

Testable Bounded Degree Graph Properties Are Random Order Streamable

Morteza Monemizadeh* S. Muthukrishnan† Pan Peng‡ Christian Sohler§

Abstract

We study which property testing and sublinear time algorithms can be transformed into graph streaming algorithms for random order streams. Our main result is that for bounded degree graphs, any property that is constant-query testable in the adjacency list model can be tested with *constant space* in a single-pass in random order streams. Our result is obtained by estimating the distribution of local neighborhoods of the vertices on a random order graph stream using constant space.

We then show that our approach can also be applied to constant time approximation algorithms for bounded degree graphs in the adjacency list model: As an example, we obtain a constant-space single-pass random order streaming algorithms for approximating the size of a maximum matching with additive error ϵn (n is the number of nodes).

Our result establishes for the first time that a large class of sublinear algorithms can be simulated in random order streams, while $\Omega(n)$ space is needed for many graph streaming problems for adversarial orders.

*Department of Computer Science, Goethe-Universität Frankfurt, Germany. Partially supported by DFG grants ME 2088/3-(1/2) and ME 2088/4-1. Email: monemi@ae.cs.uni-frankfurt.de.

†Rutgers University, Piscataway, NJ, USA. Email: muthu@cs.rutgers.edu.

‡Faculty of Computer Science, University of Vienna, Austria. The research leading to these results has received funding from the European Research Council under the European Union's Seventh Framework Programme (FP/2007-2013) / ERC Grant Agreement no. 340506. Email: pan.peng@univie.ac.at.

§Department of Computer Science, TU Dortmund, Germany. Supported by ERC Starting Grant 307696. Email: christian.sohler@tu-dortmund.de.

1 Introduction

Very large and complex networks abound. Some of the prominent examples are gene regulatory networks, health/disease networks, and online social networks like Facebook, Google+, LinkedIn and Twitter. The interconnectivity of neurons in human brain, relations in database systems, and chip designs are some further examples. Some of these networks can be quite large and it may be hard to store them completely in the main memory and some may be too large to be stored at all. However, these networks contain valuable information that we want to reveal. For example, social networks can provide insights into the structure of our society, and the structure in gene regulatory networks might yield insights into diseases. Thus, we need algorithms that can analyze the structure of these networks quickly.

One way to approach this problem is to design graph streaming algorithms [HRR98, AMS96]. A graph streaming algorithm gets access to a stream of edges in some order and exactly or approximately solves problems on the graph defined by the stream. The challenge is that a graph streaming algorithm should use space sublinear in the size of the graph. We will focus on algorithms that make only *one pass* over the graph stream, unless we explicitly say otherwise. It has been shown that many natural graph problems require $\Omega(n)$ space in the *adversarial order* model where n is the number of nodes in the graph and the edges can arrive in arbitrary order (see eg., [FKM⁺05, FKM⁺08]), and thus most of previous work has focused on the *semi-streaming* model, in which the algorithms are allowed to use $O(n \cdot \text{poly log } n)$ space. However, in many interesting applications, the graphs are sparse and so they can be fully stored in the semi-streaming model making this model useless in this setting. This raises the question *whether there are at least some natural conditions under which one can solve graph problems with space $o(n)$, possibly even $\log^{O(1)} n$ or constant.*

One such condition that recently received increasing attention is that the edges arrive in *random order*, i.e. in the order of a uniformly random permutation of the edges (e.g., [CCM08, KMM12, KKS14]). Uniformly random or near-uniformly random ordering is a natural assumption and can arise in many contexts. Indeed, previous work has shown that some problems that are hard for adversarial streams can be solved in the random order model. Konrad et al. [KMM12] gave single-pass semi-streaming algorithms for maximum matching for bipartite and general graphs with approximation ratio strictly larger than 1/2 in the random order semi-streaming model, while no such approximation algorithm is known in the adversary order model. Kapralov et al. [KKS14] gave a polylogarithmic approximation algorithm in polylogarithmic space for estimating the size of maximum matching of an unweighted graph in one pass over a random order stream. Assadi et al. [AKL17] recently showed that in the adversarial order and *dynamic* model where edges can be both inserted and deleted, any polylogarithmic approximation algorithm of maximum matching size requires $\tilde{\Omega}(n)$ space. On the other hand, Chakrabarti et al. [CCM08] presented an $\Omega(n)$ space lower bound for any single pass algorithm for graph connectivity in the random order streaming model, which is very close to the optimal $\Omega(n \log n)$ space lower bound in the adversarial order model [SW15]. In general, it is unclear which graph problems can be solved in random order streams using much smaller space than what is required for adversarially ordered streams.

An independent area of research is *property testing*, where with certain *query* access to an object (eg., random vertices or neighbors of a vertex for graphs), there are algorithms that can determine if the object satisfies a certain property, or is far from having such a property [RS96, GGR98, GR02]. The area of property testing has seen fundamental results, including testing various general graph properties. For example, it has been shown that many interesting properties (including connectivity, planarity, minor-freeness, hyperfiniteness) of bounded degree graphs can be tested with a constant number of queries [GR02, BSS10, NS13]. Another very related area of research is called *constant-time* (or in general, *sublinear-time*) *approximation* algorithms, where we are given query access

to an object (for example a graph) and the goal is to approximate the objective value of an optimal solution. For example, in bounded degree graphs, one can approximate the cost of the optimal solution with constant query complexity for some fundamental optimization problems (e.g., minimum spanning tree weight [CRT05], maximal matching size [NO08]; see also Section 1.3).

A fundamental question is if such results from property testing and constant-time approximation algorithms will lead to better graph streaming algorithms. Huang and Peng [HP16] recently considered the problem of estimating the minimum spanning tree weight and property testing for general graphs in dynamic and adversarial order model. They showed that a number of properties (e.g., connectivity, cycle-freeness) of general n -vertex graphs can be tested with space complexity $O(n^{1-\epsilon})$ and one can $(1 + \epsilon)$ -approximate the weight of minimum spanning tree with similar space guarantee. Furthermore, there exist $\Omega(n^{1-O(\epsilon)})$ space lower bounds for these problems that hold even in the insertion-only model [HP16].

1.1 Overview of Results

In this paper we provide a general framework that transforms bounded-degree graph property testing to very space-efficient random order streaming algorithms.

To formally state our main result, we first review some basic definitions of graph property testing. A *graph property* is a property that is invariant under graph isomorphism. Let $G = (V, E)$ be a graph with maximum degree upper bounded by a constant d , and we also call G a *d -bounded graph*. In the *adjacency list model* for (bounded-degree) graph property testing, we are given query access to the adjacency list of the input d -bounded graph $G = (V, E)$. That is, for any vertex $v \in V$ and index $i \leq d$, one can query the i th neighbor (if exists) of vertex v in constant time. Given a property Π , we are interested in testing if a graph G satisfies Π or is ϵ -far from satisfying Π while making as few queries as possible, where G is said to be ϵ -far from satisfying Π if one has to insert/delete more than ϵdn edges to make it satisfy Π . We call a property *constant-query testable* if there exists a testing algorithm (also called *tester*) for this property such that the number of performed queries depends only on parameters ϵ, d and is independent of the size of the input graph.

Given a graph property Π , we are interested in *approximately* testing it in a single-pass stream with a goal similar to the above. That is, the algorithm uses little space and with high constant probability, it accepts the input graph G if it satisfies P and rejects G if it is ϵ -far from satisfying P (see Section 4 for formal definitions). Our main result is as follows.

Theorem 1.1. *Any d -bounded graph property that is constant-query testable in the adjacency list model can be tested in the uniformly random order streaming model with constant space.*

To the best of our knowledge, this is the first non-trivial graph streaming algorithm with constant space complexity (measured in the number of *words*, where a word is a space unit large enough to encode an ID of any vertex in the graph.) By the constructions in [HP16], there exist graph properties (e.g., connectivity and cycle-freeness) of d -bounded graphs such that any single-pass streaming algorithm in the insertion-only and *adversary order* model must use $\Omega(n^{1-O(\epsilon)})$ space. In contrast to this lower bound, our main result implies that d -bounded connectivity and cycle-freeness can be tested in constant space in the random order stream model, since they are constant-query testable in the adjacency list model [GR02].

Our approach also works for simulating *constant-time approximation* algorithms as graph streaming algorithms with constant space. For a minimization (resp., maximization) optimization problem P and an instance I , we let $\text{OPT}(I)$ denote the value of some optimal solution of I . We call a value x an (α, β) -approximation for the problem P , if for any instance I , it holds that

$\text{OPT}(I) \leq x \leq \alpha \cdot \text{OPT}(I) + \beta$ (resp., $\frac{\text{OPT}(I)}{\alpha} - \beta \leq x \leq \text{OPT}(I)$). For example, it is known that there exists a constant-query algorithm for $(1, \varepsilon n)$ -approximating the maximal matching size of any n -vertex d -bounded graph [NO08]. That is, the number of queries made by the algorithm is independent of n and only depends on ε, d . As an application, we show:

Theorem 1.2. *Let $0 < \varepsilon < 1$ and d be constants. Then there exists an algorithm that uses constant space in the random order model, and with probability $2/3$, $(1, \varepsilon n)$ -approximates the size of some maximal matching in d -bounded graphs.*

We also remark that in a similar way, many other sublinear time algorithms for bounded degree graphs can be simulated in random order streams. Finally, our results can actually be extended to a model which requires weaker assumptions on the randomness of the order of edges in the stream, but we describe our results for the uniformly random order model, and leave the remaining details for later.

1.2 Technical Overview

The local neighborhood of depth k of a vertex v is the subgraph rooted at v and induced by all vertices of distance at most k from v . We call such a rooted subgraph a k -disc. Suppose that we are given a sufficiently large graph G whose maximum degree d is constant. This means that for any constant k , a k -disc centered at an arbitrary vertex v in G has constant size. Now assume that there exists an algorithm \mathcal{A} that, independent of the labeling of the vertices of G , accesses G by querying random vertices and exploring their k -discs. We observe that any constant-query property tester (see for example [GR11, CPS16]) falls within the framework of such an algorithm. If instead of the graph G we are given the distribution of k -discs of the vertices of G , we can use this distribution to simulate the algorithm \mathcal{A} and output with high probability the same result as executing the algorithm \mathcal{A} on G itself. Thus, the problem of developing constant-query property testers in random order streams can be reduced to the problem of designing streaming algorithms that approximate the distribution of k -discs in G .

The main technical contribution of this paper is an algorithm that given a random order stream S of edges of an underlying d -bounded degree graph G , approximates the distribution of k -discs of G up to an additive error of δ . We would like to mention that if the edges arrive in adversarial order, any algorithm that approximates the distribution of k -discs of G requires almost linear space [VY11, HP16], hence the assumption of random order streams (or something similar) is necessary to obtain our result.

Now in order to approximate the distribution of k -discs of the graph G we do the following. We proceed by sampling vertices uniformly at random and then perform a BFS for each sampled vertex using the arrival of edges along the stream S . Note that the new edges of the stream S that do not connect to the currently explored vertices are discarded. Let us call the k -disc that is observed by doing such a BFS from some vertex v to be Δ_1 . Due to possibility of missing edges during the BFS, this subgraph may be different from the true k -disc Δ_2 rooted at v .

If we are allowed to use two passes of the stream, then one can collect the k -disc of v in the first pass, and then verify if the collected disc is the true k -disc of v in the second pass (see Section 3.1). However, if we are restricted to a single pass, then it is more challenging to detect or verify if some edges have been missed in a collected disc. Fortunately, since the edges arrive in a uniformly random order, we can infer the conditional probability $\Pr[\Delta_1 | \Delta_2]$. That is, given the true rooted subgraph Δ_2 , we can compute the conditional probability of seeing a rooted subgraph Δ_1 in a random order stream when the true k -disc is Δ_2 .

We define the partial order on the set of k -discs given by $\Delta_1 \preceq \Delta_2$ whenever Δ_1 is a root-preserving isomorphic subgraph of Δ_2 . For every two k -discs Δ_1 and Δ_2 with $\Delta_1 \preceq \Delta_2$ we compute the conditional probability $\Pr[\Delta_1|\Delta_2]$. Using the set of all conditional probabilities $\Pr[\Delta_1|\Delta_2]$ we can estimate or approximate the distribution of k -discs of the graph G whose edges are revealed according to the stream S . In order to simplify the analysis of our algorithm, we require a natural independence condition for non-intersecting k -discs. Finally, we use the approximated distribution of k -discs to simulate the algorithm \mathcal{A} by the machinery that we explained above.

We remark that the idea of using a partial order to compute a distribution of k -discs in bounded degree graphs has first been used in [CPS16]. However, the setting in [CPS16] was quite different as it dealt with directed graphs where an edge can only be seen from one side (and the sample sizes required in that paper were only slightly sublinear in n).

1.3 Other Related Work

Feigenbaum et al. [FKSV02] initiated the study of property testing in streaming model, and they gave efficient testers for some properties of a sequence of data items (rather than graphs as we consider here). Bury and Schwiegelshohn [BS15] gave a lower bound of $n^{1-O(\varepsilon)}$ on the space complexity of any algorithm that $(1-\varepsilon)$ -approximates the size of maximum matching in adversarial streams. Kapralov et al. [KKS15] showed that in random streams, $\tilde{\Omega}(\sqrt{n})$ space is necessary to distinguish if a graph is bipartite or $1/2$ -far from being bipartite. Previous work has extensively studied streaming graph algorithms in both the insertion-only and dynamic models, see the recent survey [McG14].

In the framework of d -bounded graph property testing, it is now known that many interesting properties are constant-query testable in the adjacency list model, including k -edge connectivity, cycle-freeness, subgraph-freeness [GR02], k -vertex connectivity [YI08], minor-freeness [HKNO09, BSS10], matroids related properties [ITY12, TY15], hyperfinite properties [NS13], subdivision-freeness [KY13]. Constant-time approximation algorithms in d -bounded graphs are known to exist for a number of fundamental optimization problems, including $(1+\varepsilon)$ -approximating the weight of minimum spanning tree [CRT05], $(1, \varepsilon n)$ -approximating the size of maximal/maximum matching [NO08, YYI12], $(2, \varepsilon n)$ -approximating the minimum vertex cover size [PR07, MR09, ORRR12], $(O(\log d), \varepsilon n)$ -approximating the minimum dominating set size [PR07, NO08]. For d -bounded minor-free graphs, there are constant-time $(1, \varepsilon n)$ -approximation algorithms for the size of minimum vertex cover, minimum dominating set and maximum independent set [HKNO09].

2 Preliminaries

Let $G = (V, E)$ be an n -vertex graph with maximum degree upper bounded by some constant d , where we often identify V as $[n] := \{1, \dots, n\}$. We also call such a graph d -bounded graph. In this paper, we will assume the algorithms have the knowledge of n, d . We assume that G is represented as a sequence of edges, which we denote as $\text{STREAM}(G)$.

Graph k -discs. Let $k \geq 1$. The k -disc around a vertex v is the subgraph rooted at vertex v and induced by the vertices within distance at most k from v . Note that for an n -vertex graph, there are exactly n k -discs. Let $\mathcal{H}_{d,k} = \{\Delta_1, \dots, \Delta_N\}$ be the set of all k -disc isomorphism types, where $N = N_{d,k}$ is the number of all such types (and is thus a constant). In the following, we will refer to a k -disc of some vertex v in the graph G as $\text{disc}_{k,G}(v)$ and a k -disc type as Δ . Note that for every vertex v , there exists a unique k -disc type $\Delta \in \mathcal{H}_{d,k}$ such that $\text{disc}_{k,G}(v)$ is isomorphic to Δ , denoted as $\text{disc}_{k,G}(v) \cong \Delta$. (Throughout the paper, we call two rooted graphs H_1, H_2 isomorphic

to each other if there is a root-preserving mapping from the vertex set of H_1 to the vertex set of H_2 .)

We further assume that all the elements in $\mathcal{H}_{d,k}$ are ordered according to the natural partial order among k -disc types. More specifically, for any two k -disc types Δ_i, Δ_j , we let $\Delta_i \succcurlyeq \Delta_j$ (or equivalently, $\Delta_j \preccurlyeq \Delta_i$) denote that Δ_j is root-preserving isomorphic to some subgraph of Δ_i . Then we order all the k -disc types $\Delta_1, \dots, \Delta_N$ such that if $\Delta_i \succcurlyeq \Delta_j$, then $i \leq j$. Let $\mathcal{G}(j)$ denote all the indices i , except j itself, such that $\Delta_i \succcurlyeq \Delta_j$.

Locally random order streams. Let Σ_E denote the set of all permutations (or orderings) over the edge set E . Note that each $\sigma \in \Sigma_E$ determines the order of edges arriving from the stream. Let $\mathcal{D} = \mathcal{D}(\Sigma_E)$ denote a probability distribution over Σ_E . In particular, we let $\mathcal{U} = \mathcal{U}(\Sigma_E)$ denote the uniform distribution over Σ_E . Given a stream σ of edges, we define the *observed k -disc of v from the stream*, denoted as $\text{disc}_k(v, \sigma)$, to be the subgraph rooted at v and induced by all edges that are sequentially collected from the stream and the endpoints of which are within distance at most k to v . This is formally defined in the following algorithm `STREAM- k -DISC`.

Algorithm 1 The observed k -disc of v from the stream

```

1: procedure STREAM- $k$ -DISC(STREAM( $G$ ), $k$ , $v$ )
2:    $U \leftarrow \{v\}, \ell_v = 0, F \leftarrow \emptyset$ 
3:   for ( $u, w$ )  $\leftarrow$  next edge in the stream do
4:     if exactly one of  $u, w$ , say  $u$ , is contained in  $U$  then
5:       if  $\ell_u \leq k - 1$  then
6:          $U \leftarrow U \cup \{w\}, F \leftarrow F \cup \{(u, w)\}$ 
7:         for  $x \in U$  do
8:            $\ell_x \leftarrow$  the distance between  $x$  and  $v$  in the graph  $G' = (U, F)$ 
9:         end for
10:        end if
11:       else if both  $u, w$  are contained in  $U$  then
12:          $F \leftarrow F \cup \{(u, w)\}$ 
13:         for  $x \in U$  do
14:            $\ell_x \leftarrow$  the distance between  $x$  and  $v$  in the graph  $G' = (U, F)$ 
15:         end for
16:       end if
17:     end for
18:   return  $\text{disc}_k(v, \sigma) \leftarrow$  the subgraph rooted at  $v$  and induced by all edges in  $F$ 
19: end procedure

```

Now we formally define a locally random distribution on the order of edges.

Definition 2.1. Let $d, k > 0$. Let $G = (V, E)$ be a d -bounded graph. Let \mathcal{D} be a distribution over all the orderings of edges in E . Let $\Lambda_k = \{\lambda(\Delta_i | \Delta_j) : 0 \leq \lambda(\Delta_i | \Delta_j) \leq 1, \Delta_j \succcurlyeq \Delta_i, 1 \leq i, j \leq N\}$ be a set of real numbers in $[0, 1]$. We call \mathcal{D} a locally random Λ_k -distribution over G with respect to k -disc types, if for σ sampled from \mathcal{D} , the following conditions are satisfied:

1. (Conditional probabilities) For any vertex v with k -disc isomorphic to Δ_j , the probability that its observed k -disc $\text{disc}_k(v, \sigma) \cong \Delta_i$ is $\lambda(\Delta_i | \Delta_j)$, for any i such that $\Delta_j \succcurlyeq \Delta_i$.
2. (Independence of disjoint k -discs) For any two disjoint k -discs $\text{disc}_{k,G}(v)$ and $\text{disc}_{k,G}(u)$, their observed k -discs $\text{disc}_k(v, \sigma)$ and $\text{disc}_k(u, \sigma)$ are independent.

Note that the set Λ_k cannot be an arbitrary set, as there might be no distribution satisfying the above condition. On the other hand, if there indeed exists a distribution satisfying the condition with numbers in Λ_k , then we call the set Λ_k *realizable*. In the following, we call a stream a *locally random order stream* if there exists a family of realizable sets $\Lambda = \{\Lambda_k\}_{k \geq 1}$, such that the edge order is sampled from some locally random Λ_k -distribution with respect to k -disc types, for any integer $k \geq 1$. We have the following lemma.

Lemma 2.2. *Let $d \geq 1$. For any $k \geq 1$, there exists $n_0 = n_0(k, d)$, such that for $n \geq n_0$, any d -bounded n -vertex graph $G = (V, E)$, the uniform permutation \mathcal{U} over E is a locally random Λ_k -distribution over G with respect to k -disc types, for some realizable $\Lambda_k := \{\lambda(\Delta_i|\Delta_j) : 0 \leq \lambda(\Delta_i|\Delta_j) \leq 1, \Delta_j \succ \Delta_i, 1 \leq i, j \leq N\}$. Furthermore, if we let $\kappa := \max_{i, j: \Delta_j \succ \Delta_i} \frac{\lambda(\Delta_i|\Delta_j)}{\lambda(\Delta_i|\Delta_i)}$, $\lambda_{\min} := \min_{i \leq N} \lambda(\Delta_i|\Delta_i)$, then $\kappa \leq 2^{2d^{k+1}}$, $\lambda_{\min} \geq \frac{1}{(2d^{k+1})!}$.*

Proof. Note that for any vertex v with $\text{disc}_{k,G}(v) \cong \Delta_j$, the probability that the observed k -disc of v is isomorphic to Δ_i is exactly the fraction of orderings σ such that $\text{disc}_k(v, \sigma) \cong \Delta_i$, where $\Delta_j \succ \Delta_i$. We use such a fraction, which is a fixed real number, to define $\lambda(\Delta_i|\Delta_j)$. Observe that for an ordering σ sampled from \mathcal{U} , it directly satisfies the second condition Item 2 in Definition 2.1. Since there are at most $2d^{k+1}$ edges in any k -disc, the probability of observing a full k -disc is at least $\frac{1}{(2d^{k+1})!}$, that is, $\lambda_{\min} \geq \frac{1}{(2d^{k+1})!}$. Furthermore, since the k -disc type Δ_j might contain at most $\binom{|E(\Delta_j)|}{|E(\Delta_i)|} \leq 2^{2d^{k+1}}$ different subgraphs that are isomorphic to Δ_i , it holds that $\lambda(\Delta_i|\Delta_j) \leq \sum_{\substack{F: F \text{ subgraph of } \Delta_j \\ F \cong \Delta_i}} \lambda(\Delta_i|\Delta_i) \leq 2^{2d^{k+1}} \lambda(\Delta_i|\Delta_i)$ for any i, j such that $\Delta_j \succ \Delta_i$. This completes the proof of the lemma. \square

The above lemma shows that the uniformly random order stream is a special case of a locally random order stream. Another natural class of locally random order stream is ℓ -wise independent permutation of edges for any $\ell = \omega_n(1)$ (i.e., any function that tends to infinity as n goes to infinity) for n -vertex bounded degree graphs, but for our qualitative purposes here, it suffices to consider uniformly random order streams.

3 Approximating the k -Disc Type Distribution

In this section, we show how to approximate the distribution of k -disc types of any d -bounded graph in locally random order streams.

Recall that for any k, d , we let $N = N_{d,k}$ be the constant denoting the number of all possible k -disc isomorphism types. For any $i \leq N$, let V_i be the set of vertices from V with k -disc isomorphic to Δ_i in the input graph G , that is, $V_i := \{v|v \in V, \text{disc}_{k,G}(v) \cong \Delta_i\}$. Note that $f_i = \frac{|V_i|}{n}$ is the fraction of vertices with k -disc isomorphic to Δ_i .

3.1 A Two-Pass Algorithm

We start with a discussion of a two-pass algorithm for approximating the distribution of k -disc types. The main idea is that in the first pass we can collect or observe the k -disc from any vertex u , and then in the second pass, we check if the observed k -disc is the true k -disc of u or not. We can then use the statistics of the observed true k -discs to estimate the distribution of k -disc types.

Slightly more formally, we first sample a large constant number of vertices and let S denote the set of sampled vertices. Then in the first pass, for each vertex $u \in S$, we invoke the algorithm `STREAM- k -DISC` to collect the observed k -disc of u , denoted as H_u , from the stream. In the second

pass, for each vertex $w \in V(H_u)$, we collect all the incident edges to w . Then we let H'_u denote the subgraph spanned by all edges (collected in the second pass) incident to vertices within distance at most k to u . We check if H_u is isomorphic to H'_u . It is not hard to see that the true k -disc of u is observed if and only if H_u is isomorphic to H'_u . For each k -disc type Δ_i , we could then use the fraction of vertices v in S such that the true k -disc of v is observed and is isomorphic to Δ_i , to define an estimator for f_i . One should note that the naive estimator needs to be normalized appropriately by some probabilities and that there are dependencies between different variables, if one samples more than one starting vertex. Similar technical challenges also appear in our single-pass algorithm, for which we give detailed analysis in the following section. We omit further discussion on the two-pass algorithm here.

3.2 A Single-Pass Algorithm

In the following, we present our single-pass algorithm for approximating the distribution of k -disc types. We have the following lemma.

Lemma 3.1. *Let $G = (V, E)$ be a d -bounded graph presented in a locally random order stream defined by a Λ_k -distribution \mathcal{D} over G with respect to k -disc types, for some integer k . Let $\kappa := \max_{i,j:\Delta_j \succ \Delta_i} \frac{\lambda(\Delta_i|\Delta_j)}{\lambda(\Delta_i|\Delta_i)}$, $\lambda_{\min} := \min_{i \leq N} \lambda(\Delta_i|\Delta_i)$. Then for any constant $\delta > 0$, there exists a single-pass streaming algorithm that uses $O\left(\frac{\kappa^{2N} \cdot d^{3k+2} \cdot 3^{3N+1}}{\delta^2 \lambda_{\min}}\right)$ space, and with probability $\frac{2}{3}$, for any $i \leq N$, approximates the fraction f_i of vertices with k -disc isomorphic to Δ_i in G with additive error δ .*

Proof. Our algorithm is as follows. We first sample a constant number of vertices, which are called centers. Then for each center v , we collect the observed k -disc of v from the stream. Then we postprocess all the collected edges and use the corresponding empirical distribution of k -disc types of all centers to estimate the distribution of k -disc types of the input graph. The formal description is given in Algorithm 2.

Algorithm 2 Approximating the distribution of k -disc types

```

1: procedure  $k$ -DISC-DISTRIBUTION( $\text{STREAM}(G), \Lambda_k, n, d, k, \delta$ )
2:   sample a set  $A$  of  $s := \frac{8\kappa^{2N} \cdot d^{2k+1} \cdot 3^{3N+1}}{\delta^2 \lambda_{\min}}$  vertices uniformly at random
3:   for each  $v \in A$  do
4:      $H_v \leftarrow \text{STREAM}_k\text{-DISC}(\text{STREAM}(G), v, k)$  ▷ to collect observed  $k$ -disc of  $v$ 
5:   end for
6: end procedure
7:
8: procedure POSTPROCESSING
9:    $H \leftarrow$  the graph spanned by  $\cup_{v \in A} H_v$ 
10:  for  $i = 1$  to  $N$  do
11:     $Y_i \leftarrow |\{v : v \in A, \text{disc}_{k,H}(v) \cong \Delta_i\}|/s$ 
12:     $X_i \leftarrow (Y_i - \sum_{j \in \mathcal{G}(i)} X_j \cdot \lambda(\Delta_i|\Delta_j)) \cdot \lambda^{-1}(\Delta_i|\Delta_i)$ .
13:  end for
14:  return  $X_1, \dots, X_N$ 
15: end procedure

```

Note that since there are $s = \frac{8\kappa^{2N} \cdot d^{2k+1} \cdot 3^{3N+1}}{\delta^2 \lambda_{\min}}$ vertices in A and only edges that belong to the k -discs of these vertices will be collected by our algorithm, the space complexity of the algorithm is $O(sd^{k+1}) = O\left(\frac{\kappa^{2N} \cdot d^{3k+2} \cdot 3^{3N+1}}{\delta^2 \lambda_{\min}}\right)$, which is constant.

Now we show the correctness of the algorithm.

We let $A \sim \mathcal{U}_V$ denote that A is the set of s vertices sampled uniformly at random from V . For any $i \leq N$, let A_i be the set of vertices from A with k -disc isomorphic to Δ_i in the input graph G , that is, $A_i := \{v \mid v \in A, \text{disc}_{k,G}(v) \cong \Delta_i\}$. Note that $\mathbb{E}_{A \sim \mathcal{U}_V}[|A_i|] = s \cdot \frac{|V_i|}{n}$.

Let $\beta_i = 3^{i-N-2}$, $\theta_i = (3\kappa)^{i-N-1}$. By Chernoff bound and our setting of s which satisfy that $s \geq \Omega(\frac{1}{(\delta\theta_i)^2\beta_i})$, we have the following claim.

Claim 3.2. *For any $i \leq N$, $\Pr_{A \sim \mathcal{U}_V}[|\frac{|A_i|}{s} - \frac{|V_i|}{n}| \leq \delta\theta_i] \geq 1 - \beta_i$.*

We assume for now that A is a fixed set with s vertices. We let $\sigma \sim \mathcal{D}$ denote that the edge ordering σ is sampled from \mathcal{D} . For any $v \in A$, let $Z_{v,i}$ be the indicator random variable of the event that the observed k -disc $\text{disc}_k(v, \sigma)$ of v is isomorphic to Δ_i for $\sigma \sim \mathcal{D}$. Note that $\Pr_{\sigma \sim \mathcal{D}}[Z_{v,i} = 1] = \lambda(\Delta_i | \Delta_j)$ if $\text{disc}_{k,G}(v) \cong \Delta_j$. Let $Y_i^{(\sigma)} := \frac{|\{v: v \in A, \text{disc}_k(v, \sigma) \cong \Delta_i\}|}{s}$ denote the fraction of vertices in A with observed k -disc isomorphic to Δ_i . By definition, it holds that $Y_i^{(\sigma)} = \frac{1}{s} \sum_{\substack{v \in A_j \\ j \in \mathcal{G}(i) \cup \{i\}}} Z_{v,i}$, and furthermore, $\mathbb{E}_{\sigma \sim \mathcal{D}}[Y_i^{(\sigma)}] = \frac{1}{s} \sum_{j \in \mathcal{G}(i) \cup \{i\}} |A_j| \cdot \lambda(\Delta_i | \Delta_j)$. Let $X_i^{(\sigma)} = (Y_i^{(\sigma)} - \sum_{j \in \mathcal{G}(i)} X_j^{(\sigma)}) \cdot \lambda(\Delta_i | \Delta_j) \cdot \lambda^{-1}(\Delta_i | \Delta_i)$.

We have the following claim.

Claim 3.3. *For any $i \leq N$, it holds that $\mathbb{E}_{\sigma \sim \mathcal{D}}[X_i^{(\sigma)}] = \frac{|A_i|}{s}$.*

Proof. We prove the claim by induction. For $i = 1$, it holds that $\mathbb{E}_{\sigma \sim \mathcal{D}}[X_1^{(\sigma)}] = \mathbb{E}_{\sigma \sim \mathcal{D}}[Y_1^{(\sigma)}] \cdot \lambda^{-1}(\Delta_1 | \Delta_1) = \frac{|A_1|}{s} \cdot \lambda(\Delta_1 | \Delta_1) \cdot \lambda^{-1}(\Delta_1 | \Delta_1) = \frac{|A_1|}{s}$. Assuming that the claim holds for $i - 1$, and we prove it holds for i as well. By definition, we have that

$$\begin{aligned} \mathbb{E}_{\sigma \sim \mathcal{D}}[X_i^{(\sigma)}] &= \mathbb{E}_{\sigma \sim \mathcal{D}}[(Y_i^{(\sigma)} - \sum_{j \in \mathcal{G}(i)} X_j^{(\sigma)} \cdot \lambda(\Delta_i | \Delta_j)) \cdot \lambda^{-1}(\Delta_i | \Delta_i)] \\ &= \left(\sum_{j \in \mathcal{G}(i) \cup \{i\}} \frac{|A_j|}{s} \cdot \lambda(\Delta_i | \Delta_j) - \sum_{j \in \mathcal{G}(i)} \mathbb{E}_{\sigma \sim \mathcal{D}}[X_j^{(\sigma)}] \cdot \lambda(\Delta_i | \Delta_j) \right) \cdot \lambda^{-1}(\Delta_i | \Delta_i) \\ &= \left(\sum_{j \in \mathcal{G}(i) \cup \{i\}} \frac{|A_j|}{s} \cdot \lambda(\Delta_i | \Delta_j) - \sum_{j \in \mathcal{G}(i)} \frac{|A_j|}{s} \cdot \lambda(\Delta_i | \Delta_j) \right) \cdot \lambda^{-1}(\Delta_i | \Delta_i) = \frac{|A_i|}{s}, \end{aligned}$$

where the second to last equation follows from the induction. \square

We can now bound the variance of $Y_i^{(\sigma)}$ as shown in the following claim.

Claim 3.4. *For any $i \leq N$, it holds that $\text{Var}_{\sigma \sim \mathcal{D}}[Y_i^{(\sigma)}] \leq \frac{1}{s^2} \cdot d^{2k+1} \sum_{j \in \mathcal{G}(i) \cup \{i\}} |A_j| \cdot \lambda(\Delta_i | \Delta_j)$.*

Proof. Recall that $Y_i^{(\sigma)} = \frac{1}{s} \sum_{\substack{v \in A_j \\ j \in \mathcal{G}(i) \cup \{i\}}} Z_{v,i}$. Note that for each $v \in A$, by the independence assumption on \mathcal{D} , the random variable $Z_{v,i}$ can only correlate with the corresponding variables for vertices that are within distance at most $2k$ from v . The number of such vertices is at most $1 + d + d^2 + \dots + d^{2k} < d^{2k+1}$. Let $\text{dt}(u, v)$ denote the distance between u, v in the graph G . Then

we have that

$$\begin{aligned}
& \mathbb{E}_{\sigma \sim \mathcal{D}}[(\sum_{\substack{v \in A_j \\ j \in \mathcal{G}(i) \cup \{i\}}} Z_{v,i})^2] = \mathbb{E}_{\sigma \sim \mathcal{D}}[\sum_{\substack{v \in A_j \\ j \in \mathcal{G}(i) \cup \{i\}}} \sum_{\substack{u \in A_j \\ j \in \mathcal{G}(i) \cup \{i\}}} Z_{v,i} \cdot Z_{u,i}] \\
&= \mathbb{E}_{\sigma \sim \mathcal{D}}[\sum_{\substack{v \in A_j \\ j \in \mathcal{G}(i) \cup \{i\}}} (\sum_{\substack{u \in A_j \\ j \in \mathcal{G}(i) \cup \{i\} \\ \text{dt}_G(u,v) \leq 2k}} Z_{v,i} \cdot Z_{u,i} + \sum_{\substack{u \in A_j \\ j \in \mathcal{G}(i) \cup \{i\} \\ \text{dt}_G(u,v) > 2k}} Z_{v,i} \cdot Z_{u,i})] \\
&\leq \mathbb{E}_{\sigma \sim \mathcal{D}}[\sum_{\substack{v \in A_j \\ j \in \mathcal{G}(i) \cup \{i\}}} \sum_{\substack{u \in A_j \\ j \in \mathcal{G}(i) \cup \{i\} \\ \text{dt}_G(u,v) \leq 2k}} Z_{v,i}] + \left(\sum_{j \in \mathcal{G}(i) \cup \{i\}} [|A_j|] \cdot \lambda(\Delta_i | \Delta_j) \right)^2 \\
&\leq d^{2k+1} \mathbb{E}_{\sigma \sim \mathcal{D}}[\sum_{\substack{v \in A_j \\ j \in \mathcal{G}(i) \cup \{i\}}} Z_{v,i}] + (\mathbb{E}_{\sigma \sim \mathcal{D}}[\sum_{\substack{v \in A_j \\ j \in \mathcal{G}(i) \cup \{i\}}} Z_{v,i}])^2 \\
&= d^{2k+1} \cdot \sum_{j \in \mathcal{G}(i) \cup \{i\}} |A_j| \cdot \lambda(\Delta_i | \Delta_j) + (\mathbb{E}_{\sigma \sim \mathcal{D}}[\sum_{\substack{v \in A_j \\ j \in \mathcal{G}(i) \cup \{i\}}} Z_{v,i}])^2,
\end{aligned}$$

where the first inequality follows from the fact that $Z_{u,i} \leq 1$, and that for any two vertices u, v with $\text{dt}(u, v) > 2k$, $Z_{u,i}, Z_{v,i}$ are independent.

Then we have that

$$\begin{aligned}
& \text{Var}_{\sigma \sim \mathcal{D}}[Y_i^{(\sigma)}] = \frac{1}{s^2} \cdot \text{Var}_{\sigma \sim \mathcal{D}}[\sum_{\substack{v \in A_j \\ j \in \mathcal{G}(i) \cup \{i\}}} Z_{v,i}] \\
&= \frac{1}{s^2} \left(\mathbb{E}_{\sigma \sim \mathcal{D}}[(\sum_{\substack{v \in A_j \\ j \in \mathcal{G}(i) \cup \{i\}}} Z_{v,i})^2] - (\mathbb{E}_{\sigma \sim \mathcal{D}}[\sum_{\substack{v \in A_j \\ j \in \mathcal{G}(i) \cup \{i\}}} Z_{v,i}])^2 \right) \\
&\leq \frac{1}{s^2} \cdot d^{2k+1} \sum_{j \in \mathcal{G}(i) \cup \{i\}} |A_j| \cdot \lambda(\Delta_i | \Delta_j).
\end{aligned}$$

□

We next prove that each $X_i^{(\sigma)}$ is concentrated around its expectation with high probability.

Claim 3.5. *For any $i \leq N$, it holds that $\Pr_{\sigma \sim \mathcal{D}}[|X_i^{(\sigma)} - \mathbb{E}_{\sigma \sim \mathcal{D}}[X_i^{(\sigma)}]| \leq \theta_i \delta] \geq 1 - \beta_i$.*

Proof. We prove the claim by induction. For $i = 1$, it holds that

$$\begin{aligned}
& \Pr_{\sigma \sim \mathcal{D}}[|X_1^{(\sigma)} - \mathbb{E}_{\sigma \sim \mathcal{D}}[X_1^{(\sigma)}]| \leq \theta_1 \delta] \leq \Pr_{\sigma \sim \mathcal{D}}[|Y_1^{(\sigma)} - \mathbb{E}_{\sigma \sim \mathcal{D}}[Y_1^{(\sigma)}]| \cdot \lambda^{-1}(\Delta_1 | \Delta_1) \geq \delta \theta_1] \\
&\leq \frac{\text{Var}_{\sigma \sim \mathcal{D}}[Y_1^{(\sigma)}]}{(\delta \theta_1)^2 \cdot \lambda^2(\Delta_1 | \Delta_1)} \leq \frac{d^{2k+1} |A_1| \cdot \lambda(\Delta_1 | \Delta_1)}{s^2 \cdot (\delta \theta_1)^2 \cdot \lambda^2(\Delta_1 | \Delta_1)} \leq \frac{d^{2k+1}}{s(\delta \theta_1)^2 \cdot \lambda(\Delta_1 | \Delta_1)} \leq \beta_1,
\end{aligned}$$

where the last inequality follows from our choice of β_1, θ_1 and s which satisfy that $s \geq \frac{d^{2k+1}}{(\delta \theta_1)^2 \beta_1 \cdot \lambda(\Delta_1 | \Delta_1)}$. Now let us consider arbitrary $i \geq 2$, assuming that the claim holds for any $j \leq i - 1$. First, with

probability (over the randomness that $\sigma \sim \mathcal{D}$) at least $1 - \sum_{j=1}^{i-1} \beta_j = 1 - \sum_{j=1}^{i-1} 3^{j-N-2} \geq 1 - \frac{\beta_i}{2}$, it holds that for all $j \leq i-1$, $|X_j^{(\sigma)} - \mathbb{E}_{\sigma \sim \mathcal{D}}[X_j^{(\sigma)}]| \leq \theta_j \delta$. This further implies that with probability at least $1 - \frac{\beta_i}{2}$,

$$\begin{aligned} & \left| \sum_{j \in \mathcal{G}(i)} X_j^{(\sigma)} \cdot \frac{\lambda(\Delta_i | \Delta_j)}{\lambda(\Delta_i | \Delta_i)} - \mathbb{E}_{\sigma \sim \mathcal{D}} \left[\sum_{j \in \mathcal{G}(i)} X_j^{(\sigma)} \cdot \frac{\lambda(\Delta_i | \Delta_j)}{\lambda(\Delta_i | \Delta_i)} \right] \right| \\ & \leq \sum_{j \in \mathcal{G}(i)} |X_j^{(\sigma)} - \mathbb{E}_{\sigma \sim \mathcal{D}}[X_j^{(\sigma)}]| \cdot \frac{\lambda(\Delta_i | \Delta_j)}{\lambda(\Delta_i | \Delta_i)} \\ & \leq \sum_{j \in \mathcal{G}(i)} \delta \theta_j \cdot \frac{\lambda(\Delta_i | \Delta_j)}{\lambda(\Delta_i | \Delta_i)} \leq \kappa \cdot \sum_{j \in \mathcal{G}(i)} \delta \theta_j \leq \kappa \cdot \sum_{j=1}^{i-1} \delta (3\kappa)^{j-N} \leq \frac{\theta_i \delta}{2}. \end{aligned}$$

Now note that

$$\begin{aligned} & \Pr_{\sigma \sim \mathcal{U}} \left[|Y_i^{(\sigma)} - \mathbb{E}[Y_i^{(\sigma)}]| \cdot \lambda(\Delta_i | \Delta_i)^{-1} \geq \frac{\theta_i \delta}{2} \right] \leq \frac{4 \cdot \text{Var}_{\sigma \sim \mathcal{D}}[Y_i^{(\sigma)}]}{(\delta \theta_i)^2 \cdot \lambda(\Delta_i | \Delta_i)^2} \\ & \leq \frac{4 \cdot d^{2k+1} \sum_{j \in \mathcal{G}(i) \cup \{i\}} |A_j| \cdot \lambda(\Delta_i | \Delta_j)}{s^2 \cdot (\delta \theta_i)^2 \cdot \lambda(\Delta_i | \Delta_i)^2} \leq \frac{4 \cdot d^{2k+1} \cdot \kappa}{s \cdot (\delta \theta_i)^2 \cdot \lambda(\Delta_i | \Delta_i)} \leq \frac{\beta_i}{2}, \end{aligned}$$

where the last inequality follows from our choice of β_i, θ_i and s which satisfy that $s \geq \frac{8\kappa \cdot d^{2k+1}}{(\delta \theta_i)^2 \beta_i \cdot \lambda(\Delta_i | \Delta_i)}$.

Therefore, with probability (over $\sigma \sim \mathcal{D}$) at least $1 - \frac{\beta_i}{2} - \frac{\beta_i}{2} = 1 - \beta_i$, it holds that

$$\begin{aligned} & |X_i^{(\sigma)} - \mathbb{E}_{\sigma \sim \mathcal{D}}[X_i^{(\sigma)}]| \\ & = \left| \frac{Y_i^{(\sigma)} - \sum_{j \in \mathcal{G}(i)} X_j^{(\sigma)} \cdot \lambda(\Delta_i | \Delta_j)}{\lambda(\Delta_i | \Delta_i)} - \mathbb{E}_{\sigma \sim \mathcal{D}} \left[\frac{Y_i^{(\sigma)} - \sum_{j \in \mathcal{G}(i)} X_j^{(\sigma)} \cdot \lambda(\Delta_i | \Delta_j)}{\lambda(\Delta_i | \Delta_i)} \right] \right| \\ & = \left| \frac{(Y_i^{(\sigma)} - \mathbb{E}_{\sigma \sim \mathcal{D}}[Y_i^{(\sigma)}])}{\lambda(\Delta_i | \Delta_i)} - \left(\sum_{j \in \mathcal{G}(i)} X_j^{(\sigma)} \cdot \frac{\lambda(\Delta_i | \Delta_j)}{\lambda(\Delta_i | \Delta_i)} - \mathbb{E}_{\sigma \sim \mathcal{D}} \left[\sum_{j \in \mathcal{G}(i)} X_j^{(\sigma)} \cdot \frac{\lambda(\Delta_i | \Delta_j)}{\lambda(\Delta_i | \Delta_i)} \right] \right) \right| \\ & \leq \frac{\delta \theta_i}{2} + \frac{\delta \theta_i}{2} = \delta \theta_i. \end{aligned}$$

□

Now with probability (over both $A \sim \mathcal{U}_V$ and $\sigma \sim \mathcal{D}$) at least $1 - \beta_i - \beta_i$, it holds that

$$\begin{aligned} & \left| X_i^{(\sigma)} - \frac{|V_i|}{n} \right| \leq \left| X_i^{(\sigma)} - \mathbb{E}_{\sigma \sim \mathcal{D}}[X_i^{(\sigma)}] \right| + \left| \mathbb{E}_{\sigma \sim \mathcal{D}}[X_i^{(\sigma)}] - \frac{|V_i|}{n} \right| \\ & = \left| X_i^{(\sigma)} - \mathbb{E}_{\sigma \sim \mathcal{D}}[X_i^{(\sigma)}] \right| + \left| \frac{|A_i|}{s} - \frac{|V_i|}{n} \right| \leq \delta \theta_i + \delta \theta_i = 2\delta \theta_i. \end{aligned}$$

Finally, with probability at least $1 - 2 \sum_{j=1}^N \beta_j = 1 - 2 \sum_{j=1}^N 3^{j-N-2} \geq 1 - \frac{1}{3}$, it holds that for all $i \leq N$, $|X_i - \frac{|V_i|}{n}| \leq 2\theta_i \delta \leq \delta$. This completes the proof of the lemma. □

4 Constant-Space Property Testing

In this section, we show how to transform constant-query property testers in the adjacency list model to constant-space property testers in the random order stream model in a single pass and prove our main result Theorem 1.1. (Our transformation also works in the locally random order model as defined in Definition 2.1, but for simplicity, we only state our result in the uniformly random order model.)

Definition 4.1. Let $\Pi = (\Pi_n)_{n \in \mathbb{N}}$ be a property of d -bounded graphs, where Π_n is a property of graphs with n vertices. We say that Π is testable with query complexity q , if for every ε, d and n , there exists an algorithm that performs $q = q(n, \varepsilon, d)$ queries to the adjacency list of the graph, and with probability at least $2/3$, accepts any n -vertex d -bounded graph G satisfying Π , and rejects any n -vertex d -bounded graph that is ε -far from satisfying Π . If $q = q(\varepsilon, d)$ is a function independent of n , then we call Π constant-query testable.

Similarly, we can define constant-space testable properties in graph streams.

Definition 4.2. Let $\Pi = (\Pi_n)_{n \in \mathbb{N}}$ be a property of d -bounded graphs, where Π_n is a property of graphs with n vertices. We say that Π is testable with space complexity q , if for every ε, d and n , there exists an algorithm that performs a single pass over an edge stream of an n -vertex d -bounded graph G , uses $q = q(n, \varepsilon, d)$ space, and with probability at least $2/3$, accepts G if it satisfies Π , and rejects G if it is ε -far from satisfying Π . If $q = q(\varepsilon, d)$ is a function independent of n , then we call Π constant-space testable.

The proof of Theorem 1.1 is based on the following known fact: every constant-query property tester can be simulated by some canonical tester which only samples a constant number of vertices, and explores the k -discs of these vertices, and then makes deterministic decisions based on the explored subgraph. This implies that it suffices to approximate the distribution of k -disc types of the input graph to test the corresponding property. Formally, we will use the following lemma relating the constant-time testable properties and their k -disc distributions. For any graph G , let $S_{G,k}$ denote the subgraph spanned by the union of k -discs rooted at k uniformly sampled vertices from G . The following lemma is implied by Lemma 3.2 in [CPS16] (which was built on [GT03] and [GR11]). (The result in [CPS16] is stated for d -bounded directed graphs, while it also holds in the undirected case.)

Lemma 4.3. Let $\Pi = (\Pi_n)_{n \in \mathbb{N}}$ be any d -bounded graph property that is testable with $q = q(\varepsilon, d)$ query complexity in the adjacency list model. Then there exist integer n_0 , $k = c \cdot q$ for some large universal constant c , and an infinite sequence of $\mathcal{F} = \{\mathcal{F}_n\}_{n \geq n_0}$ such that for any $n \geq n_0$, \mathcal{F}_n is a set of graphs, each being a union of k disjoint k -discs, and for any n -vertex graph G ,

- if G satisfies Π_n , then with probability at most $\frac{5}{12}$, $S_{G,k}$ is isomorphic to one of the members in \mathcal{F}_n .
- if G is ε -far from satisfying Π_n , then with probability at least $\frac{7}{12}$, $S_{G,k}$ is isomorphic to one of the members in \mathcal{F}_n .

Now we are ready to give the proof of Theorem 1.1, which follows almost directly from the proof of Theorem 1.1 in [CPS16]. For the sake of completeness, we present the full proof here.

Proof of Theorem 1.1. Let $\Pi = (\Pi_n)_{n \in \mathbb{N}}$ be any property that is testable with query complexity $q = q(\varepsilon, d)$ in the adjacency list model. We set $k = c \cdot q$ and let \mathcal{F}_n be the set of graphs as

guaranteed in Lemma 4.3. Note that each subgraph $F = (\Gamma_1, \dots, \Gamma_k) \in \mathcal{F}_n$ is a multiset of k -discs. Set $N = N(d, k)$, $\delta = \frac{1}{48(2kN)^k}$. Let Λ_k be the set of probabilities as guaranteed in Lemma 2.2. Let $n_1 := n_1(d, k)$ be some sufficiently large constant.

Now let us describe our random order streaming algorithm for testing if an n -vertex d -bound graph G satisfies Π_n or is ε -far from satisfying Π_n . If $n < n_1$, we trivially test Π_n with constant space by storing the whole graph which contains at most $O(dn_1)$ edges. If $n \geq n_1$, we first invoke the algorithm k -DISC_DISTRIBUTION(STREAM(G), $\Lambda_k, n, d, k, \delta$) to get estimators X_1, \dots, X_N for the fraction f_1, \dots, f_N of vertices whose k -discs are isomorphic to $\Delta_1, \dots, \Delta_N$, respectively. Then for each $F = (\Gamma_1, \dots, \Gamma_k) \in \mathcal{F}_n$, we calculate its *empirical frequency* as $\Psi(F) = \frac{\prod_{i=1}^N \binom{X_i \cdot n}{x_i}}{\binom{n}{k}}$, where x_i is the number of copies among $\Gamma_1, \dots, \Gamma_k$ that are of the same type as Δ_i , for $1 \leq i \leq N$. Finally, we accept the graph if and only if $\sum_{F \in \mathcal{F}_n} \Psi(F) < \frac{1}{2}$.

Note that the space used by the algorithm is a constant. More precisely, the space complexity is $O(\max\{dn_1, \frac{\kappa^{2N} \cdot d^{3k+2} \cdot 3^{3N+1}}{\delta^2 \lambda_{\min}}\}) = O(\max\{dn_1, \frac{\kappa^{2N} \cdot d^{3k+2} \cdot 3^{3N+1} \cdot (2kN)^{2k}}{\lambda_{\min}}\})$, where the equation follows from Lemma 3.1 and our setting of δ .

Now we show the correctness of the algorithm. Note that we only need to consider the case that $n \geq n_1$. By Lemma 3.1, with probability at least $2/3$, it holds that for any $i \leq N$, $|X_i - f_i| \leq \delta$. In the following, we will condition on this event and we will prove that

- if G satisfies Π_n , then $\sum_{F \in \mathcal{F}_n} \Psi(F) < \frac{1}{2}$, and
- if G is ε -far from satisfying Π_n , then $\sum_{F \in \mathcal{F}_n} \Psi(F) \geq \frac{1}{2}$.

This would complete the proof.

For every $F = \{\Gamma_1, \dots, \Gamma_k\} \in \mathcal{F}_n$ and the relevant x_1, \dots, x_N , we will study $\psi(\Gamma_1, \dots, \Gamma_k) := \frac{\prod_{i=1}^N \binom{f_i \cdot n}{x_i}}{\binom{n}{k}}$, from which we will obtain the required bounds for $\sum_{F \in \mathcal{F}_n} \Psi(F)$.

Observe that for any multiset $\{\Gamma_1, \dots, \Gamma_k\}$, the probability that the k -discs of k vertices sampled uniformly at random (without replacement) span a subgraph isomorphic to the subgraph corresponding to $\{\Gamma_1, \dots, \Gamma_k\}$ has the *multivariate hypergeometric distribution* with parameters $n, f_1 \cdot n, \dots, f_N \cdot n, k$. That is, if for every $i \leq N$, there are exactly x_i copies in the multiset $\{\Gamma_1, \dots, \Gamma_k\}$ that are of the same isomorphic type as Γ_i (note that $x_1 + \dots + x_N = k$ for any $1 \leq i \leq N$), then the probability that the subgraph $S_{G,k}$ spanned by k -discs of k uniformly sampled vertices is isomorphic to $\{\Gamma_1, \dots, \Gamma_k\}$ is equal to $\psi(\Gamma_1, \dots, \Gamma_k) = \frac{\prod_{i=1}^N \binom{f_i \cdot n}{x_i}}{\binom{n}{k}}$, where we assumed $\binom{L}{M} = 0$ for $L < M$.

To study the relation between $\Psi(F)$ and $\psi(F)$, we begin with the following auxiliary claim.

Claim 4.4. *For any i , if $|X_i - f_i| \leq \delta$, it holds that $|\binom{X_i \cdot n}{x_i} - \binom{f_i \cdot n}{x_i}| \leq 4\delta n^{x_i}$.*

Proof. Let us first observe that the inequality trivially holds for $x_i = 0$, and it also easily holds for $x_i = 1$: $|\binom{X_i \cdot n}{x_i} - \binom{f_i \cdot n}{x_i}| = |X_i \cdot n - f_i \cdot n| \leq \delta n \leq 4\delta n^{x_i}$. Therefore, let us assume now that $x_i \geq 2$.

Let us recall a binomial identity: $\binom{L}{M} = \sum_{K=M-1}^{L-1} \binom{K}{M-1}$, which gives for $M \leq J \leq L$ the

following: $\binom{L}{M} = \binom{J}{M} + \sum_{K=J}^{L-1} \binom{K}{M-1}$. Using this identity, that $f_i \leq 1$, and $x_i \geq 2$, we obtain,

$$\begin{aligned}
\binom{X_i \cdot n}{x_i} &\leq \binom{f_i \cdot n + \lceil \delta n \rceil}{x_i} = \binom{f_i \cdot n}{x_i} + \sum_{j=1}^{\lceil \delta n \rceil} \binom{f_i \cdot n + j - 1}{x_i - 1} \\
&\leq \binom{f_i \cdot n}{x_i} + \lceil \delta n \rceil \cdot \binom{f_i \cdot n + \lceil \delta n \rceil - 1}{x_i - 1} \\
&\leq \binom{f_i \cdot n}{x_i} + 2\delta n (f_i \cdot n + \delta n)^{x_i - 1} = \binom{f_i \cdot n}{x_i} + 2\delta n ((1 + \delta)n)^{x_i - 1} \\
&= \binom{f_i \cdot n}{x_i} + 2\delta(1 + \delta)^{x_i - 1} n^{x_i} \leq \binom{f_i \cdot n}{x_i} + 4\delta n^{x_i} ,
\end{aligned}$$

where in the last inequality, we used the fact that $(1 + \delta)^{x_i - 1} \leq (1 + \delta)^k \leq 2$.

Similarly, if $f_i \cdot n \geq \lceil \delta n \rceil + k$, we have $f_i \cdot n \geq \lceil \delta n \rceil + x_i$, and we obtain,

$$\begin{aligned}
\binom{X_i}{x_i} &\geq \binom{f_i \cdot n - \lceil \delta n \rceil}{x_i} = \binom{f_i \cdot n}{x_i} - \sum_{j=1}^{\lceil \delta n \rceil} \binom{f_i \cdot n - j}{x_i - 1} \\
&\geq \binom{f_i \cdot n}{x_i} - \lceil \delta n \rceil \binom{f_i \cdot n}{x_i - 1} \geq \binom{f_i \cdot n}{x_i} - 2\delta n \binom{n}{x_i - 1} \\
&\geq \binom{f_i \cdot n}{x_i} - 2\delta n \cdot n^{x_i - 1} = \binom{f_i \cdot n}{x_i} - 2\delta n^{x_i} .
\end{aligned}$$

On the other hand, if $f_i \cdot n \leq \lceil \delta n \rceil + k$, we note that $\binom{f_i \cdot n}{x_i} \leq \binom{\lceil \delta n \rceil + k}{x_i} \leq (\lceil \delta n \rceil + k)^{x_i} \leq (2\delta n)^{x_i} \leq 4\delta n^{x_i}$, where the third inequality follows from the fact that $n \geq n_1$ and that n_1 is a sufficiently large constant. Therefore since $\binom{X_i}{x_i} \geq 0$, we have $\binom{X_i}{x_i} \geq \binom{f_i \cdot n}{x_i} - 4\delta n^{x_i}$.

Now we can combine all the bounds above and obtain that for $x_2 \geq 2$, the following holds,

$$\binom{f_i \cdot n}{x_i} - 4\delta n^{x_i} \leq \binom{X_i}{x_i} \leq \binom{f_i \cdot n}{x_i} + 4\delta n^{x_i} ,$$

what yields the claim. □

Next, consider any $F = \{\Gamma_1, \dots, \Gamma_k\}$ and the corresponding frequencies x_1, \dots, x_N . Note that there are at most k indices i with $x_i > 0$, and that $x_1 + \dots + x_N = k$. Let $\mathcal{I} = \{i : x_i > 0, 1 \leq i \leq N\}$ and thus $|\mathcal{I}| \leq k$ and $\prod_{i \in \mathcal{I}} n^{x_i} = n^k$. We have the following auxiliary claim.

Claim 4.5. *For any i , conditioned on $|X_i - f_i| \leq \delta$, the following inequalities hold:*

$$\begin{aligned}
\prod_{i \in \mathcal{I}} \left(\binom{f_i \cdot n}{x_i} + 4\delta n^{x_i} \right) &< \prod_{i \in \mathcal{I}} \binom{f_i \cdot n}{x_i} + 4\delta 2^k n^k , \\
\prod_{i \in \mathcal{I}} \left(\binom{f_i \cdot n}{x_i} - 4\delta n^{x_i} \right) &> \prod_{i \in \mathcal{I}} \binom{f_i \cdot n}{x_i} - 4\delta 2^k n^k .
\end{aligned}$$

Proof. For any $i \in \mathcal{I}$, we let $y_{i,0} = \binom{f_i \cdot n}{x_i}$ and $y_{i,1} = 4\delta n^{x_i}$. Then

$$\begin{aligned} \prod_{i \in \mathcal{I}} \left(\binom{f_i \cdot n}{x_i} + 4\delta n^{x_i} \right) &= \prod_{i \in \mathcal{I}} (y_{i,0} + y_{i,1}) = \sum_{i \in \mathcal{I}, j_i \in \{0,1\}} \prod_{i \in \mathcal{I}} y_{i,j_i} \\ &= \prod_{i \in \mathcal{I}} y_{i,0} + \sum_{\substack{i \in \mathcal{I}, j_i \in \{0,1\}, \\ \text{there exists } j_i = 1}} \prod_{i \in \mathcal{I}} y_{i,j_i} \\ &= \prod_{i \in \mathcal{I}} \binom{f_i \cdot n}{x_i} + \sum_{\substack{i \in \mathcal{I}, j_i \in \{0,1\}, \\ \text{there exists } j_i = 1}} \prod_{i \in \mathcal{I}} y_{i,j_i} . \end{aligned}$$

Now note that for any $i \in \mathcal{I}$, $y_{i,0} = \binom{f_i \cdot n}{x_i} \leq n^{x_i}$. Therefore, for any sequence $\{j_i\}_{i \in \mathcal{I}}$ with at least one element equal to 1, we have the following bound $\prod_{i \in \mathcal{I}} y_{i,j_i} \leq 4\delta \prod_{i \in \mathcal{I}} n^{x_i} = 4\delta n^k$. Since the total number of such indices is $2^k - 1 < 2^k$, we have

$$\prod_{i \in \mathcal{I}} \binom{f_i \cdot n}{x_i} + \sum_{\substack{i \in \mathcal{I}, j_i \in \{0,1\}, \\ \text{there exists } j_i = 1}} \prod_{i \in \mathcal{I}} y_{i,j_i} < \prod_{i \in \mathcal{I}} \binom{f_i \cdot n}{x_i} + 4\delta n^k \cdot 2^k ,$$

which completes the proof of the first inequality. The proof of the second inequality is analogues. \square

Using Claims 4.4 and 4.5, we can prove the following relation between $\Psi(F)$ and $\psi(F)$.

Claim 4.6. *If $|X_i - f_i| \leq \delta$ for every i , then $|\Psi(F) - \psi(F)| \leq 4\delta(2k)^k$ for every $F \in \mathcal{F}_n$.*

Proof. Let $F = \{\Gamma_1, \dots, \Gamma_k\} \in \mathcal{F}_n$. By Claims 4.4 and 4.5, we have

$$\begin{aligned} \Psi(\Gamma_1 \dots \Gamma_k) &= \frac{\prod_{i \in \mathcal{I}} \binom{X_i \cdot n}{x_i}}{\binom{n}{k}} \leq \frac{\prod_{i \in \mathcal{I}} \left(\binom{f_i \cdot n}{x_i} + 4\delta n^{x_i} \right)}{\binom{n}{k}} < \frac{\prod_{i \in \mathcal{I}} \binom{f_i \cdot n}{x_i} + 4\delta 2^k n^k}{\binom{n}{k}} \\ &\leq \psi(\Gamma_1, \dots, \Gamma_k) + 4\delta(2k)^k , \end{aligned}$$

where the last inequality follows from that $\binom{n}{k} \geq \left(\frac{n}{k}\right)^k$. Similarly, by Claims 4.4 and 4.5, we have,

$$\begin{aligned} \Psi(\Gamma_1, \dots, \Gamma_k) &\geq \frac{\prod_{i \in \mathcal{I}} \left(\binom{f_i \cdot n}{x_i} - 4\delta n^{x_i} \right)}{\binom{n}{k}} \geq \frac{\prod_{i \in \mathcal{I}} \binom{f_i \cdot n}{x_i} - 4\delta 2^k n^k}{\binom{n}{k}} \\ &\geq \psi(\Gamma_1, \dots, \Gamma_k) - 4\delta(2k)^k . \end{aligned}$$

\square

Now consider the case that G satisfies II. Then, by Lemma 4.3, with probability at most $\frac{5}{12}$, the subgraph $S_{G,k}$ spanned by the k -discs of k vertices that are sampled uniformly at random without replacement is isomorphic to some member in \mathcal{F}_n , that is, $\sum_{F \in \mathcal{F}_n} \psi(F) \leq \frac{5}{12}$. Therefore, by Claim 4.6, we have,

$$\begin{aligned} \sum_{F \in \mathcal{F}_n} \Psi(F) &< \sum_{F \in \mathcal{F}_n} \psi(F) + \sum_{F \in \mathcal{F}_n} 4\delta(2k)^k \leq \sum_{F \in \mathcal{F}_n} \psi(F) + N_{d,k}^k \cdot 4\delta(2k)^k \\ &\leq \frac{5}{12} + \frac{1}{12} = \frac{1}{2} . \end{aligned}$$

Similarly, by Lemma 4.3, if G is ε -far from satisfying Π , then with probability at least $\frac{7}{12}$, the k -discs rooted at k vertices that are sampled uniformly at random span a subgraph in \mathcal{F}_n . Hence, Claim 4.6 gives

$$\begin{aligned} \sum_{F \in \mathcal{F}_n} \Psi(F) &\geq \sum_{F \in \mathcal{F}_n} \psi(F) - \sum_{F \in \mathcal{F}_n} 4\delta(2k)^k \geq \sum_{F \in \mathcal{F}_n} \psi(F) - N_{d,k}^k \cdot 4\delta(2k)^k \\ &\geq \frac{7}{12} - \frac{1}{12} = \frac{1}{2}. \end{aligned}$$

These inequalities conclude the analysis of our algorithm and the proof of Theorem 1.1. \square

5 Constant-Time Approximation Algorithms

As we mentioned in the introduction, to simulate any constant-time algorithm that is independent of the labeling of the vertices, and accesses the graph by sampling random vertices and exploring neighborhoods (or k -discs for some k) of these vertices, it suffices to have the distribution of k -disc types. Now we explain slightly more about this simulation and sketch the proof of Theorem 1.2. In order to approximate the size of the solution of an optimization problem (e.g., maximum matching, minimum vertex cover), it has been observed by Parnas and Ron [PR07] that it suffices to have efficient oracle \mathcal{O}_S access to a solution S . This is true since one can attain a good estimator for the size of S by sampling a constant number of vertices, performing corresponding queries to the oracle \mathcal{O}_S and then returning the fraction of vertices that belong to S based on the returned answers from \mathcal{O}_S . Nguyen and Onak [NO08] implemented such an oracle via an elegant approach of locally simulating the classical greedy algorithm. In particular, they showed the following result.

Lemma 5.1 ([NO08]). *There exist $q = q(\varepsilon, d)$, an oracle \mathcal{O}_M to a maximal matching M , and an algorithm that queries \mathcal{O}_M about all the edges incident to a set of $s = O(1/\varepsilon^2)$ randomly sampled vertices and with probability at least $2/3$, returns an estimator that is $(1, \varepsilon n)$ -approximation of the size of M , and each query to \mathcal{O}_M performs at most q queries to the adjacency list of the graph.*

A key observation is that the algorithm in Lemma 5.1 can be viewed as first sampling s q -discs from the graph and then perform \mathcal{O}_M queries on each of these q -discs. It is easy to see that with high probability 0.99, all these q -discs are disjoint. Furthermore, the answer of the above oracle only depends on the structure of the corresponding neighborhood of the starting vertex v and the random ordering of the edges belonging to this neighborhood.

Now we can approximate the size of a maximal matching in the random order streaming model as follows: we first invoke Algorithm 2 to get an estimator for the distribution of q -discs. Then we can simulate the oracle on this distribution.

Acknowledgment

We would like to thank G. Cormode, H. Jowhari for helpful discussions.

References

- [AKL17] Sepehr Assadi, Sanjeev Khanna, and Yang Li. On estimating maximum matching size in graph streams. In *Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 1723–1742. SIAM, 2017.

- [AMS96] Noga Alon, Yossi Matias, and Mario Szegedy. The space complexity of approximating the frequency moments. In *Proceedings of the twenty-eighth annual ACM symposium on Theory of computing*, pages 20–29. ACM, 1996.
- [BS15] Marc Bury and Chris Schwiegelshohn. Sublinear estimation of weighted matchings in dynamic data streams. In *Algorithms-ESA 2015*, pages 263–274. Springer, 2015.
- [BSS10] Itai Benjamini, Oded Schramm, and Asaf Shapira. Every minor-closed property of sparse graphs is testable. *Advances in mathematics*, 223(6):2200–2218, 2010.
- [CCM08] Amit Chakrabarti, Graham Cormode, and Andrew McGregor. Robust lower bounds for communication and stream computation. In *Proceedings of the fortieth annual ACM symposium on Theory of computing*, pages 641–650. ACM, 2008.
- [CPS16] Artur Czumaj, Pan Peng, and Christian Sohler. Relating two property testing models for bounded degree directed graphs. In *Proceedings of the 48th Annual ACM SIGACT Symposium on Theory of Computing*, STOC 2016, pages 1033–1045. ACM, 2016.
- [CRT05] Bernard Chazelle, Ronitt Rubinfeld, and Luca Trevisan. Approximating the minimum spanning tree weight in sublinear time. *SIAM J. Comput.*, 34(6):1370–1379, 2005.
- [FKM⁺05] Joan Feigenbaum, Sampath Kannan, Andrew McGregor, Siddharth Suri, and Jian Zhang. On graph problems in a semi-streaming model. *Theoretical Computer Science*, 348(2-3):207–216, 2005.
- [FKM⁺08] Joan Feigenbaum, Sampath Kannan, Andrew McGregor, Siddharth Suri, and Jian Zhang. Graph distances in the data-stream model. *SIAM Journal on Computing*, 38(5):1709–1727, 2008.
- [FKSV02] Joan Feigenbaum, Sampath Kannan, Martin Strauss, and Mahesh Viswanathan. Testing and spot-checking of data streams. *Algorithmica*, 34(1):67–80, 2002.
- [GGR98] Oded Goldreich, Shafi Goldwasser, and Dana Ron. Property testing and its connection to learning and approximation. *J. ACM*, 45(4):653–750, 1998.
- [GR02] Oded Goldreich and Dana Ron. Property testing in bounded degree graphs. *Algorithmica*, 32:302–343, 2002.
- [GR11] Oded Goldreich and Dana Ron. On proximity-oblivious testing. *SIAM Journal on Computing*, 40(2):534–566, 2011.
- [GT03] Oded Goldreich and Luca Trevisan. Three theorems regarding testing graph properties. *Random Structures & Algorithms*, 23(1):23–57, 2003.
- [HKNO09] Avinatan Hassidim, Jonathan A Kelner, Huy N Nguyen, and Krzysztof Onak. Local graph partitions for approximation and testing. In *Foundations of Computer Science, 2009. FOCS'09. 50th Annual IEEE Symposium on*, pages 22–31. IEEE, 2009.
- [HP16] Zengfeng Huang and Pan Peng. Dynamic graph stream algorithms in $o(n)$ space. In *43rd International Colloquium on Automata, Languages, and Programming, ICALP 2016, July 11-15, 2016, Rome, Italy*, pages 18:1–18:16, 2016.

- [HRR98] Monika Rauch Henzinger, Prabhakar Raghavan, and Sridhar Rajagopalan. Computing on data streams. In *External Memory Algorithms, Proceedings of a DIMACS Workshop, New Brunswick, New Jersey, USA, May 20-22, 1998*, pages 107–118, 1998.
- [ITY12] Hiro Ito, Shin-Ichi Tanigawa, and Yuichi Yoshida. Constant-time algorithms for sparsity matroids. In *International Colloquium on Automata, Languages, and Programming*, pages 498–509. Springer, 2012.
- [KKS14] Michael Kapralov, Sanjeev Khanna, and Madhu Sudan. Approximating matching size from random streams. In *Proceedings of the Twenty-Fifth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 734–751. Society for Industrial and Applied Mathematics, 2014.
- [KKS15] Michael Kapralov, Sanjeev Khanna, and Madhu Sudan. Streaming lower bounds for approximating max-cut. In *Proceedings of the Twenty-Sixth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 1263–1282. Society for Industrial and Applied Mathematics, 2015.
- [KMM12] Christian Konrad, Frédéric Magniez, and Claire Mathieu. Maximum matching in semi-streaming with few passes. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*, pages 231–242. Springer, 2012.
- [KY13] Ken-ichi Kawarabayashi and Yuichi Yoshida. Testing subdivision-freeness: property testing meets structural graph theory. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*, pages 437–446. ACM, 2013.
- [McG14] Andrew McGregor. Graph stream algorithms: a survey. *ACM SIGMOD Record*, 43(1):9–20, 2014.
- [MR09] Sharon Marko and Dana Ron. Approximating the distance to properties in bounded-degree and general sparse graphs. *ACM Transactions on Algorithms (TALG)*, 5(2):22, 2009.
- [NO08] Huy N Nguyen and Krzysztof Onak. Constant-time approximation algorithms via local improvements. In *Foundations of Computer Science, 2008. FOCS'08. IEEE 49th Annual IEEE Symposium on*, pages 327–336. IEEE, 2008.
- [NS13] Ilan Newman and Christian Sohler. Every property of hyperfinite graphs is testable. *SIAM Journal on Computing*, 42(3):1095–1112, 2013.
- [ORRR12] Krzysztof Onak, Dana Ron, Michal Rosen, and Ronitt Rubinfeld. A near-optimal sublinear-time algorithm for approximating the minimum vertex cover size. In *Proceedings of the twenty-third annual ACM-SIAM symposium on Discrete Algorithms*, pages 1123–1131. SIAM, 2012.
- [PR07] Michal Parnas and Dana Ron. Approximating the minimum vertex cover in sublinear time and a connection to distributed algorithms. *Theoretical Computer Science*, 381(1):183–196, 2007.
- [RS96] Ronitt Rubinfeld and Madhu Sudan. Robust characterizations of polynomials with applications to program testing. *SIAM Journal on Computing*, 25(2):252–271, 1996.

- [SW15] Xiaoming Sun and David P Woodruff. Tight bounds for graph problems in insertion streams. In *LIPICs-Leibniz International Proceedings in Informatics*, volume 40. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2015.
- [TY15] Shin-Ichi Tanigawa and Yuichi Yoshida. Testing the supermodular-cut condition. *Algorithmica*, 71(4):1065–1075, 2015.
- [VY11] Elad Verbin and Wei Yu. 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*, pages 11–25. SIAM, 2011.
- [YI08] Yuichi Yoshida and Hiro Ito. Property testing on k-vertex-connectivity of graphs. In *Automata, Languages and Programming*, pages 539–550. Springer, 2008.
- [YYI12] Yuichi Yoshida, Masaki Yamamoto, and Hiro Ito. Improved constant-time approximation algorithms for maximum matchings and other optimization problems. *SIAM J. Comput.*, 41(4):1074–1093, 2012.