] S. Leonardi. List of Open Problems in Sublinear Algorithms: Problem 11. <http://sublinear.info/11>
- [18] M. Manjunath, K. Mehlhorn, K. Panagiotou, H. Sun. Approximate Counting of Cycles in Streams. *ESA 2011*: 677–688
- [19] S. Muthukrishnan. Data Streams: Algorithms and Applications. *Foundations and Trends in Theoretical Computer Science*, Vol. 1, Issue 2, 2005
- [20] A. Pagh, R. Pagh. Uniform Hashing in Constant Time and Optimal Space. *SIAM J. Comput.* 38(1): 85–96 (2008)
- [21] R. Pagh, C. E. Tsourakakis. Colorful triangle counting and a MapReduce implementation. *Inf. Process. Lett.* 112(7): 277–281 (2012)
- [22] A. Pavan, K. Tangwongsan, S. Tirthapura, K.-L. Wu. Counting and Sampling Triangles from a Graph Stream. *PVLDB* 6(14): 1870–1881 (2013)
- [23] C. Seshadhri, A. Pinar, and T. Kolda. Triadic Measures on Graphs: The Power of Wedge Sampling. *SDM 2013*, to appear.
- [24] M. Pătraşcu, M. Thorup. The Power of Simple Tabulation Hashing. *J. ACM* 59(3): 14 (2012)
- [25] M. Thorup, Y. Zhang. Tabulation based 4-universal hashing with applications to second moment estimation. *SODA 2004*: 615–624
- [26] C. E. Tsourakakis, U. Kang, G. L. Miller, C. Faloutsos. DOULION: counting triangles in massive graphs with a coin. *KDD 2009*: 837–846
- [27] C. E. Tsourakakis, M. N. Kolountzakis, G. L. Miller. Triangle Sparsifiers. *J. of Graph Algorithms and Appl.* 15(6): 703–726 (2011)
- [28] V. Vassilevska Williams. Multiplying matrices faster than Coppersmith-Winograd. *STOC 2012*, 887–898