



Dynamic Graph Stream Algorithms in $o(n)$ Space

Zengfeng Huang¹ · Pan Peng² 

Received: 26 October 2016 / Accepted: 15 September 2018 / Published online: 25 September 2018
© The Author(s) 2018

Abstract

In this paper we study graph problems in the dynamic streaming model, where the input is defined by a sequence of edge insertions and deletions. As many natural problems require $\Omega(n)$ space, where n is the number of vertices, existing works mainly focused on designing $O(n \cdot \text{poly log } n)$ space algorithms. Although sublinear in the number of edges for dense graphs, it could still be too large for many applications (e.g., n is huge or the graph is sparse). In this work, we give single-pass algorithms beating this space barrier for two classes of problems. We present $o(n)$ space algorithms for estimating the number of connected components with additive error εn of a general graph and $(1 + \varepsilon)$ -approximating the weight of the minimum spanning tree of a connected graph with bounded edge weights, for any small constant $\varepsilon > 0$. The latter improves upon the previous $O(n \cdot \text{poly log } n)$ space algorithm given by Ahn et al. (SODA 2012) for the same class of graphs. We initiate the study of approximate graph property testing in the dynamic streaming model, where we want to distinguish graphs satisfying the property from graphs that are ε -far from having the property. We consider the problem of testing k -edge connectivity, k -vertex connectivity, cycle-freeness and bipartiteness (of planar graphs), for which, we provide algorithms using roughly $O(n^{1-\varepsilon} \cdot \text{poly log } n)$ space, which is $o(n)$ for any constant ε . To complement our algorithms, we present $\Omega(n^{1-O(\varepsilon)})$ space lower bounds for these problems, which show that such a dependence on ε is necessary.

Keywords Dynamic graph streams · Graph sketching · Property testing · Minimum spanning tree

A preliminary version of this paper appeared in the Proceedings of the 43rd International Colloquium on Automata, Languages and Programming (ICALP), Rome, Italy, July 2016.

Zengfeng Huang: Supported by Shanghai Science and Technology Commission (Grant No. 17JC1420200) and Australian Research Council Discovery Grant (DP150102728).

Pan Peng: Work done while the author was at TU Dortmund, Germany and supported by ERC grant No. 307696.

Extended author information available on the last page of the article

1 Introduction

Graphs or networks are a natural way to describe structural information. For example, users of Facebook and the acquaintance relations among them form a social network, the proteins together with interactions between them define a biological network, and web-pages and hyperlinks give rise to a huge web graph. Due to the rapid development of information technology, many such graphs become extremely large, and are constantly changing, which poses great challenges for analyzing their structures. Over the last decade, the data stream model [34] has proven to be successful in dealing with *big data*. In this model, the algorithm should make only one pass (or a few passes) over the stream, and use sublinear working space. The time required to output the final answer and process each element is also important. There is a growing body of work studying graph problems over data streams. Graph streams were first considered by Henzinger et al. [24], and later have been extensively studied in the *insertion-only* model (e.g., [17,18,34]), where there is no edge deletion in the stream. Recently, starting from the seminal works of Ahn, Guha and McGregor [2,3], the interest has shifted to the *dynamic streaming model*, where the edges can be both inserted and deleted (see e.g., [1,5–7,9,10,14,23,28,29,31,33]). In this setting, most algorithms designed are linear sketch-based, which is also an effective technique for processing distributed graphs. For more information about graph streaming algorithms see the recent survey by McGregor [32].

For graph streams, both insertion-only and dynamic, the research in the past has mostly focused on the *semi-streaming model*, in which the algorithms are allowed to use $\tilde{O}(n)$ space, where n is the number vertices in the graph. (For notational convenience, we will use $\tilde{O}(g)$ and $\tilde{\Omega}(g)$ to hide $\text{poly}(\log(g))$ factors.) The reason behind this is that even in the insertion-only model, many natural graph problems require $\Omega(n)$ space (e.g., testing if the graph is connected [18]). Note that the allowed space in semi-streaming model is sublinear in the input size as the number of edges of the graph might be as large as $\Omega(n^2)$. However, in many real applications (e.g., the input graph is already very sparse), an $\tilde{O}(n)$ space algorithm might be even worse than just storing all the edges. From this perspective, one may naturally ask the question *which kind of problems can be solved with even less space, i.e., $o(n)$ space*.

To the best of our knowledge, very few results are known in this direction. Chitnis et al. [10] and Fafianie and Kratsch [16] introduced parameterized graph stream algorithms which may only use $o(n)$ space with some promise of the size of the solution. This parameterized setting has been further investigated in [9]. In addition, it has been shown that the size of the maximum matching can be approximated within constant factor in $\tilde{O}(n^{4/5})$ space for graphs with bounded arboricity [7,9,15].

In this paper, we study two classes of graph problems that admit single-pass $o(n)$ space algorithms in the dynamic streaming model. The first class contains the problems of estimating the number of connected components and the weight of the minimum spanning tree (MST). We show that one can estimate the number of connected components within an *additive* error of εn with $o(n)$ space and post-processing time, for any constant $\varepsilon > 0$. We also present an algorithm to $(1 + \varepsilon)$ -approximate the weight of the MST with $o(n)$ space and post-processing time for connected graphs with bounded edge weights, which improves the best known algorithm with $\tilde{O}(n)$ space in the same

setting given by Ahn et al. [2]. It is worth noting that the problem of estimating the number of connected components within small *multiplicative* error requires $\Omega(n)$ space, as it is generally harder than the problem of (exactly) testing graph connectivity; and that estimating the weight of MST for graphs with arbitrarily large edge weights (e.g., $\Omega(\log n)$) requires $\Omega(n)$ space (see Theorem 10). Previously these two problems have been studied in the framework of *sublinear time algorithms* (see eg. [8,40]).

The second class consists of problems that are relaxations of deciding graph properties. Given a huge graph, it is very useful to know whether the graph has some predetermined property, such as k -connectivity, bipartiteness, cycle-freeness and etc., which provide valuable information about the graph. However, besides the requirement of $\Omega(n)$ space, exactly testing properties sometimes is too strong a requirement for analyzing highly dynamic graphs, since the answer may change in the next second due to an insertion or deletion of a single edge. In this paper, we initiate the study of *approximate graph property testing* in the dynamic streaming model: we want to test whether a graph satisfies some property or one has to modify a small constant fraction of edges to make it have the property. This notion of approximation is adapted from the framework of *property testing* [21,22,36], and a large number of existing literatures have given efficient testing algorithms (called *testers*) for many properties under different query models (see surveys [20,39]). We show that some fundamental properties can be tested in both $o(n)$ space and post-processing time in the dynamic streaming model and we also present close lower bounds for these problems which hold even in the insertion-only model. We remark that McGregor [32] also suggested to study the (approximate) property testers in graph streaming model, and asked whether more space-efficient algorithms exist for these problems, and we thus give an affirmative answer to this question.

1.1 Our Results

Now we formally state our main results. Our results regarding estimating the number of connected components and the MST weight are as follows.

- *Estimating the Number of Connected Components* We present a dynamic streaming algorithm that estimates the number of connected components within additive error εn in $\tilde{O}(n^{1-\varepsilon+\varepsilon^{q+1}})$ space and post-processing time for any constant $q \geq 1$. We note that a lower bound of $\Omega(n^{1-O(\varepsilon)})$ for this problem follows from the work [42].
- *Estimating the Weight of the Minimum Spanning Tree (MST)* In this problem, we want to estimate the weight of the MST of a connected graph with edge weights in the set $\{1, 2, \dots, W\}$. We give a dynamic streaming algorithm that computes a $(1+\varepsilon)$ -approximation of the MST weight and uses space and post-processing time $\tilde{O}(Wn^{1-\frac{\varepsilon}{W-1}+\frac{\varepsilon^t}{(W-1)^t}})$ for any constant $t \geq 1$. By an argument in [8], the result can be extended to non-integral weights, as long as the ratio between the largest and the smallest weight is bounded. A space lower bound of $\Omega(n^{1-\frac{4\varepsilon}{W-1}})$ is shown for this problem.

Table 1 Upper and lower bounds of streaming testers

	Space \tilde{O}	Space lower bound Ω
Connectivity	$n^{1-\varepsilon}$	$n^{1-8\varepsilon}$
k -Edge connectivity	$k^{1+\varepsilon} \cdot n^{1-\varepsilon}$	
k -Vertex connectivity	$\frac{k^{1+\varepsilon/4}}{\varepsilon} \cdot n^{1-\varepsilon/4}$	
Cycle-freeness	$n^{1-\varepsilon+\varepsilon^2}$	$n^{1-8\varepsilon}$
Bipartiteness of planar graphs	$n^{1-\Omega(\varepsilon^2)}$	$n^{1-4\varepsilon}$

We also present approximate testing algorithms for a number of fundamental graph properties. Before stating the performance of these algorithms, we first introduce some definitions. Given a graph property Π , an m -edge graph G is called ε -far from having Π if one has to modify more than εm edges of G to get a graph G' satisfying Π . This distance definition is adapted from [36] and is most suitable for general graphs where neither edge density nor maximum degree is restricted. We call an algorithm a (*dynamic*) *streaming tester* for Π , if it makes a single-pass over a stream of edge insertions and deletions, with probability at least $2/3$, accepts any graph satisfying Π , and rejects any graph that is ε -far from having Π .

We give sketch-based streaming testers for properties of being connected, k -edge connected, k -vertex connected, cycle-freeness and bipartite (for planar graphs). The performance of our testers are summarized in Table 1. We stress that most of our testers have (asymptotically) the same post-processing time as the space they used except for testing k -edge connectivity when $k \geq \Omega(n^{\varepsilon/(1+\varepsilon)})$ and k -vertex connectivity when $k \geq \Omega(n^{\varepsilon/(4+\varepsilon)})$.

1.2 Our Techniques

To estimate the number of connected components with small additive error εn , we note that it is sufficient to estimate the number $\text{scc}(G)$ of connected components of small size (i.e., $O(1/\varepsilon)$), since the number of components of size larger than this is at most $O(\varepsilon n)$ (see also [8]). To estimate $\text{scc}(G)$, the following vertex sampling framework is used: we sample a sufficiently large set of vertices S by sampling each vertex in G with some probability p , and then use the statistics of the sampled connected components of the original graph to estimate $\text{scc}(G)$. For any small connected component C in G , it is likely that all the vertices in C will be sampled out. Conditioned on this, we add $1/p^{|C|}$ to our final estimator, which is the reciprocal of the probability that C is entirely sampled out. Now the task is then to identify which subsets of S are connected components in the original graph. A trivial way is to check all subsets of S , which takes too much time. A more efficient way is to only check all the connected components in $G[S]$, since a sampled component of G must also form a component in $G[S]$. We carefully use a set of linear sketches to do this. More specifically, we first recover all connected components in $G[S]$ by invoking a sketch-based streaming algorithm

given in [2], which only needs space near-linear in $|S|$. Then we use (different) linear sketches to check if any of these components is indeed a connected component of the original graph. We remark that the first set of linear sketches of a vertex v sketch its neighborhood information in $G[S]$, while the second set sketch its neighborhood information in G . Our $o(n)$ space streaming algorithm for $(1 + \varepsilon)$ -approximating the weight of MST follows via a connection between the number of connected components and the weight of MST established in [8].

To give testers for some graph property Π in dynamic streaming model, we start from the observation that if a graph G is far from having Π , then typically, there exist many small disjoint subgraphs, each of which is a witness that the graph G does not satisfy Π . (For example, if Π is connectivity, then there exists at least $\Omega(\varepsilon m)$ connected components of size at most $O(1/\varepsilon)$ in a graph that is ε -far from being connected.) This implies that by sampling a sufficient large set of vertices, with high probability, one of such subgraphs will be entirely sampled. Checking which vertices form a witness of the original graph can then be done by using the aforementioned framework. Different sketches will be used for testing different properties.

To prove lower bounds for our studied problems, we give reductions from *Boolean Hidden Hypermatching (BHH)* problem that was studied in [42]. Our reductions share similarity with the reduction in [42] to the cycle-counting problem and the reductions in [27,30] to the approximate max-cut problem.

1.3 Related Work

Ahn et al. [2] initiated the study of graph sketches, and gave dynamic semi-streaming algorithms for computing a spanning forest (which can be used to count the exact number of connected components), and $(1 + \varepsilon)$ -approximate the weight of MST. They also proposed algorithms to *exactly* testing of a set of properties, including testing connectivity, k -edge connectivity, and bipartiteness. Recently, Guha et al. gave dynamic streaming algorithms for exactly testing of k -vertex connectivity [23]. All these algorithms use $\tilde{O}(n)$ space ($\tilde{O}(kn)$ for k -connectivity). On the other hand, the randomized space lower bounds for these exact testing problems were known to be $\Omega(n)$ in the insertion-only model [17,18]. Recently, Sun and Woodruff improved these lower bounds to $\Omega(n \log n)$ [41]. Verbin and Yu [42] proved a lower bound for cycle-counting, which implied a lower bound of $\Omega(n^{1-O(\varepsilon)})$ for estimating the number of components.

In the *random order* insertion-only model Kapralov et al. [26] gave a one pass streaming algorithm that estimates the maximum matching size with polylogarithmic approximation ratio in polylogarithmic space. Although sublinear in n , the model considered is very different from ours.

Sublinear time algorithms for estimating the number of connected component and the weight of MST were first given by Chazelle et al. [8]. Later these two problems have been further considered in geometric settings [11,13,19]. In particular, Frahling et al. studied the problem of $(1 + \varepsilon)$ -approximating the weight of MST in dynamic geometric data stream [19].

There has been a rich line of work on graph property testing in the query model (see surveys [20,39]) and the goal there is to design fast algorithms that make as few queries as possible. The query models that are mostly related to ours are bounded degree model and general graph model. In particular, our definition of ϵ -far is adapted from the general graph model. Goldreich and Ron [22] initiated the study of property testers in bounded degree graph model, and gave testers for connectivity, k -edge connectivity, 2, 3-vertex connectivity, cycle-freeness, Eulerianity. Testing k -vertex connectivity in bounded degree graphs for arbitrary constant k was given in [43]. These testers have later been generalized to general graph model [35,36]. Testing bipartiteness in planar graphs was studied in [12].

After having submitted the paper, we became aware that Hossein Jowhari [25] has independently studied the problem of estimating the number of connected components and provided similar results as ours, while he did not consider the streaming property testers considered here. Furthermore, subsequent to our work, Peng and Sohler [37] showed that in *random order* streams, approximating the number of connected components with additive error ϵn and $(1 + \epsilon)$ -approximating the weight of the MST of a connected graph with bounded edge weights can be solved in a single-pass and constant space (in terms of words), i.e., the space complexity only depends on ϵ and is independent of the size of the graph.

2 Preliminaries

2.1 Notations

Let $[n] := \{1, \dots, n\}$. We use $V := [n]$ to denote the vertex set of the graph G defined by the stream, and let m denote the number of edges of G . For an undirected graph $G = ([n], E)$ and a vertex $i \in [n]$, we let $\Gamma(i)$ denote all the neighbors of i . For a set $C \subseteq [n]$, let $\Gamma(C)$ denote the set of vertices in $V \setminus C$ that have at least one neighbor in C , that is, $\Gamma(C) = \cup_{i \in C} \Gamma(i) \setminus C$. Let $E(C, V \setminus C)$ denote the set of edges crossing C and $V \setminus C$. We will use $G[C]$ to denote the subgraph induced by C .

For each vertex i , we define two vectors $\Delta^i \in \{-1, 0, 1\}^{\binom{n}{2}}$ and $\Lambda^i \in \{0, 1\}^n$ to encode the neighborhood information of i as follows:

$$\Delta^i_{j,k} = \begin{cases} 1 & \text{if } i = j < k \text{ and } (j, k) \in E \\ -1 & \text{if } j < k = i \text{ and } (j, k) \in E \\ 0 & \text{otherwise} \end{cases} \quad \Lambda^i_j = \begin{cases} 1 & \text{if } j \in \Gamma(i) \text{ or } j = i \\ 0 & \text{otherwise} \end{cases}$$

By simple induction arguments, it is easy to prove that for any vertex set $C \subset V$, the nonzero entries in the vector $\Delta^C := \sum_{i \in C} \Delta^i$ corresponds to the edges between C and its complement $V \setminus C$. The nonzero entries in $\sum_{i \in C} \Lambda^i$ corresponds exactly to vertices in $C \cup \Gamma(C)$.

2.2 Linear Sketches

Linear sketch (or sketch for short) is a powerful tool widely used in the streaming model and other areas. Given a large vector $\mathbf{x} \in \mathbb{R}^n$, we want to construct a small sketch $\mathcal{L}(\mathbf{x})$, from which certain properties of \mathbf{x} can be recovered. We call \mathcal{L} a linear sketch if $\mathcal{L}(\mathbf{x} + \mathbf{y}) = \mathcal{L}(\mathbf{x}) + \mathcal{L}(\mathbf{y})$ for all \mathbf{x}, \mathbf{y} , and this additive property make it trivial to implement linear sketches in the dynamic streaming model. As in the previous works, we will use linear sketches as our main tool.

AGM sketch We will use a dynamic streaming algorithm for constructing a spanning forest of a graph by Ahn, Guha and McGregor [2], which is summarized as follows.

Theorem 1 (AGM sketch [2]) *There exists a single-pass sketch-based dynamic streaming algorithm that uses $O(n \log^3 n)$ space, and recovers a spanning forest of the graph with probability 0.99. The recovery time of the algorithm is $\tilde{O}(n)$, and the update time is $\text{poly } \log n$.*

AMS sketch To check whether the input vector \mathbf{x} is $\mathbf{0}$ or not, one can simply maintain a constant approximation of its *second frequency moment*, that is $F_2(\mathbf{x}) := \sum_i x_i^2$. In particular, by using the classical *AMS sketch* that was introduced by Alon, Matias and Szegedy [4], one can approximate $F_2(\mathbf{x})$ within a multiplicative factor of c using $O(\log(1/\delta) \log n)$ bits of memory with probability at least $1 - \delta$, for any $0 < \delta < 1$ and constant $c > 1$.

Exact k -sparse recovery We call a vector k -sparse if $|\text{supp}(\mathbf{x})| \leq k$. Given a non-zero vector $\mathbf{x} \in \mathbb{R}^n$, the goal here is to recover \mathbf{x} if \mathbf{x} is k -sparse, otherwise outputs **Fail**. We have the following result from [38].

Lemma 1 [38] *There exists an $O(k \log n \log_k \delta^{-1})$ space sketch-based algorithm that takes as input a non-zero vector $\mathbf{x} \in \mathbb{R}^n$, and with probability $1 - \delta$, recovers \mathbf{x} if \mathbf{x} is k -sparse, otherwise outputs **Fail**. The update time is $O(\text{poly } \log n)$ and the recovery time is $O(k \cdot \text{poly } \log n)$.*

3 Estimating the Number of Connected Components and MST Weight

In this section, we present and analyze our algorithms for estimating the number of the connected components in a graph and $(1 + \varepsilon)$ -approximating the weight of the MST.

3.1 Estimating the Number of Connected Components

Our first observation is that, to estimate the number of connected components within additive error εn , we can simply ignore all the large components (see also [8]). In particular, the number of components of size larger than $\Omega(1/\varepsilon)$ is at most $O(\varepsilon n)$. Thus it will be sufficient to estimate the number of components of small size, for which we have the following theorem.

Theorem 2 *For any constant $t \geq 1$, there exists a one-pass dynamic streaming algorithm that uses $\tilde{O}(e^t n^{1-\varepsilon})$ space and post-processing time to estimates the number of*

connected components of size at most $1/\varepsilon$ within an additive error $\varepsilon^t n$. The update time is $O(\text{poly } \log n)$.

By invoking Theorem 2 with parameter $\varepsilon' = (1 - \varepsilon^q)\varepsilon$ and $t = (q + 1)$, we get an estimator for the number of connected components of size smaller than $1/\varepsilon'$ within additive error at most $\varepsilon^{q+1}n$. Since the number of components of size at least $1/\varepsilon'$ is at most $\varepsilon'n = \varepsilon n - \varepsilon^{1+q}n$, the estimator also approximates the total number of connected components within additive error at most εn . The space of the algorithm is $\tilde{O}(e^{q+1}n^{1-\varepsilon+\varepsilon^{q+1}})$, and we have the following result.

Theorem 3 *Let $q \geq 1$ be a constant. There exists a one-pass dynamic streaming algorithm that with constant success probability, estimates the number of connected components of a graph within an additive error εn in $\tilde{O}(e^{q+1}n^{1-\varepsilon+\varepsilon^{q+1}})$ space and post-processing time.*

Now we give the proof of Theorem 2. Recall that the vectors Δ^C encode the information of the number of edges between C and $V \setminus C$.

Proof of Theorem 2 Let $\text{scc}(G)$ denote the number of connected components of size at most $1/\varepsilon$ in G . Our algorithm for estimating $\text{scc}(G)$ is as follows. We first sample each vertex with probability $p := (\varepsilon^{2t}n/16)^{-\varepsilon}$. Let S be the set of sampled vertices. We then use the AGM sketch from Theorem 1 to maintain a spanning forest F of the subgraph induced by S . Then for each component C in F , we test whether C is actually a connected component in G by testing whether the vector $\Delta^C := \sum_{v \in C} \Delta^v$ is $\mathbf{0}$, which can be done by the AMS sketch. If $\Delta^C = \mathbf{0}$, we set $X_C = 1$, otherwise set $X_C = 0$. Our estimator is then defined as $\sum_C \frac{X_C}{|C|}$, where C ranges over all components of F with size at most $\frac{1}{\varepsilon}$. See Algorithm 1 for the details.

Algorithm 1 EstimateNumSCC

- 1: Sample each vertex with probability $p := (\varepsilon^{2t}n/16)^{-\varepsilon}$. If more than $16np$ vertices are sampled, then abort and output **Fail**. Let S denote the set of sampled vertices.
 - 2: Maintain an AGM sketch of $G[S]$ using Theorem 1.
 - 3: For each $v \in S$, maintain an AMS sketch $\text{AMS}(\Delta^v)$, sketching the neighborhood of v in G .
 - 4: **Post-Processing:**
 - 5: Use the AGM sketch to recover a spanning forest F of $G[S]$ using Theorem 1.
 - 6: For each component $C \in F$, estimate $F_2(\Delta^C)$ using the AMS sketch $\text{AMS}(\Delta^C) = \sum_{v \in C} \text{AMS}(\Delta^v)$, and set $X_C = 1$ if $F_2 = 0$, otherwise set $X_C = 0$. For each $1 \leq \ell \leq \frac{1}{\varepsilon}$, let $X_\ell := \sum_{C:|C|=\ell} X_C$.
 - 7: Output $Y := \sum_{\ell \leq \frac{1}{\varepsilon}} \frac{X_\ell}{\ell}$.
-

Note that the algorithm samples at most $16np = O(\varepsilon^{-2t\varepsilon} \cdot n^{1-\varepsilon})$ vertices and we maintained an AGM sketch on $G[S]$ and an AMS sketch for each sampled vertex, which imply that the space complexity of the algorithm is $O(\varepsilon^{-2t\varepsilon} n^{1-\varepsilon} \cdot \text{poly } \log n)$. By simple calculus, for any ε , it holds that $\varepsilon^{-2\varepsilon} \leq e^{2/e} < e$, so the space is at most $\tilde{O}(e^t n^{1-\varepsilon})$. The post-processing time is near linear in the space, and the update time is $O(\text{poly } \log n)$.

Now we prove the correctness of the above algorithm. First we note that the expected number of sampled vertices in Step (1) is np , and thus by Markov inequality, the probability that more than $16np$ vertices are sampled is at most $\frac{1}{16}$. Also note that with probability at least $1 - \frac{1}{16}$, the AGM sketch returns a true spanning forest of $G[S]$. In addition, since the number of components in F is at most n , we will query the AMS sketch at most n times. Thus if we set the error probability of the AMS sketch to be $\delta = \frac{1}{16n}$, then with probability at least $1 - \frac{1}{16n} \cdot n = 1 - \frac{1}{16}$, all invocations of AMS sketches (with $\log^2 n$ bits of space per sketch) for testing if $\Delta^C = \mathbf{0}$ will give the correct answer. Conditioned on this event, X_ℓ defined in Step (6) is exactly the number of connected components B of size ℓ in G such that all vertices in B are sampled out, which is true since for any component $C \in F$, $F_2(\Delta^C) = \mathbf{0}$ if and only if C is a connected component in G .

Let $B_1, \dots, B_{\text{scc}(G)}$ be the connected components of size at most $\frac{1}{\varepsilon}$ of G .

For any integer $\ell \leq \frac{1}{\varepsilon}$, let \mathcal{B}_ℓ denote the set of connected components of size ℓ in G , that is, $\mathcal{B}_\ell = \{B_i : 1 \leq i \leq \text{scc}(G), |B_i| = \ell\}$. Let $b_\ell := |\mathcal{B}_\ell|$. Note that $\text{scc}(G) = \sum_{\ell \leq \frac{1}{\varepsilon}} b_\ell$. For any set B , let Z_B denote the indicator random variable that all the vertices in B have been sampled. Note that $\Pr[Z_B = 1] = p^{|B|}$. Now by the above argument, $X_\ell = \sum_{B \in \mathcal{B}_\ell} Z_B$, and $E[X_\ell] = b_\ell \cdot p^\ell$. Furthermore, we have $Y = \sum_{\ell \leq \frac{1}{\varepsilon}} \frac{X_\ell}{p^\ell} = \sum_{\ell \leq \frac{1}{\varepsilon}} \frac{\sum_{B \in \mathcal{B}_\ell} Z_B}{p^\ell}$, and thus $E[Y] = \sum_{\ell \leq \frac{1}{\varepsilon}} b_\ell = \text{scc}(G)$.

Note that all Z_{B_i} 's are mutually independent for all i , so it holds that

$$\begin{aligned} \text{Var}[Y] &= \sum_{\ell \leq \frac{1}{\varepsilon}} \frac{\sum_{B \in \mathcal{B}_\ell} \text{Var}[Z_B]}{p^{2\ell}} = \sum_{\ell \leq \frac{1}{\varepsilon}} \frac{b_\ell(p^\ell - p^{2\ell})}{p^{2\ell}} \leq \sum_{\ell \leq \frac{1}{\varepsilon}} \frac{b_\ell}{p^\ell} \\ &\leq \frac{\sum_{\ell \leq \frac{1}{\varepsilon}} b_\ell}{p^{1/\varepsilon}} = \frac{\text{scc}(G)}{p^{1/\varepsilon}} \leq \frac{n}{p^{1/\varepsilon}} = \varepsilon^{2t} n^2 / 16, \end{aligned} \tag{1}$$

where we use the fact that $\text{scc}(G) \leq n$, and $p = (\varepsilon^{2t} n / 16)^{-\varepsilon}$. Then by Chebyshev's inequality,

$$\Pr[|Y - \text{scc}(G)| \geq \varepsilon^t n] = \Pr[|Y - E[Y]| \geq \varepsilon^t n] \leq \frac{\text{Var}[Y]}{\varepsilon^{2t} n^2} \leq 1/16.$$

By the union bound, the algorithm will succeed with probability at least $\frac{2}{3}$. □

3.2 Approximating the Weight of Minimum Spanning Tree

We use the previous algorithm on estimating the number of connected components to approximate the weight of a minimum spanning tree of a weighted graph. Let $W \geq 2$ be an integer, G be a connected graph with integer edge weights from $[W] := \{1, \dots, W\}$, and $c(\text{MST})$ be the weight of an MST of G . For any $1 \leq \ell \leq W$, let $G^{(\ell)}$ denote the subgraph of G consisting of all edges of weight at most ℓ . Let $\text{cc}^{(\ell)}$ denote the number

of connected components of $G^{(\ell)}$. Chazelle et al. [8] give the following lemma relating the weight of MST to the number of connected components of $G^{(\ell)}$.

Lemma 2 [8] *It holds that $c(\text{MST}) = n - W + \sum_{\ell=1}^{W-1} cc^{(\ell)}$.*

For a connected graph with integer edge weights, the weight of any MST is at least $n - 1$, so it is sufficient to estimate $cc^{(\ell)}$ within an additive error of $\varepsilon n / (W - 1)$ for each ℓ . To do this, we can simply run $W - 1$ parallel instances of Theorem 3, each of which sketches a subgraph $G^{(\ell)}$. Then the space of the algorithm will be $\tilde{O}(Wn^{1 - \frac{\varepsilon}{W-1}})$.

Theorem 4 *Let $t \geq 1$ be any constant. There exists a single-pass dynamic streaming algorithm that uses space and post-processing time $\tilde{O}(e^t Wn^{1 - \frac{\varepsilon}{W-1} + \frac{\varepsilon^t}{(W-1)^t}})$ to compute a $(1 + \varepsilon)$ -approximation of the weight of the MST.*

We remark that Ahn et al. [2] have given a dynamic streaming algorithm for this problem for any graph with maximum edge weight upper bounded by $O(\text{poly}(n))$, and their algorithm uses space $O(n \cdot \text{poly} \log n)$. Our algorithm uses $o(n)$ space for any connected graph with maximum edge weight bounded by $o(\log n)$ (for constant ε), which improves the algorithm of [2] in this setting. We also note that $\Omega(n)$ space is necessary for estimating the weight of MST for graphs with maximum edge weight at least $c \log n$ for constant ε and some large universal constant c (see Theorem 10). Finally, we remark that the algorithm can also be extended to the setting where non-integral weights are allowed (see [8] for more details).

4 Dynamic Streaming Testers

In this section, we give our streaming testers for a number of graph properties, including k -edge connectivity, k -vertex connectivity, cycle-freeness, planar graph bipartiteness, and Eulerianity.

4.1 Testing k -Edge Connectivity

A graph is k -edge connected if the minimum cut of the graph has size at least k . We start from the simplest case, i.e., $k = 1$, which is equivalent to the problem of testing connectivity.

4.1.1 Connectivity

It is clear that if G is ε -far from being connected, one must add at least εm edges to make it connected, which implies that there are at least $\varepsilon m + 1$ connected components in G [22,36]. Therefore, we can distinguish a connected graph from any graph that is ε -far from being connected by estimating the number of connected components with an additive error $\Theta(\varepsilon n)$. However, by a more careful analysis, we can reduce the space by a factor of $O(n^{O(\varepsilon)})$.

Theorem 5 *There exists a dynamic streaming tester for 1-edge connectivity that runs in $\tilde{O}(n^{1-\varepsilon})$ post-processing time and space.*

Proof First observe that one can simply reject the input graph if $m < n - 1$, since in this case, the graph is disconnected. Thus, in the following we assume $m \geq n - 1$ and our tester is described in Algorithm 2.

Algorithm 2 TestConnectivity

- 1: Sample each vertex with probability $p := (\varepsilon n/10)^{-\varepsilon}$. If more than $16np$ vertices are sampled, abort and output **Fail**. Let S denote the set of sampled vertices.
 - 2: For each $v \in S$, maintain an AMS sketch $AMS(\Delta^v)$, sketching the neighborhood of v in G .
 - 3: Maintain an AGM sketch of $G[S]$ using Theorem 1.
 - 4: **Post-Processing:**
 - 5: Use the above sketch to construct a spanning forest F of $G[S]$ as guaranteed by Theorem 1.
 - 6: For each connected component $C \in F$, estimate $F_2(\Delta^C)$ using the AMS sketch $AMS(\Delta^C) = \sum_{v \in C} AMS(\Delta^v)$. If the answer $\tilde{F}_2 = 0$, **Reject**.
 - 7: **Accept**.
-

It is easy to see that Algorithm 2 only uses $\tilde{O}(|S|)$ space, which is bounded by $\tilde{O}(np) = \tilde{O}(\varepsilon^{-\varepsilon} n^{1-\varepsilon}) = \tilde{O}(n^{1-\varepsilon})$. The post-processing time is nearly linear in the size of S , since the AGM algorithm needs $\tilde{O}(|S|)$ post-processing time, and we invoke at most $|S|$ AMS queries, each of which takes $\tilde{O}(1)$ time. The update time is poly log n .

For the correctness of the algorithm, we condition on the event that the number of sampled vertices is at most $16np$, which occurs with probability at least $1 - \frac{1}{16}$, and on the event that the spanning forest F is constructed correctly, which occurs with probability 0.99. By setting the error probability of the AMS sketch to be $1/n^2$ (with an extra log n factor in space), with probability 0.99, all the answers from AMS sketches are all correct, and we also condition on this.

If G is connected, then it will always be accepted, since for each $C \in F$, $\Delta^C \neq \mathbf{0}$, and conditioned on the correctness of the AMS sketch, \tilde{F}_2 will never be 0. On the other hand, if the graph is ε -far from being connected, the number of connected components in G , denoted as $cc(G)$, is at least $1 + \varepsilon m \geq \varepsilon n$. Let $B_1, \dots, B_{cc(G)}$ denote all connected components in G . Let $p_i = p^{|B_i|}$ for $1 \leq i \leq cc(G)$. Using the inequality $1 - x \leq e^{-x}$ for all x , the probability that none of the components is entirely sampled out is $(1 - p_1) \cdot (1 - p_2) \cdot \dots \cdot (1 - p_{cc(G)}) \leq e^{-\sum_i p_i}$. Then by the AM-GM inequality, this probability is at most

$$e^{-cc(G) \cdot (\prod_i p_i)^{1/cc(G)}} = e^{-cc(G) \cdot p^{n/cc(G)}} \leq e^{-cc(G) \cdot p^{1/\varepsilon}} \leq e^{-\varepsilon n \cdot p^{1/\varepsilon}} \leq 1/16,$$

where we use the fact that $p = (\varepsilon n/10)^{-\varepsilon}$ and $cc(G) \geq \varepsilon n$. So the probability that at least one of the components is sampled out is at least $15/16$. Conditioned on this, $F_2(\Delta^C) = 0$ for some component in $G[S]$ and the algorithm will output **Reject**. By union bound, our algorithm will succeed with probability $1 - \frac{1}{16} - 0.01 - 0.01 - \frac{1}{16} > 3/4$. □

4.1.2 k -Edge Connectivity: $k \geq 2$

By using a slightly more involved argument and replacing AMS sketches with $(k - 1)$ -sparse recovery sketches, we can generalize the above idea to testing k -edge connectivity for $k \geq 2$. We have the following theorem on testing k -edge connectivity.

Theorem 6 *Let $k \leq O(n^{\varepsilon/(1+\varepsilon)})$. There exists a single-pass dynamic streaming tester for k -edge connectivity with post-processing time and space $\tilde{O}(k^{1+\varepsilon} \cdot n^{1-\varepsilon})$.*

In order to prove Theorem 6, we will use the following result by Orenstein and Ron [35], who have given, for any $k \geq 2$, a characterization of graphs that are ε -far from being k -edge connected (which simplifies the corresponding result in [22]). We define a subset C to be ℓ -extreme if $|E(C, V \setminus C)| = \ell < k$ and for any $C' \subset C$, $|E(C', V \setminus C')| > \ell$.

Lemma 3 (Corollary 14 and Claim 16 in [35]) *If G is ε -far from being k -edge connected, then there are at least $\frac{2\varepsilon m}{k}$ disjoint subsets with an edge-cut smaller than k . For each such a subset C , it contains a minimal subset $C' \subseteq C$ that is ℓ -extreme for some $\ell < k$.*

Now we present the proof of Theorem 6.

Proof of Theorem 6 It is clear that $m \geq nk/2$ for any k -connected graph, and thus we can safely reject whenever $m < nk/2$. In the following, we will only consider the case that $m \geq nk/2$. Our tester is then described in Algorithm 3.

Algorithm 3 TestKEdgeConnectivity

- 1: Sample each vertex with probability $p := (\varepsilon n/4k)^{-\varepsilon}$. If more than $16np$ vertices are sampled, abort and output **Fail**. Let S denote the set of sampled vertices.
 - 2: For each $v \in S$, maintain a $(k - 1)$ -sparse recovery sketch $S_{k-1}(\Delta^v)$.
 - 3: Maintain an AGM sketch of $G[S]$ using Theorem 1.
 - 4: **Post-Processing:**
 - 5: Use the above sketch to recover a spanning forest F of $G[S]$ using Theorem 1.
 - 6: For each component $C \in F$, recover Δ^C from $S_{k-1}(\Delta^C)$, and if it succeeds, **Reject**.
 - 7: **Accept**.
-

Note that the AGM sketch use space $\tilde{O}(|S|) = \tilde{O}(np) = \tilde{O}(k^\varepsilon n^{1-\varepsilon})$. In addition, each sampled vertex only needs to store a k -sparse recovery sketch, so the space complexity of the algorithm is $\tilde{O}(k) \cdot np = \tilde{O}(k^{1+\varepsilon} n^{1-\varepsilon})$. The post-processing time is near linear in the space, and the update time is $O(\text{poly } \log n)$.

For the correctness of the algorithm, we first note that if G is k -edge connected, then G will be accepted as long as there is no error happening when querying the k -sparse recovery sketches. This happens with probability $1 - 1/n$ by setting the error probability of the sketch to be $1/n^2$, and we will condition on this event.

Now if G is ε -far from being k -edge connected, then from Lemma 3, it follows that there are at least $\frac{2\varepsilon m}{k} \geq \varepsilon n$ disjoint ℓ -extreme subsets. Let B_1, \dots, B_s be the set

of these ℓ -extreme subsets where $s \geq \varepsilon n$. Observe that for any ℓ -extreme subset B , the induced subgraph $G[B]$ is connected. This is true since otherwise, there exists a subset $B' \subset B$ satisfying $|E(B', B \setminus B')| = 0$, which implies that $|E(B', V \setminus B')| \leq |E(B, V \setminus B)| = \ell$, contradicting to the assumption that B is ℓ -extreme.

Let \mathcal{E}_i be the event that B_i is entirely sampled out, and \mathcal{F}_i be the event that none of the vertices in $\Gamma(B_i)$ is sampled.

Note that our algorithm will **reject** if $\mathcal{E}_i \wedge \mathcal{F}_i$ happens for some i , and thus our theorem will follow from the inequality that

$$\Pr \left[\bigvee_i (\mathcal{E}_i \wedge \mathcal{F}_i) \right] \geq \frac{3}{4}. \tag{2}$$

Now we prove inequality (2). Note that the events $(\mathcal{E}_i \wedge \mathcal{F}_i)$ are not necessarily independent across i since two different ℓ -extreme subsets may contain neighbors of each other or share neighbors. We have the following simple lemma to deal with this issue.

Lemma 4 *There exists a set $I \subset [s]$, with $|I| = s/k$, such that:*

1. $|E(B_i, B_j)| = 0$ for all $i, j \in I$ and $i \neq j$, and
2. $\sum_{i \in I} |B_i| \leq \sum_{j=1}^s |B_j|/k$.

Proof We say B_i and B_j are neighbors, if $|E(B_i, B_j)| > 0$. We iteratively construct the index set $I \subset [s]$ as follows. We start from the empty set $I_0 = \emptyset$ and add one index at each step. Let I_t denote the set that at the end of step t . In the $(t + 1)$ -th step, we find the smallest set $B_{i_{t+1}}$ that is not a neighbor of B_{i_h} for any $h \leq t$ and add the index i_{t+1} , i.e., $I_{t+1} = I_t \cup \{i_{t+1}\}$. Note that since each ℓ -extreme set has at most $k - 1$ neighbors, we can always find such a set if $t < s/k$. Let $I = I_{s/k}$. Then Item 1 of the lemma follows by our construction. Since the set B_{i_t} that we found in the t -th step may intersect with at most k sets, and B_{i_t} is the smallest set that has no intersection with all sets found in the first $t - 1$ steps, there must exist a partition of $[s]$ into s/k sets $\{P_1, P_2, \dots, P_{s/k}\}$, such that for any $t \leq s/k$ and $j \in P_t$, $|B_j| \geq |B_{i_t}|$. Item 2 of the lemma then follows from our construction of the index subset $I = \{i_j : 1 \leq j \leq s/k\}$, and the fact that $[s] = \bigcup_{t=1}^{s/k} P_t$. \square

Now we give a lower bound for $\Pr[\bigvee_{i \in I} \mathcal{E}_i]$. Let $p_i = p^{|B_i|}$ be the probability that all vertices in B_i are sampled. Using the fact $1 - x \leq e^{-x}$ for all x and the AM-GM inequality, we have

$$\begin{aligned} \ln \prod_{i \in I} (1 - p_i) &\leq - \sum_{i \in I} p_i \leq -|I| \cdot \left(\prod_{i \in I} p_i \right)^{1/|I|} = -|I| \cdot p^{\sum_{i \in I} |B_i|/|I|} \\ &\leq -(s/k) \cdot p^{\sum_i |B_i|/s} \\ &\leq -\frac{\varepsilon n}{k} \cdot p^{1/\varepsilon}. \end{aligned}$$

Thus we have

$$\Pr \left[\bigvee_{i \in I} \mathcal{E}_i \right] = 1 - \left(\prod_{i \in I} (1 - p_i) \right) = 1 - e^{\ln \prod_{i \in I} (1 - p_i)} \geq 1 - e^{-\frac{\varepsilon n}{k} \cdot p^{1/\varepsilon}} \geq 15/16,$$

since we set $p = (\varepsilon n / 4k)^{-\varepsilon}$.

Now by the property of I as guaranteed in Lemma 4, it follows that \mathcal{F}_j and \mathcal{E}_i are independent for all $i, j \in I$. Hence, conditioned on the event $\bigvee_{i \in I} \mathcal{E}_i$, the probability of $\bigvee_{i \in I} (\mathcal{E}_i \wedge \mathcal{F}_i)$ happening is

$$\begin{aligned} \Pr \left[\bigvee_{i \in I} (\mathcal{E}_i \wedge \mathcal{F}_i) \mid \bigvee_{i \in I} \mathcal{E}_i \right] &\geq \min_{j \in I} \Pr[\mathcal{F}_j] = \min_{j \in I} (1 - p)^{|\Gamma(B_j)|} \geq (1 - p)^k \\ &\geq e^{-pk - p^2k} \\ &\geq 0.8, \end{aligned}$$

where in the penultimate inequality, we used the basic inequality that $1 - x \geq e^{-x - x^2}$ for $x \leq 0.5$; the last inequality holds for $k \leq 0.1/p$ or equivalent $k \leq O(n^{\varepsilon/(1+\varepsilon)})$.

Finally, we have

$$\begin{aligned} \Pr \left[\bigvee_i (\mathcal{E}_i \wedge \mathcal{F}_i) \right] &\geq \Pr \left[\bigvee_{i \in I} (\mathcal{E}_i \wedge \mathcal{F}_i) \right] = \Pr \left[\left(\bigvee_{i \in I} (\mathcal{E}_i \wedge \mathcal{F}_i) \right) \wedge \left(\bigvee_{i \in I} \mathcal{E}_i \right) \right] \\ &= \Pr \left[\bigvee_{i \in I} \mathcal{E}_i \right] \cdot \Pr \left[\bigvee_{i \in I} (\mathcal{E}_i \wedge \mathcal{F}_i) \mid \bigvee_{i \in I} \mathcal{E}_i \right] \\ &\geq \frac{15}{16} \cdot \frac{4}{5} \geq 3/4. \end{aligned}$$

□

We remark that the problem can still be solved in space $\tilde{O}(kn^{1-\varepsilon})$ for larger k by testing the neighborhood of all subsets of size smaller than $1/\varepsilon$ in S , however the post-processing time will be $\tilde{O}(kn^{O(1/\varepsilon)})$. Also, $k \leq O(n^\varepsilon)$ is the most interesting case for us, since we are mostly interested in $o(n)$ space algorithms.

4.2 k -Vertex Connectivity

A graph is k -vertex connected if the minimum vertex cut of the graph has size at least k , i.e. it remains connected whenever fewer than k vertices are removed. The following lemma on the structure of graphs that are ε -far from being k -vertex connected can be directly deduced from Corollary 19 in [35].

Lemma 5 *If the graph is ε -far from k -vertex connected, then there exists at least $\frac{\varepsilon m}{2k}$ subsets C of size at most $\frac{2kn}{\varepsilon m}$ such that $G[C]$ is connected and $\Gamma(C) < k$.*

Proof (sketch) In Corollary 19 in [35], it is proven that for any directed graph G that is ε -from k -vertex connected, then there exists at least $\frac{\varepsilon m}{2k}$ subsets C of size at most $\frac{2kn}{\varepsilon m}$, and either $|\Gamma^+(C)| < k$ or $|\Gamma^-(C)| < k$, where $\Gamma^+(C) := \{v \in V \setminus C : \langle v, u \rangle \in E(G), u \in C\}$ (resp., $\Gamma^-(C) := \{v \in V \setminus C : \langle u, v \rangle \in E(G), u \in C\}$) denotes the set of vertices in $V \setminus C$ that are endpoints of incoming (resp., outgoing) edges incident to C .

On the other hand, in Sect. 5.3 in [35], it is proven that if an undirected graph G is ε -far from k -vertex connected, then the corresponding directed graph G' that is obtained by turning each undirected edge (u, v) into directed edges $\langle u, v \rangle$ and $\langle v, u \rangle$ is ε -far from being k -vertex connected. Therefore, there exists at least $\frac{\varepsilon m}{2k}$ subsets C in G' of size at most $\frac{2kn}{\varepsilon m}$, and either $|\Gamma_{G'}^+(C)| < k$ or $|\Gamma_{G'}^-(C)| < k$. This directly implies that the corresponding set C in G satisfies that $|\Gamma_G(C)| < k$. Finally, if $G[C]$ is not connected, then we can replace C by one maximal subset $C' \subset C$ such that $G[C']$ is connected. Note that $|\Gamma_G(C')| \leq |\Gamma_G(C)| < k$. This completes the proof of the lemma. \square

Now we use the above lemma to show our k -vertex connectivity tester.

Theorem 7 *Let $k \leq O(n^{\frac{\varepsilon}{4+\varepsilon}})$. There exists a single-pass dynamic streaming tester for k -vertex connectivity with post-processing time and space complexity $\tilde{O}(\frac{k^{1+\varepsilon/4}}{\varepsilon} \cdot n^{1-\varepsilon/4})$.*

Proof (sketch) We can also simply consider the case that $m \geq nk/2$, since otherwise the graph cannot be k -vertex connected and we can directly reject. Our approach for testing k -connectivity is similar to testing k -edge connectivity. The difference here is that now we cannot use the $(k - 1)$ -sparse recovery sketch for the vector $\mathbf{\Lambda}^v$. Instead, for each vertex $v \in S$, we will maintain an exact k' -sparse recovery sketch of the vector $\mathbf{\Lambda}^v$ (defined in Sect. 2.1), $\mathcal{S}_{k'}(\mathbf{\Lambda}^v)$, for $k' = \frac{4}{\varepsilon} + k$. Then for each detected connected component C of size smaller than $4/\varepsilon$ in $G[S]$ (by AGM sketch), recover $\mathbf{\Lambda}^C := \sum_{v \in C} \mathbf{\Lambda}^v$ from the sketch $\mathcal{S}_{k'}(\mathbf{\Lambda}^C) = \sum_{v \in C} \mathcal{S}_{k'}(\mathbf{\Lambda}^v)$. If it succeeds, we get the set $C \cup \Gamma(C)$, and since we know C , we get $\Gamma(C)$. If $|\Gamma(C)| < k$, we **reject**. For any k -vertex connected graph, the tester will never **reject** if all the sparse recover sketches return correctly, which happens with high probability. On the other hand, if G is ε -far from k -vertex connected, by similar analysis as in k -edge connectivity together with Lemma 5, we know that with high probability, there is a subset $C \subseteq S$ such that $G[C]$ is a connected component in $G[S]$, $|\Gamma(C)| < k$ and $|C| \leq 4/\varepsilon$, and conditioned on this the algorithm will successfully recover $\Gamma(C)$, and reject with high probability. Here to make the analysis work, we have to set the sampling probability $p := (\varepsilon n / 16k)^{-\varepsilon/4}$, so the space used is $\tilde{O}(k' \cdot k^{\varepsilon/4} \cdot n^{1-\varepsilon/4}) = \tilde{O}(\frac{k^{1+\varepsilon/4}}{\varepsilon} \cdot n^{1-\varepsilon/4})$. Since the analysis is almost the same as k -edge connectivity, we omit the details here. \square

4.3 Testing Cycle-Freeness

Now we consider the problem of testing cycle-freeness, which is equivalent to testing if the graph is a forest. Let $cc(G)$ denote the number of connected components of the

input graph G . Let $B_1, \dots, B_{cc(G)}$ be the connected components in G . Note that if G is cycle-free, then for each $i \leq cc(G)$, $|E(B_i)| = |B_i| - 1$, and thus the total number of edges in G is

$$m = \sum_{i=1}^{cc(G)} |E(B_i)| = \sum_{i=1}^{cc(G)} (|B_i| - 1) = n - cc(G),$$

that is, $cc(G) = n - m$. If G is ε -far from being cycle-free, i.e., one has to delete more than εm edges to make it cycle-free, then $cc(G) > n - m + \varepsilon m$. Therefore, to test cycle-freeness of a graph, it will be sufficient to approximate the number of connected components with additive error $\varepsilon m/2$. One may try to directly invoke Algorithm 1 with parameter $\varepsilon' = \frac{\varepsilon m}{2n}$. However, m could be much smaller than n and we do not know m in advance. We overcome this obstacle by a case analysis.

Theorem 8 *There exists a single-pass dynamic streaming algorithm that tests cycle-freeness of a graph with space and post-processing time $\tilde{O}(n^{1-\varepsilon+\varepsilon^2})$.*

Proof Note that if $m > n - 1$, then the graph must contain at least one cycle, and thus we can safely reject the graph. In the following, we assume that $m \leq n - 1$. Note that if $\varepsilon \leq 1/(10 \log n)$, then we can simply store whole graph and test if it is cycle-free, as the size of the graph is $O(n) = O(n^{1-\varepsilon+\varepsilon^2}) \cdot \text{poly} \log n$. In the following, we will assume that $\varepsilon > 1/(10 \log n)$. Our algorithm for testing cycle-freeness depends on the construction of AGM sketch, in which each vertex u maintains a linear sketch of Δ^u (denoted as $\mathcal{A}(\Delta^u)$). Each such sketch has size $\text{poly} \log n$ and the property that $\mathcal{A}(\mathbf{0}) = \mathbf{0}$ (it consists of $O(\log n)$ l_0 -samplers, see [2] for details). Our main idea is to maintain a sparse recovery sketch for the AGM sketch (i.e. a composition of sparse recovery sketch and AGM sketch). Now we describe our algorithm in Algorithm 4.

Note that the space used by the algorithm is $\max\{\tilde{O}(np_0), \tilde{O}(np), k \cdot \text{poly} \log n\} = \tilde{O}(n^{1-\varepsilon+\varepsilon^2} + n^{2\varepsilon/(1+\varepsilon+\varepsilon^2)}/\varepsilon^6) = \tilde{O}(n^{1-\varepsilon+\varepsilon^2})$ as $\varepsilon > 1/(10 \log n)$, and the post-processing time is near linear in space.

Now we prove the correctness of the algorithm. We define $G' \subseteq G$ to be a subgraph which consists of all the vertices of positive degree. Let $n' = |G'|$. Note that $m \geq n'/2$.

If $n' \leq n^{1-\eta}$, then the vector \mathbf{Y} is $\tilde{O}(n^{1-\eta})$ -sparse, since for all isolated vertices u , we have $\mathcal{A}(\Delta^u) = \mathbf{0}$, and thus we can recover the entire \mathbf{Y} exactly. Then by Step 2-(c) and Theorem 1, we can get the exact number of components of G' . Since the number of vertices of G' is $|Y|$, and $\lambda = m$ is the total number of edges, then the graph is cycle-free if and only if $\tilde{c}_1 = |Y| - \lambda$.

If $n' > n^{1-\eta}$, then conditioned on the event that all $AMS(\Delta^v)$ for $v \in S_0$ are correct (which occurs with high probability), our estimator \tilde{c}_0 approximates the number of isolated vertices in G , denoted by c_0 , with an additive error $(\varepsilon^3/16)n^{1-\eta}$ with probability at least $1 - \frac{1}{16}$ (by our choice that $p_0 = \frac{1}{4096\varepsilon^6 n^{1-2\eta}}$ and similar analysis for the proof of Theorem 2). We will condition on this event. Since $n' > n^{1-\eta}$ and $m \geq n'/2$, we have that $|c_0 - \tilde{c}_0| \leq (\varepsilon^3/8)m$.

Now note that by Theorem 2, \tilde{c}_2 is an estimator for the number, denoted by c_2 , of components in G' of size smaller than $1/\eta$ with additive error $\eta^t \sqrt{n'n^{1-\eta}}$. This

Algorithm 4 TestCycleFreeness

- 1: Maintain a count λ of the number of edges.
- 2: Let $\eta = \varepsilon/(1 + \varepsilon + \varepsilon^2)$. Let $k = n^{1-\eta}$ poly log n . Perform the following steps (a),(b),(c) in parallel:
 - (a) Maintain an exact k -sparse recovery sketch \mathcal{S} of the vector $\Upsilon := (\mathcal{A}(\Delta^u))_{u \in V}$ using Lemma 1.
 - (b) Sample each vertex with probability $p_0 := \frac{1}{4096\varepsilon^6 n^{1-2\eta}}$. For each sampled vertex v , maintain the AMS sketch $AMS(\Delta^v)$. Let S_0 denote the resulting sample set.
 - (c) Run **Algorithm 1** with parameter $p = (\eta^{2t} n^{1-\eta}/16)^{-\eta}$, while in step (6) of **Algorithm 1**, ignore all the isolated vertices that are sampled out (i.e., set $X_C = 0$ whenever $|C| = 1$).
- 3: **Post-Processing:**
- 4: Recover Υ from \mathcal{S} .
- 5: **if** The recovery does not fail **then**
- 6: Use Υ to construct a spanning forest on vertex set $Y := \{u : \mathcal{A}(\Delta^u) \neq \mathbf{0}\}$ using Theorem 1. Let \tilde{c}_1 denote the number of connected components of this forest. If $\tilde{c}_1 = |Y| - \lambda$, **Accept**; otherwise, **Reject**.
- 7: **else**
- 8: Let X_0 denote the number of vertices v in S_0 with $F_2(\Delta^v) = 0$ (by the AMS sketches $AMS(\Delta^v)$). Let $\tilde{c}_0 := \frac{X_0}{p}$.
- 9: Let \tilde{c}_2 be the resulting estimator of **Algorithm 1** in Step 2-(c). If $\tilde{c}_0 + \tilde{c}_2 \leq n - \lambda + \frac{\varepsilon^3}{4}\lambda$, **Accept**; otherwise, **Reject**.
- 10: **end if**

follows by the upper bound $\eta^{2t} n^{1-\eta} n'/16$ of the variance of the estimator (which can be shown similarly to inequality (1) in Sect. 3) and the Chebyshev’s inequality. Now note that the additive error is at most $\eta^t n' \leq \varepsilon^3 m/8$ for some constant t since $n' > n^{1-\eta}$ and $m \geq n'/2$. That is, with high probability, $|c_2 - \tilde{c}_2| \leq \varepsilon^3 m/8$. In the following, we condition on this event.

Let L be the number of components in G' of size larger than $1/\eta$. Note that each such component has at least $1/\eta - 1$ edges. Thus, $m \geq L \cdot (1/\eta - 1)$, which gives that $L \leq \frac{\eta}{1-\eta} m = \frac{\varepsilon}{1+\varepsilon^2} m$ by our choice of η .

If the original graph G is cycle-free, then the number of connected components of G equals $n - m$, i.e., $c_0 + L + c_2 = n - m$. Thus, we have that $\tilde{c}_0 + \tilde{c}_2 \leq c_0 + \varepsilon^3 m/8 + c_2 + \varepsilon^3 m/8 = n - m - L + \varepsilon^3 m/4 \leq n - m + \frac{\varepsilon^3}{4} m$. The algorithm will output **Accept**.

If G is ε -far from being cycle-free, then $c_0 + L + c_2 > n - m + \varepsilon m$. Thus, $\tilde{c}_0 + \tilde{c}_2 \geq c_0 - \varepsilon^3 m/8 + c_2 - \varepsilon^3 m/8 > n - m + \varepsilon m - L - \varepsilon^3 m/4 \geq n - m + \frac{\varepsilon^3}{1+\varepsilon^2} m - \varepsilon^3 m/4 > n - m + \varepsilon^3 m/4$. The algorithm will output **Reject**.

Thus, our algorithm can distinguish cycle-free graphs from those graphs that are ε -far from being cycle-free with probability at least $2/3$. This completes the proof of the theorem. □

4.4 Testing Bipartiteness of the Planar Graphs

Now we consider the problem of testing if a planar graph is bipartite or ε -far from bipartite. Here a planar graph is ε -far from bipartite if one has to delete at least εm edges to get a bipartite graph. Czumaj et al. [12] showed the following result¹.

¹ In [12], ε -far is expressed as εn edges, rather than εm edges as in our definition, that has to be deleted to obtain a bipartite graph. However, Lemma 6 directly follows from their proof.

Lemma 6 [12] *For any (simple) planar graph G that is ε -far from bipartite, then G has at least $\varepsilon m/q(\varepsilon)$ edge-disjoint odd-length cycles of length at most $q(\varepsilon)/2$ each, where $q(\varepsilon) = O(1/\varepsilon^2)$.*

By the above lemma, we only need to sample each *edge* independently with some probability (rather than vertices as we did before) of the graph so that with high probability the resulting sampled graph contains at least one short odd-length cycle. The edge-sampling process can be done by using hash functions (see e.g. [3]). Similar to our previous analysis, it will be sufficient to set the sample probability to $p = O_\varepsilon(n^{-q(\varepsilon)})$, which implies that the space used is $\tilde{O}(n^{1-\Omega(\varepsilon^2)})$. We omit the details here.

4.5 Testing Eulerianity

Note that the algorithm for connectivity testing can be directly used to testing Eulerianity. A graph G is Eulerian if there is a path in the graph that traverses each edge exactly once, or equivalently, if G is connected and the degrees of all vertices are even or exactly two vertices have odd degrees. Note that if graph G is ε -far from being Eulerian then either G has $\Omega(\varepsilon n)$ connected components (i.e. far from being connected) or has $\Omega(\varepsilon n)$ vertices of odd degree (cf., [22,36]). Then one can test Eulerianity by first invoking the previous algorithm on testing connectivity, and then sample $O(1/\varepsilon)$ vertices and check if some sampled vertex has odd degree. The post-processing time and space complexity of the final algorithm are $\tilde{O}(n^{1-c\varepsilon})$ for some universal constant c .

5 Lower Bounds

In this section we present lower bounds, which hold in the insertion-only model. Our proofs are based on the reductions to the *Boolean Hidden Hypermatching (BHH)* problem (See [42]), which are in the same spirit as the lower bound proof for the *Cycle Counting* problem in [42]. We first give the definition of the boolean hidden hypermatching problem.

Definition 1 (BHH_n^t) In this problem, Alice gets a boolean vector $x \in \{0, 1\}^n$, where $n = 2kt$ for some integer k . Bob gets a partition (or hypermatching) of the set $[n]$, $\{m_1, \dots, m_{n/t}\}$, where the size of each m_i is t , and a vector $w \in \{0, 1\}^{n/t}$. For convenience, we will also use the corresponding n -dimensional boolean indicator vector M_i to represent m_i , and let M be a $n/t \times n$ matrix, the i row of which is M_i . The promise of the input is either $Mx + w = \mathbf{1}$ or $Mx + w = \mathbf{0}$, where all the operations are modulo 2. The goal of the problem is to output 1 when $Mx + w = \mathbf{1}$, and output 0 otherwise.

We have the following lower bound from [42].

Theorem 9 [42] *The randomized one-way communication complexity of BHH_n^t when $n = 2kt$ for some integer $k \geq 1$ is $\Omega(n^{1-1/t})$.*

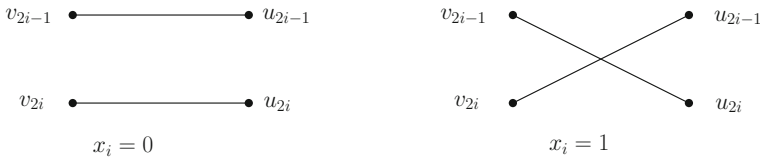


Fig. 1 Parallel (left) and crossing (right) matching according to the value of x_i

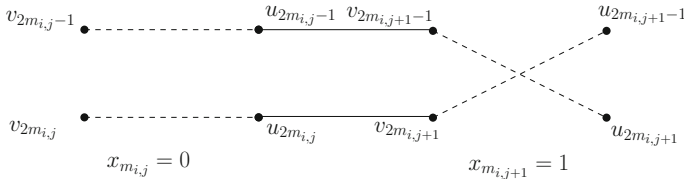


Fig. 2 Bob connects $(u_{2m_i,j-1}, v_{2m_i,j+1-1})$ and $(u_{2m_i,j}, v_{2m_i,j+1})$ for each $j \in [t - 1]$

Our lower bounds will be built upon the following basic construction.

Construction of $G(x, M)$. Given vector x and matrix M respectively, Alice and Bob construct a bipartite graph $G(x, M) = (U, V, E)$, where $U = \{u_1, \dots, u_{2n}\}$ and $V = \{v_1, \dots, v_{2n}\}$, as follows. Given $x \in \{0, 1\}^n$, Alice adds a perfect matching between U and V . For each $i \in [n]$, if $x_i = 0$, she adds two parallel edges (u_{2i-1}, v_{2i-1}) and (u_{2i}, v_{2i}) ; otherwise if $x_i = 1$, she adds two crossing edges (u_{2i-1}, v_{2i}) and (u_{2i}, v_{2i-1}) (see Fig. 1).

Given M , Bob will do the following. For each $i \in [n/t]$ and the hyperedge $m_i \subset [n]$ (that corresponds to the i th row M_i), we use $m_{i,j} \in [n]$ to denote the j th element in m_i and we let $S_i := \{x \mid x = v_{2m_{i,j-1}} \text{ or } v_{2m_{i,j}} \text{ or } u_{2m_{i,j-1}} \text{ or } u_{2m_{i,j}}, j \in [t]\}$. For each $i \in [n/t]$ and $j \in [t - 1]$, Bob adds two edges $(u_{2m_{i,j-1}}, v_{2m_{i,j+1-1}})$ and $(u_{2m_{i,j}}, v_{2m_{i,j+1}})$ (see Fig. 2).

Observe that the edges added by Alice and Bob form two paths p_{2i-1}, p_{2i} over vertex set S_i , where p_{2i-1} starts from $v_{2m_{i,1-1}}$ and p_{2i} starts from $v_{2m_{i,1}}$ for each i . The entire graph $G(x, M)$ consists of $2n/t$ disjoint paths $\{p_1 \dots, p_{2n/t}\}$. It also has the following property.

Fact 1 *Based on the value of $(Mx)_i$, we have: 1) if $(Mx)_i = 0$, then p_{2i-1} is a path from $v_{2m_{i,1-1}}$ to $u_{2m_{i,t-1}}$ and p_{2i} is a path from $v_{2m_{i,1}}$ to $u_{2m_{i,t}}$; 2) if $(Mx)_i = 1$, then p_{2i-1} is a path from $v_{2m_{i,1-1}}$ to $u_{2m_{i,t}}$ and p_{2i} is a path from $v_{2m_{i,1}}$ to $u_{2m_{i,t-1}}$.*

5.1 Minimum Spanning Tree

Theorem 10 *In the insertion-only model, if all edges of the graph have weights in $[W]$, any algorithm that $(1 \pm \varepsilon)$ -approximates the weight of the MST must use $\Omega(n^{1 - \frac{4\varepsilon}{W-1}})$ bits of space.*

Proof Given x and M , Alice and Bob first construct the graph $G(x, M)$ as described above. Next Bob adds $(u_{2m_{i,t-1}}, v_{2m_{i,1-1}})$ and $(u_{2m_{i,t}}, v_{2m_{i,1}})$ if $w_i = 0$; adds $(u_{2m_{i,t-1}}, v_{2m_{i,1}})$ and $(u_{2m_{i,t}}, v_{2m_{i,1-1}})$ if $w_i = 1$. The weight of all the edges added

so far is 1. Next, regardless of the value of w_i , Bob places edges $(v_{2m_{i,t}}, v_{2m_{i+1,t}})$ with weight 1 for $i = 1, \dots, n/t - 1$ and edges $(v_{2m_{i,t}}, u_{2m_{i,t}})$ with weight W for each $i \in [n/t]$, so that the graph become connected. By similar argument as above, if $Mx + w = \mathbf{0}$, all the edges $(v_{2m_{i,t}}, u_{2m_{i,t}})$ must be picked in any minimum spanning tree, since each of these edges forms a cut, and thus the weight of any MST is $nW/t + 4n - n/t - 1 = 4n\varepsilon + 4n - 1$, where we set $t = (W - 1)/4\varepsilon$. On the other hand, when $Mx + w = \mathbf{1}$, the weight of the MST is $4n - 1$, since in this case, the graph is already connected without those edges with weight W . So if the algorithm can compute an $(1 + \varepsilon)$ -approximation of the weight of the minimum spanning tree, it solves the BHH_n^t problem. This completes the proof. \square

5.2 Testing Connectivity

Theorem 11 *In the insertion-only model, to distinguish whether a graph of $4n$ vertices is connected or $\frac{1}{8t+1}$ -far from being connected, any algorithm must use $\Omega(n^{1-1/t})$ bits of space.*

Proof Given x and M , Alice and Bob first construct the graph $G(x, M)$. Next Bob adds another set of edges based on vector w . If $w_i = 0$, he adds $(u_{2m_{i,t}-1}, v_{2m_{i,t}-1})$ and $(u_{2m_{i,t}}, v_{2m_{i,t}})$; if $w_i = 1$, he adds $(u_{2m_{i,t}-1}, v_{2m_{i,t}})$ and $(u_{2m_{i,t}}, v_{2m_{i,t}-1})$. So when $(Mx)_i + w_i = 0$, p_{2i-1} and p_{2i} become 2 disjoint cycles. On the other hand, when $(Mx)_i + w_i = 1$, p_{2i-1} and p_{2i} together form a larger cycle. Now Bob places $(v_{2m_{i,t}}, v_{2m_{i+1,t}})$ in E for $i = 1, \dots, n/t - 1$ which connect p_{2i} with $p_{2(i+1)}$ for all $i \in [n/t - 1]$, i.e. all the paths in $G(x, M)$ with even indices become a connected component. The total number of edges is $8n + n/t$. When $Mx + w = \mathbf{0}$, the graph has $n/t + 1$ components which is $\frac{1}{8t+1}$ -far from connected; when $Mx + w = \mathbf{1}$ the graph is connected. So if a streaming algorithm can distinguish whether a graph of size $4n$ is connected or $1/8t$ -far from being connected, it solves BHH_n^t , since Alice can first run the algorithm on her part of the graph and send the memory to Bob, and then Bob continues to run the algorithm on his part and output the answer. Therefore, the communication lower bound of BHH_n^t implies a space lower bound of testing connectivity. \square

5.3 Testing Cycle-Freeness

As in the proof of Theorem 11, given x and M , Alice and Bob first construct $G(x, M)$. Then, for $i \in [n/t]$, Bob adds $(u_{2m_{i,t}-1}, v_{2m_{i,t}-1})$ if $w_i = 0$; adds $(u_{2m_{i,t}-1}, v_{2m_{i,t}})$ if $w_i = 1$. The total number of edges is less than $8n$. Through similar arguments, it is easy to verify that if $Mx + w = \mathbf{0}$, the graph has exactly n/t cycles and n/t paths, which is $1/8t$ -far from cycle-free. On the contrary, if $Mx + w = \mathbf{1}$, the graph has n/t paths and no cycle. So if an algorithm can distinguish whether a graph of size $4n$ is cycle-free or $1/8t$ -far from cycle-free, it solves BHH_n^t .

Theorem 12 *In the insertion-only model, any algorithm that can distinguish whether a graph of $4n$ vertices is cycle-free or $1/8t$ -far from being cycle-free, must use $\Omega(n^{1-1/t})$ bits of space.*

5.4 Testing Bipartiteness of Planar Graphs

Alice and Bob first construct the graph $G(x, M)$. Next, for each $i \in [n/t]$, Bob adds edges $(v_{2m_{i,1}-1}, \xi_1)$ and $(v_{2m_{i,1}}, \xi_2)$, where ξ_1, ξ_2 are new vertices. For $i \in [n/t]$, Bob also adds $(u_{2m_{i,t}-1}, \xi_1)$ and $(u_{2m_{i,t}}, \xi_2)$ if $w_i = 0$; adds $(u_{2m_{i,t}-1}, \xi_2)$ and $(u_{2m_{i,t}}, \xi_1)$ if $w_i = 1$. For this problem we assume t is odd. So by similar arguments, we can easily verify that, if $Mx + w = \mathbf{0}$, the graph contains $2n/t$ edge-disjoint cycles of length $2t + 1$, and if $Mx + w = \mathbf{1}$, the graph has no odd cycle, and thus bipartite. The graph constructed is planar and has $4n + 2$ vertices and $8n + 4n/t$ edges, so we have the following lower bound for testing bipartiteness.

Theorem 13 *In the insertion-only model, any algorithm that can distinguish whether a planar graph of $4n + 2$ vertices is bipartite or $\frac{1}{4t+2}$ -far from being bipartite, must use $\Omega(n^{1-1/t})$ bits space.*

Open Access This article is distributed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made.

References

- Ahn, K.J., Cormode, G., Guha, S., McGregor, A., Wirth, A.: Correlation clustering in data streams. In: Proceedings of the 32nd International Conference on Machine Learning, ICML, pp. 6–11 (2015)
- Ahn, K.J., Guha, S., McGregor, A.: Analyzing graph structure via linear measurements. In: Proceedings of the Twenty-third Annual ACM-SIAM Symposium on Discrete Algorithms, pp. 459–467. SIAM, Philadelphia (2012)
- Ahn, K.J., Guha, S., McGregor, A.: Graph sketches: sparsification, spanners, and subgraphs. In: Proceedings of the 31st Symposium on Principles of Database Systems, pp. 5–14. ACM, New York (2012)
- Alon, N., Matias, Y., Szegedy, M.: The space complexity of approximating the frequency moments. In: Proceedings of the Twenty-Eighth Annual ACM Symposium on Theory of Computing, pp. 20–29. ACM, New York (1996)
- Assadi, S., Khanna, S., Li, Y., Yaroslavtsev, G.: Maximum matchings in dynamic graph streams and the simultaneous communication model. In: Proceedings of the Twenty-Seventh Annual ACM-SIAM Symposium on Discrete Algorithms, SODA '16, pp. 1345–1364. SIAM, Philadelphia (2016)
- Bhattacharya, S., Henzinger, M., Nanongkai, D., Tsourakakis, C.E.: Space-and time-efficient algorithm for maintaining dense subgraphs on one-pass dynamic streams. In: ACM Symposium on Theory of Computing (2015)
- Bury, M., Schwiegelshohn, C.: Sublinear estimation of weighted matchings in dynamic data streams. ESA (2015)
- Chazelle, B., Rubinfeld, R., Trevisan, L.: Approximating the minimum spanning tree weight in sub-linear time. SIAM J. Comput. **34**(6), 1370–1379 (2005)
- Chitnis, R., Cormode, G., Esfandiari, H., Hajiaghayi, M., McGregor, A., Monemizadeh, M., Vorotnikova, S.: Kernelization via sampling with applications to dynamic graph streams. SODA (2016)
- Chitnis, R., Cormode, G., Hajiaghayi, M., Monemizadeh, M.: Parameterized streaming: maximal matching and vertex cover. In: Proceedings of the Twenty-Sixth Annual ACM-SIAM Symposium on Discrete Algorithms, pp. 1234–1251. SIAM, Philadelphia (2015)
- Czumaj, A., Ergün, F., Fortnow, L., Magen, A., Newman, I., Rubinfeld, R., Sohler, C.: Approximating the weight of the euclidean minimum spanning tree in sublinear time. SIAM J. Comput. **35**(1), 91–109 (2005)
- Czumaj, A., Monemizadeh, M., Onak, K., Sohler, C.: Planar graphs: random walks and bipartiteness testing. In: Foundations of Computer Science (FOCS), 2011 IEEE 52nd Annual Symposium on, pp. 423–432. IEEE (2011)

13. Czumaj, A., Sohler, C.: Estimating the weight of metric minimum spanning trees in sublinear time. *SIAM J. Comput.* **39**(3), 904–922 (2009)
14. Esfandiari, H., Hajiaghayi, M., Woodruff, D.P.: Brief announcement: applications of uniform sampling: densest subgraph and beyond. In: *Proceedings of the 28th ACM Symposium on Parallelism in Algorithms and Architectures, SPAA 2016*, pp. 397–399 (2016)
15. Esfandiari, H., Hajiaghayi, M.T., Liaghat, V., Monemizadeh, M., Onak, K.: Streaming algorithms for estimating the matching size in planar graphs and beyond. In: *Proceedings of the Twenty-Sixth Annual ACM-SIAM Symposium on Discrete Algorithms*, pp. 1217–1233. SIAM, Philadelphia (2015)
16. Fafianie, S., Kratsch, S.: Streaming kernelization. In: *Mathematical Foundations of Computer Science 2014*, pp. 275–286. Springer, Berlin (2014)
17. Feigenbaum, J., Kannan, S., McGregor, A., Suri, S., Zhang, J.: On graph problems in a semi-streaming model. *Theor. Comput. Sci.* **348**(2), 207–216 (2005)
18. Feigenbaum, J., Kannan, S., McGregor, A., Suri, S., Zhang, J.: Graph distances in the data-stream model. *SIAM J. Comput.* **38**(5), 1709–1727 (2008)
19. Frahling, G., Indyk, P., Sohler, C.: Sampling in dynamic data streams and applications. In: *Proceedings of the Twenty-First Annual Symposium on Computational Geometry*, pp. 142–149. ACM, New York (2005)
20. Goldreich, O.: Introduction to testing graph properties. In: *Property Testing*, pp. 105–141. Springer, Berlin (2011)
21. Goldreich, O., Goldwasser, S., Ron, D.: Property testing and its connection to learning and approximation. *J. ACM* **45**(4), 653–750 (1998)
22. Goldreich, O., Ron, D.: Property testing in bounded degree graphs. *Algorithmica* **32**, 302–343 (2002)
23. Guha, S., McGregor, A., Tench, D.: Vertex and hyperedge connectivity in dynamic graph streams. In: *Proceedings of the 34th ACM Symposium on Principles of Database Systems*, pp. 241–247. ACM, New York (2015)
24. Henzinger, M.R., Raghavan, P., Rajagopalan, S.: Computing on data streams. In: *External Memory Algorithms, Proceedings of a DIMACS Workshop*, New Brunswick, New Jersey, USA, May 20–22, pp. 107–118 (1998)
25. Jowhari, H.: Estimating the number of connected components in graph streams. *Personal Communication*
26. Kapralov, M., Khanna, S., Sudan, M.: Approximating matching size from random streams. In: *Proceedings of the Twenty-Fifth Annual ACM-SIAM Symposium on Discrete Algorithms*, pp. 734–751. SIAM, Philadelphia (2014)
27. Kapralov, M., Khanna, S., Sudan, M.: Streaming lower bounds for approximating max-cut. In: *Proceedings of the Twenty-Sixth Annual ACM-SIAM Symposium on Discrete Algorithms*, pp. 1263–1282. SIAM, Philadelphia (2015)
28. Kapralov, M., Lee, Y.T., Musco, C., Sidford, A.: Single pass spectral sparsification in dynamic streams. In: *Foundations of Computer Science (FOCS), 2014 IEEE 55th Annual Symposium on*, pp. 561–570. IEEE (2014)
29. Kapralov, M., Woodruff, D.: Spanners and sparsifiers in dynamic streams. In: *Proceedings of the 2014 ACM symposium on Principles of distributed computing*, pp. 272–281. ACM, New York (2014)
30. Kogan, D., Krauthgamer, R.: Sketching cuts in graphs and hypergraphs. In: *Proceedings of the 2015 Conference on Innovations in Theoretical Computer Science*, pp. 367–376. ACM, New York (2015)
31. Konrad, C.: Maximum matching in turnstile streams. *ESA* (2015)
32. McGregor, A.: Graph stream algorithms: a survey. *ACM SIGMOD Rec.* **43**(1), 9–20 (2014)
33. McGregor, A., Tench, D., Vorotnikova, S., Vu, H.T.: Densest subgraph in dynamic graph streams. *MFCS* (2015)
34. Muthukrishnan, S.: Data streams: algorithms and applications. *Theor. Comput. Sci.* **1**(2), 117–236 (2005)
35. Orenstein, Y., Ron, D.: Testing eulerianity and connectivity in directed sparse graphs. *Theor. Comput. Sci.* **412**(45), 6390–6408 (2011)
36. Parnas, M., Ron, D.: Testing the diameter of graphs. *Random Struct. Algorithms* **20**(2), 165–183 (2002)
37. Peng, P., Sohler, C.: Estimating graph parameters from random order streams. In: *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*, pp. 2449–2466. SIAM, Philadelphia (2018)
38. Price, E.: Efficient sketches for the set query problem. In: *Proceedings of the twenty-second annual ACM-SIAM symposium on Discrete Algorithms*, pp. 41–56. SIAM, Philadelphia (2011)

39. Ron, D.: Algorithmic and analysis techniques in property testing: foundations and trends ®. *Theor. Comput. Sci.* **5**(2), 73–205 (2010)
40. Rubinfeld, R., Shapira, A.: Sublinear time algorithms. *SIAM J. Discrete Math.* **25**(4), 1562–1588 (2011)
41. Sun, X., Woodruff, D.P.: Tight bounds for graph problems in insertion streams. In: *The 18th International Workshop on Approximation Algorithms for Combinatorial Optimization Problems (APPROX'2015)* (2015)
42. Verbin, E., Yu, W.: The streaming complexity of cycle counting, sorting by reversals, and other problems. In: *Proceedings of the twenty-second annual ACM-SIAM symposium on Discrete Algorithms*, pp. 11–25. SIAM, Philadelphia (2011)
43. Yoshida, Y., Ito, H.: Property testing on k-vertex-connectivity of graphs. In: *Automata, Languages and Programming*, pp. 539–550. Springer, Berlin (2008)

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Affiliations

Zengfeng Huang¹ · Pan Peng² 

✉ Pan Peng
p.peng@sheffield.ac.uk

Zengfeng Huang
huangzf@fudan.edu.cn

¹ School of Data Science, Fudan University, Shanghai, China

² Department of Computer Science, University of Sheffield, Sheffield, UK