

Universally Optimal Streaming Algorithm for Random Walks in Dense Graphs

Klim Efremenko  

Ben-Gurion University, Be'er-Sheva, Israel

Gillat Kol  

Princeton University, NJ, US

Raghuvansh R. Saxena 

Tata Institute of Fundamental Research, Mumbai, India

Zhijun Zhang  

INSAIT, Sofia University “St. Kliment Ohridski”, Bulgaria

Abstract

Sampling a random walk is a fundamental primitive in many graph applications. In the streaming model, it is known that sampling an L -step random walk on an n -vertex directed graph requires $\Omega(nL)$ space, implying that no sublinear-space streaming algorithm exists for general graphs.

We show that sublinear algorithms are possible for the case of dense graphs, where every vertex has out-degree at least $\Omega(n)$. In particular, we give a one-pass turnstile streaming algorithm that uses only $\tilde{O}(L)$ memory for such graphs. More broadly, for graphs with minimum out-degree at least d , our streaming algorithm samples a random walk using $\tilde{O}(\frac{n}{d} \cdot L)$ memory.

We show that our algorithm is optimal in a strong “beyond worst-case” sense. To formalize this, we introduce the notion of universal optimality for graph streaming algorithms. Informally, a streaming algorithm is universally optimal if it performs (almost) as well as possible on every graph, assuming a worst-case choice of the streaming order. This notion of universal optimality is a key conceptual contribution of our work.

2012 ACM Subject Classification Theory of computation \rightarrow Streaming, sublinear and near linear time algorithms

Keywords and phrases Random Walk, streaming Algorithm, universal Optimality

Digital Object Identifier 10.4230/LIPIcs.ITCS.2026.55

Funding *Klim Efremenko*: Supported by the Israel Science Foundation (ISF) through grant No. 1456/18 and European Research Council Grant number: 949707.

Gillat Kol: Supported by a National Science Foundation CAREER award CCF-1750443 and by a BSF grant No. 2018325.

Raghuvansh R. Saxena: Supported by the Department of Atomic Energy, Government of India, under project no. RTI4001.

Zhijun Zhang: This research was partially done while at Princeton University and was partially funded by the Ministry of Education and Science of Bulgaria (support for INSAIT, part of the Bulgarian National Roadmap for Research Infrastructure).

1 Introduction

1.1 Background and Motivation

Sampling a random-walk in a graph is a fundamental problem with broad algorithmic applications. For instance, it underlies techniques for connectivity testing [17], clustering [2, 3, 6, 21], sampling [13], constructing random spanning trees [20], and approximate counting [12]. Another well-known application of random walk sampling is for *PageRank* estimation [19], where the rank of a webpage corresponds to the likelihood that a user,



© Klim Efremenko, Gillat Kol, Raghuvansh R. Saxena, and Zhijun Zhang; licensed under Creative Commons License CC-BY 4.0

17th Innovations in Theoretical Computer Science Conference (ITCS 2026).

Editor: Shubhangi Saraf; Article No. 55; pp. 55:1–55:20



Leibniz International Proceedings in Informatics

Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

randomly following hyperlinks, arrives at that page¹. Since these applications often involve extremely large, real-world networks, an important challenge is to develop space-efficient *graph-streaming algorithms* that process the graph edge-by-edge rather than storing it explicitly.

The folklore algorithm. A well-known folklore streaming algorithm for simulating an L -step random walk on an n -vertex directed graph from a given start vertex v_0 uses $\tilde{O}(nL)$ memory². The algorithm works as follows: As the graph's edges arrive in a stream, maintain for each vertex v a list of L outgoing edges selected independently and uniformly at random. To generate the walk, begin at v_0 and, at step $i \in [L]$, move to a new vertex v_i by taking the first sampled edge that haven't been taken before from the precomputed list of outgoing edges at v_{i-1} .³

The folklore algorithm was later shown to be optimal [14]; however, the hard instances used by the lower-bound proof crucially involve vertices of very low out-degree. This naturally raises the question of whether better algorithms might exist for *dense graphs*, which we define as graphs in which every vertex has out-degree at least d for some large $d \geq 1$.⁴ We use the *minimum out-degree* d as our measure of density rather than, say, the average degree (or equivalently, the total number of edges), as, for instance, the graph formed by the union of a clique and a sparse subgraph may still have a high average out-degree. Yet, if v_0 is in the sparse subgraph, then sampling a random walk starting from v_0 requires sampling a random walk in the sparse subgraph.

1.2 Our Results

1.2.1 Random Walk Sampling for Dense Graphs

The algorithm $\pi_{d,L}^{\text{sparse}}$. We begin by noting that if the input graph has minimum out-degree at least d and $\frac{L}{d}$ is large, the folklore algorithm can be improved to store only $\tilde{O}(\lceil \frac{L}{d} \rceil \cdot n)$ edges. Specifically, instead of storing L outgoing edges per vertex, $\pi_{d,L}^{\text{sparse}}$ stores $\tilde{O}(\lceil \frac{L}{d} \rceil)$ outgoing edges for each vertex. We show that this is enough because, with high probability, no vertex is visited more than that many times during an L -step random walk on a graph with minimum out-degree at least d . We call this algorithm $\pi_{d,L}^{\text{sparse}}$.

The *distortion* of a random walk sampling algorithm is the maximum, over all graphs G , of the statistical distance between the true distribution of random walk on G and the distribution over outputs produced by the algorithm on input G . The distortion of $\pi_{d,L}^{\text{sparse}}$ is small but nonzero, since, for example, it never outputs the rare walks in which a vertex is visited more than $\lceil \frac{L}{d} \rceil$ times.

The algorithm $\pi_{d,L}^{\text{dense}}$. Since $\pi_{d,L}^{\text{sparse}}$ must store at least one sampled edge per vertex (and we assume $\frac{L}{d}$ is large), it never achieves sublinear space. We therefore propose a second, fundamentally different, streaming algorithm, $\pi_{d,L}^{\text{dense}}$, which, when d is sufficiently large (relative to L), only uses sublinear memory.

¹ Formally, in a web graph $G = (V, E)$, PageRank values satisfy the system $\text{PageRank}(u) = \sum_{(v,u) \in E} \text{PageRank}(v)/d_{\text{out}}(v)$ for every vertex u , where $d_{\text{out}}(v)$ denotes the out-degree of v [5].

² Throughout, the notation $\tilde{O}(\cdot)$ hides polylogarithmic factors in $\max(n, L)$. Algorithms are assumed to know n in advance, but not the length of the stream (the number of edges).

³ For example, if the first two sampled edges from v_{i-1} have already been used, take the third. A single edge (v, u) may appear and be traversed multiple times if it occurs repeatedly in v 's stored list of outgoing edges.

⁴ We assume $d \geq 1$, since if a vertex (e.g., v_0) has out-degree 0, there may not be any paths of length L starting from v_0 .

The algorithm $\pi_{d,L}^{\text{dense}}$ proceeds as follows: Set $V_0 = \{v_0\}$. Sample L sets of vertices (layers of a graph) V_1, \dots, V_L , each of size $\tilde{O}(\lceil \frac{n}{d} \rceil)$ (if a vertex is sampled in multiple layers, we treat each occurrence as a distinct copy of the original vertex). Intuitively, V_i represents the possible choices for vertex v_i . Then, for each $i \in \{0, \dots, L-1\}$ and $v \in V_i$, the algorithm samples a random outgoing edge from v into V_{i+1} , thereby constructing an L -layer graph. The output is the unique path of length L starting from v_0 in this layered graph.

It is not hard to verify the correctness of $\pi_{d,L}^{\text{dense}}$: $\pi_{d,L}^{\text{dense}}$ fails in the event that some vertex in a layer has no outgoing edge to the next layer. However, the sizes of the layers are selected such that this event happens with low probability. If the algorithm does not fail and the path (v_0, v_1, \dots, v_L) is sampled, we claim that each v_{i+1} is a random (outgoing) neighbor of v_i , and therefore that the output is a random walk. Indeed, since V_{i+1} is a random subset of vertices (of a specific size), choosing a random neighbor of v_i within V_{i+1} is equivalent to selecting a random neighbor of v_i in V . Additionally, the algorithm only stores $\tilde{O}(\lceil \frac{n}{d} \rceil \cdot L)$ edges, and, in particular, on graphs with minimum out-degree $d = \Omega(n)$, its space usage reduces to $\tilde{O}(L)$, which can be sublinear in n .

Note that $\pi_{d,L}^{\text{dense}}$ enjoys the following property: It fails with a small probability, but, conditioned on non-failure, it samples from the true distribution of random walks. $\pi_{d,L}^{\text{sparse}}$ does not satisfy this property.

1.2.2 Streaming Universal Optimality

The notion of streaming universal optimality. $\pi_{d,L}^{\text{dense}}$ is essentially optimal for very dense graphs with $d = \Omega(n)$ as it uses only $\tilde{O}(L)$ memory and the output consists of L edges. But is $\pi_{d,L}^{\text{dense}}$ optimal for all values of d ?

A worst-case optimality result for a graph streaming algorithm would identify a family of hard graph instances (on which the algorithm uses the most memory) and show that no other correct algorithm can use less memory on these instances. However, for easier instances, the algorithm may be highly suboptimal. This notion of optimality can be viewed as *existential* optimality: There exist inputs within the considered class for which the algorithm is optimal. A much more ambitious goal is to seek *universal* optimality, where the algorithm achieves optimal performance on *every* instance.

Universal optimality [10], along with other “beyond worst-case” notions, has been studied for decades in various computational models, see the great book [18]. In this work, we introduce the notion of universal optimality for streaming algorithms. Informally, a streaming algorithm is *universally optimal* if it performs (nearly) as well as possible on *every* graph under *worst-case* choice of streaming order. Specifically, for any graph G , if the algorithm uses S space on some streaming order of G 's edges, then *every* other correct algorithm must also use at least (approximately) S space on some (possibly different) streaming order of G . Further discussion of the notion of universal optimality for streaming algorithms appears in Section 1.2.3.

Universally optimal random walk samplers? A priori, it is not clear that universally optimal algorithms should even exist for the task T of sampling an L -step random walk with distortion, say, 0.1, over the set of graphs \mathcal{G} consisting of directed graphs with minimum out-degree at least d . Indeed, it is not hard to see that neither $\pi_{d,L}^{\text{sparse}}$ nor $\pi_{d,L}^{\text{dense}}$ satisfies such universal optimality for all values of d and L .⁵

⁵ For $d = \Omega(n)$ and $L = \sqrt{n}$, $\pi_{d,L}^{\text{sparse}}$ uses $\Omega(n)$ memory, while $\pi_{d,L}^{\text{dense}}$ uses $\tilde{O}(\sqrt{n})$, so $\pi_{d,L}^{\text{sparse}}$ is not universally optimal. If $L = \sqrt{n}$ and \sqrt{n} vertices have out-degree $\Omega(n)$ while the rest have constant out-degrees, d is constant. In this case, $\pi_{d,L}^{\text{dense}}$ uses $\Omega(n^{1.5})$ memory, whereas $\pi_{d,L}^{\text{sparse}}$ uses only $\tilde{O}(n)$ (as each of $n - \sqrt{n}$ constant out-degree vertices only stores a constant number of edges), showing that $\pi_{d,L}^{\text{dense}}$ is not universally optimal.

Nevertheless, we show that with a careful implementation, $\pi_{d,L}^{\text{sparse}}$ and $\pi_{d,L}^{\text{dense}}$ can be combined into a universally optimal algorithm $\pi_{d,L}$ as follows: If $\frac{L}{d}$ is greater than (roughly) $\log n$, run $\pi_{d,L}^{\text{sparse}}$; otherwise, run $\pi_{d,L}^{\text{dense}}$.

► **Theorem 1.** *For $d, L \in \mathbb{N}$, $\pi_{d,L}$ solves the task of sampling an L -step random walk, with distortion 0.1, from a given vertex in any directed graph with minimum out-degree d . It uses at most $\tilde{O}(\sum_v \min(d_{\text{out}}(v), \frac{L}{d})) \leq \tilde{O}(\frac{n}{d} \cdot L)$ space.*

Furthermore, it is universally optimal over the set of directed graphs with minimum out-degree at least d , for any $d, L \geq \log^{10} n$.

We note that $\pi_{d,L}$ does not work for dynamic streams, since it relies on reservoir sampling [23] to select the outgoing edges to be stored. Nevertheless, by replacing reservoir sampling with ℓ_0 -samplers [15], we can obtain a new algorithm that works for dynamic streams and uses $\tilde{O}(\frac{n}{d} \cdot L)$ space. However, this algorithm is not (insertion-only) universally optimal with a straightforward implementation⁶. In Section 1.2.3, we will show that this turns out to be the best we can do.

As pointed out in [10], an interesting “side-effect” of universal optimality is that it precisely identifies the parameters of the problem that are inherently responsible for its complexity. In our case, the minimum out-degree d of the graph dictates the optimal space complexity (at least when d is known to the algorithm and within certain parameter regimes). In particular, once d is specified, the optimal space requirement no longer depends on other parameters (such as the average out-degree, the minimum total degree, or the length of the longest directed path in the graph).

Do we need to know d ? The algorithm $\pi_{d,L}$ requires advance knowledge of d . A natural question is whether one can design a one-pass streaming algorithm with comparable space usage that does not rely on knowing d beforehand⁷.

Specifically, let d be sufficiently large for Theorem 1 to hold, and let π_L be an algorithm for sampling an L -step random walk in any graph with minimum out-degree at least d (but unknown to π_L). Let G be a graph with minimum out-degree $d' = \Omega(n)$. Running $\pi_{d,L}$ on G requires $\Omega(\frac{n}{d} \cdot L)$ space. By the universal optimality of $\pi_{d,L}$ guaranteed by Theorem 1, it follows that any algorithm that works for all graphs with minimum out-degree at least d , including π_L , must also use $\tilde{\Omega}(\frac{n}{d} \cdot L)$ space on some streaming order of G . In contrast, $\pi_{d',L}$ requires only $\tilde{O}(L)$ space on every streaming order of G . Hence, algorithms that do not know d in advance cannot achieve the same space as $\pi_{d',L}$.

1.2.3 Further Discussion of the Notion of Universal Optimality

Universal optimality for dynamic streaming algorithms. Our definition of universal optimality considers an algorithm τ that works for *insertion-only streams*, and, for every graph, compares its space usage to that of any other algorithm τ' that works for *insertion-only streams*. One can consider the same definition where τ and τ' work for *dynamic streams*. However, we observe that notion of universal optimality for dynamic streaming algorithms coincides with standard worst-case optimality.

⁶ Intuitively, it is possible that a vertex v has $\Omega(L/d)$ outgoing edges during the stream while its actual out-degree is very small at the end. For $\pi_{d,L}^{\text{sparse}}$, this means $\Omega(L/d)$ ℓ_0 samplers are needed. In contrast, if only insertion-only streams are considered, $\pi_{d,L}^{\text{sparse}}$ can simply store all outgoing edges if $d_{\text{out}}(v) < L/d$.

⁷ This can be done in just two passes: the first pass computes the minimum out-degree d , and the second pass runs $\pi_{d,L}$ using the degree d computed in the first pass.

To see this, suppose τ is a worst-case optimal algorithm for some task T on dynamic streams. Fix n . Let S denote the maximum memory used by τ on graphs with n vertices. By (worst-case) optimality of τ , for every other algorithm τ' that solves T , there exists a stream of edge updates $E_{\tau'}$ on n vertices such that τ' requires at least S memory on $E_{\tau'}$. Now take any graph G' on n vertices and consider an algorithm τ' for T . Construct the following stream for G' : First insert/delete edges according to $E_{\tau'}$, then reverse all of them, and finally insert the edges of G' . Since τ' already uses at least S memory on the prefix $E_{\tau'}$, it uses at least S memory on this streaming order of G' . In contrast, τ never exceeds S memory on any stream. Hence, τ is also universally optimal.

Universal optimality vs. instance optimality. Our definition of universal optimality requires the algorithm to perform (nearly) as well as any competitor on every graph G , when comparing their worst-case memory usage across all possible streaming orders of G . However, this leaves open the possibility that on specific streaming orders of G , a competing algorithm could do substantially better. A stronger requirement would be to guarantee optimality on every *stream*, not just every graph. Concretely, for any edge stream I , the memory used by any competing algorithm on I must be at least as large as that used by our algorithm on the same I . This stronger guarantee corresponds to adapting the notion of *instance optimality*, where an algorithm must be optimal on each individual instance, to the context of graph streaming algorithms [9] (see also [1, 22]).

However, it turns out that instance optimality makes little sense. Suppose $d = \sqrt{n}$ and $L = n$. Consider an n -vertex graph G which is the disjoint union of \sqrt{n} cliques of size \sqrt{n} each, along with the input stream where edges are given in lexicographical order. We can construct an algorithm τ which checks if the input stream is exactly the edges of G given in lexicographical order. This can be done in only $\tilde{O}(1)$ additional space because the edges of G can be easily generated in lexicographical order in such space. As soon as the check fails, τ simulates $\pi_{d,L}$ from the beginning of the stream. This is possible because the prefix of the stream can be regenerated. If the check passes at the end, then an L -step random walk can be easily generated in $\tilde{O}(L)$ space as we know all the edges (implicitly). Clearly, τ always performs as well as $\pi_{d,L}$. Furthermore, τ is strictly better than $\pi_{d,L}$ on the specific instance of G and a lexicographical order of the edges of G .

More generally, the above “hard-coding” method can apply to any graph problem, if we can design a more space-efficient algorithm tailored to specific instances that have succinct representations and can be checked in small space. However, note that such things cannot happen with universal optimality as there is no way to check, using small space, if the input graph falls into certain categories.

1.3 Additional Related Work

In an influential work, [19] give an $\tilde{O}(\sqrt{L})$ -pass, $\tilde{O}(n)$ -space algorithm for simulating L -step random walks on directed graphs, and used it to derive space-efficient algorithms for estimating *PageRank* on graph streams. For one-pass streaming, [14] proved an $\Omega(nL)$ space lower bound, thereby proving the optimality of the folklore algorithm. For undirected graphs, [14] also gave an $\mathcal{O}(n \cdot \sqrt{L})$ -space algorithm and a matching lower bound for sufficiently large L . Recently, [7] showed that the space complexity of random walk sampling on directed graphs is $\tilde{\Theta}(n\sqrt{L})$ when *two passes* over the stream are allowed. Additionally, [16, 8] studied the random walk sampling problem in the *random-order* streaming model.

[4] studies sublinear algorithms for MAXCUT and correlation clustering. For dense graphs, constant-space algorithms are known, whereas for sparse graphs, $\Omega(n)$ space lower bounds hold in the streaming model. Their work bridges these two extremes by providing algorithms whose performance is parameterized by the average degree.

2 Proof Overview

We overview the proof of Theorem 1 in this section. As the algorithm mentioned was already described in Section 1, we mostly focus on the proof of universal optimality. However, we do give some extra intuition for the algorithm $\pi_{d,L}^{\text{dense}}$.

2.1 The Sampling Algorithm $\pi_{d,L}^{\text{dense}}$

Just like $\pi_{d,L}^{\text{sparse}}$, our algorithm $\pi_{d,L}^{\text{dense}}$ is based on the folklore algorithm but develops it from a different perspective than $\pi_{d,L}^{\text{sparse}}$. Recall that the standard view of the folklore algorithm interprets its memory as a set $\{L_v\}_{v \in V}$ of n lists, where each L_v stores L sampled outgoing edges from v . The algorithm $\pi_{d,L}^{\text{sparse}}$ is obtained by simply shortening these lists by a factor of d and a proof that such a shortening will not affect the correctness of the algorithm.

For the algorithm $\pi_{d,L}^{\text{dense}}$, we reinterpret the memory of the folklore algorithm as maintaining an L -layered graph rather than a collection of lists. Each layer contains a copy of every vertex v , and the (only) outgoing edge of v in layer i points to some vertex in layer $i + 1$. Specifically, if the i -th edge in L_v is (v, u) , then the copy of v in layer i has a single outgoing edge to the copy of u in layer $i + 1$. This perspective is equivalent to viewing the folklore algorithm's memory as an $n \times L$ matrix of edges, where the row corresponding to vertex $v \in V$ encodes L_v . In this representation, the i -th column specifies the edges connecting layer i to layer $i + 1$ in the layered graph.

In this matrix (or layered-graph) view, $\pi_{d,L}^{\text{sparse}}$ corresponds to restricting the number of columns (equivalently, the number of layers) from L to $\frac{L}{d}$. By contrast, our algorithm $\pi_{d,L}^{\text{dense}}$ reduces the number of rows (equivalently, the size of each layer) from n to $\frac{n}{d}$, using the guarantee on the minimum out-degree to show that such a reduction does not affect the correctness of the algorithm. It is easy to see that this memory restriction of the folklore algorithm yields exactly the algorithm $\pi_{d,L}^{\text{dense}}$ described in Section 1.2.1.

The property testing approach. Prior work has sometimes been motivated by a query model perspective to design algorithms for dense graphs. As such a perspective also yields an algorithm for our problem, we outline it here. However, the algorithm so obtained turns out to be suboptimal.

In the query model, a random walk can be efficiently simulated if adaptivity is allowed. Starting from v_0 , we select a random neighbor v_1 by querying random entries in the adjacency row of v_0 until a 1 is found, then proceed similarly to select a neighbor v_2 of v_1 , and so on. Consider, for instance, the case where minimum degree at least $d = \Omega(n)$. In this case, the expected number of queries needed to find the next vertex is constant. Consequently, the entire process requires roughly $\Theta(L)$ queries in expectation.

To obtain a one-pass streaming algorithm, we aim to transform the above adaptive query algorithm into a non-adaptive one, which are easily simulated by a streaming algorithm. The beautiful technique of [11] shows how to convert an adaptive query algorithm with q queries into a non-adaptive one by sampling an induced subgraph on $\mathcal{O}(q)$ vertices, which requires $\mathcal{O}(q^2)$ queries.⁸ However, applying this idea directly to our adaptive random walk sampling

⁸ While [11] state their result in the context of converting adaptive graph property testers into non-adaptive ones, their technique can be applied to convert the adaptive random walk algorithm above into a non-adaptive version.

algorithm yields a memory complexity of $\tilde{O}(L^2)$ when $d = \Omega(n)$, which is not optimal as it is worse than the bound we get.

2.2 Universal Optimality

We now overview the proof that our algorithm is universally optimal. For this, we first describe a lower bound approach designed for directed bipartite graphs that has been helpful in prior works as well. Consider a directed bipartite graph $G = (V_1 \cup V_2, E)$ and an integer $d > 0$ such that every vertex $v \in V_1$ has out-edges to an independent uniformly random subset of size d_{out} of vertices in V_2 . The vertices in V_2 have only one out-edge and all these edges lead to the same vertex $z \in V_1$ that is chosen uniformly at random.

From a “sampling” problem to a “search” problem. Because all the vertices in V_2 lead to the same vertex z , a random walk starting from any vertex in V_2 would first go to z , then take a uniformly random edge going out of z , then go to z again, and so on for L steps. This means that any such walk contains within it, a sequence of $L/2$ edges going out of z chosen independently and uniformly at random. As the number of distinct out-edges of z is d_{out} , we get that with high probability at least $\Omega(\min(d_{\text{out}}, L))$ of these edges are distinct. We conclude that any algorithm that samples a random walk in this graph can also be used to solve a search problem, where the goal is to identify a subset of edges of size $\Omega(\min(d_{\text{out}}, L))$ coming out of z , with high probability.

Lower bound for the search problem. However, it can be shown that the search problem is hard for streaming algorithms. Indeed, consider what happens when the edges in the graph are presented in a stream with all the edges going out of the vertices in V_1 appearing first in the stream followed by the edges going out of the vertices in V_2 . In this case, when the algorithm is processing the edges going out of the vertices in V_1 , it has no information about the vertex z for which it has to remember the out-edges. Thus, unless the algorithm records enough out-edges from a lot of vertices in V_1 , it is highly likely that when z is chosen uniformly at random from the vertices in V_1 , the algorithm would not have stored enough edges coming out of z .

In other words, in order to output $\Omega(\min(d_{\text{out}}, L))$, the algorithm actually needs to remember this number of edges from all the vertices in V_1 , thereby needing at least $\Omega(|V_1| \cdot \min(d_{\text{out}}, L))$ memory. Observe that this proof is an adaptation of the lower bound proof for the well-known INDEX problem.

Further discussion and extensions of this approach. Note that it is important in the above proof that every vertex in V_1 has the same number d_{out} of edges going out of it. In case the number of out-edges going out of different vertices in V_1 is vastly different, the algorithm can focus on remembering only the edges from vertices that have few out-edges, and hope that when z is sampled it happens to be one of these vertices. This could mean a worse space lower bound.

Also, observe that the instances considered in this lower bound approach are such that all vertices in V_2 have out-degree 1. However, for all $d \leq \min(d_{\text{out}}, |V_1|)$, it can be extended to the case where the out-degree of every vertex is at least d by letting all the vertices in V_2 have out-edges to the same subset $Z \subseteq V_1$ of size d of vertices in V_1 instead of having edges to exactly one vertex $z \in V_1$.

With this change, a random walk would not always return to the same vertex in V_1 and instead returns to a uniformly random vertex in Z . Thus, for any given vertex in Z , a random

walk will only return to it (roughly) $L/(2d)$ many times. In turn, we have that the algorithm only needs to store around $L/(2d)$ many edges coming out of any vertex. As this number is now smaller, this would lead to a reduced space lower bound of $\Omega(|V_1| \cdot \min(d_{\text{out}}, \frac{L}{d}))$. This is exactly where the expression in Theorem 1 comes from, but its not the end of the proof, as the input graph is not guaranteed to be bipartite or promised to have the same number of edges from all vertices in V_1 .

Extending to general graphs G . We now turn to show how the above approach can be extended to work for all graphs and give a lower bound of $\tilde{\Omega}(\sum_v \min(d_{\text{out}}(v), \frac{L}{d}))$ by embedding a hard instance of the form above in any input graph G . For this, we “clean” any input graph G in a sequence of steps, ultimately leading to a hard instance like those above.

Our first step is to identify a “dominant” degree in the graph. For this, we partition the vertices into $\log n$ levels based on which two consecutive powers of 2 its out-degree lies in. This also partitions the sum we want in our lower bound and we focus on the part that is the largest. This gives a subset of vertices $V' \subseteq V$ of vertices that have similar degrees and suffices for the lower bound up to logarithmic factors.

Having decided on a set V' to focus on, we now identify a cut in the graph G such that a constant fraction of the vertices in V' lie on the left side of the cut, and a constant fraction of the neighbors of all these vertices lie on the right side of the cut. To show that such a cut exists, note that a random cut in the graph will have the property such that in expectation, a constant fraction of the neighbors of any vertex lie on a different side of the cut as that vertex. Using probabilistic arguments, we show that a cut exists for which this is true for a constant fraction of the vertices in V' .

We are almost there, except that the vertices in V' don't have *exactly* the same degree which was assumed in our proof approach. This can be easily fixed by removing some of the edges from all vertices in V' till they have the same degree. As they had degrees within a factor of 2 to being with, this loses only a constant fraction of edges and only costs us another constant factor in our lower bound. Another difference is that we started with a graph where every vertex is promised to have out-degree at least d . However, as we removed the edges not in the cut, we may have shrunk the out-degrees of vertices. As this may make some out-degrees smaller than d , it may take the resulting graph outside the promise, effectively making it useless. To fix this, we identify another set of $\mathcal{O}(d)$ randomly chosen vertices in the graph and connect them with each other in a clique and also connect all the other vertices to them. This increases the degrees of all vertices back to at least d , satisfying the promise. Moreover, as⁹ $d = o(n)$, this only loses us another constant in the lower bound.

A hard distribution. Finally, note that the above construction only gave us a single graph while our lower bound approach required a distribution of graphs. This is a major difference, as there are algorithms that will solve the search problem on this particular graph (one can simply “hardwire” the answer for this graph in the algorithm) making a lower bound impossible. We get around this by converting our hard graph into a distribution over graphs as follows: For every vertex in V' , we only include a randomly chosen subset of half its edges in the graph and the other edges are discarded. This again only costs us a constant factor but creates for us a distribution just like we need for the lower bound approach.

⁹ When d and n are within a constant factor of each other, our bound becomes $\tilde{\mathcal{O}}(L)$ which is tight as it is the size of the output.

We finish this section by arguing why a lower bound obtained from such a hard distribution after making all these changes would also apply to the original graph G , as required for universal optimality. For this, we observe that our lower bound approach shows a lower bound on the memory needed at a *specific* point in the stream, when it has processed all the edges coming out from V_1 but before processing any edges that come out from V_2 . In our constructed distribution, the set V_1 corresponds to the set V' and thus, the inputs streams we generate will start with the edges coming out of the vertices in V' and the memory lower bound will hold for the memory state of the algorithm immediately after processing the edges. Crucially, (unlike the edges going into V' or the edges added to raise the degree) all these edges were also part of the original graph G . Thus, for all the algorithm knows, the graph might as well have been G with these edges presented first in the stream, and the lower bound we show for our constructed distribution will also apply to the original graph G , finishing the proof.

3 Preliminaries

Notation. For integer $n > 0$, we write $[n]$ as a shorthand for $\{1, \dots, n\}$. \mathbb{N} denotes the set of positive integers. We often omit ceiling symbols, and by assigning a real number x to an integer y , we mean that y gets $\lceil x \rceil$.

Graphs. For a directed graph $G = (V, E)$ and any vertex $v \in V$, we define $N_{\text{out}}^G(v)$ to be the set of all out-edges of v and define $d_{\text{out}}^G(v) = |N_{\text{out}}^G(v)|$ to be the out-degree of v . We omit the superscript G when it is clear from context. We use $n = |V|$ throughout.

Random walks. An L -step random walk in G with starting vertex $v_0 \in V$ is a (directed) path (v_0, \dots, v_L) of length L in G such that v_i is an independent and uniformly random out-neighbor of v_{i-1} for all $i \in [L]$. We say an algorithm samples an L -step random walk starting from v_0 with *distortion* $\delta \geq 0$, if the sampled distribution of (v_0, \dots, v_L) is δ -close in statistical distance to the true distribution of an L -step random walk starting from v_0 .

Graph streaming. An (insertion-only) graph streaming algorithm takes as input a sequence of edges of an underlying graph, and is required to process the edges one at a time using limited space. At the end of the stream, the algorithm needs to return an output for certain graph problems. In this paper, we consider the task of sampling random walks.

Universal optimality. A graph streaming algorithm τ is *universally optimal* for a task T over a set of graphs \mathcal{G} if:

1. τ solves T .
2. For *any* streaming algorithm τ' that also solves T , the following holds: Let G be *any* graph in \mathcal{G} and let S be the maximum memory size used by τ' over all streaming orders of G . Then, the maximum memory size used by τ over any streaming order of G is at most $S \cdot \text{poly } n$.

Chernoff bound. We use the following version of the Chernoff bound.

55:10 Universally Optimal Streaming Algorithm for Random Walks in Dense Graphs

► **Lemma 2** (Multiplicative Chernoff bound). *Suppose X_1, \dots, X_n are independent random variables taking values in $[0, 1]$. Let X denote their sum. Then,*

$$\Pr(X \geq (1 + \delta) \cdot \mathbb{E}[X]) \leq e^{-\frac{\delta^2 \cdot \mathbb{E}[X]}{2 + \delta}}, \quad \forall 0 \leq \delta,$$

$$\Pr(X \leq (1 - \delta) \cdot \mathbb{E}[X]) \leq e^{-\frac{\delta^2 \cdot \mathbb{E}[X]}{2}}, \quad \forall 0 \leq \delta \leq 1.$$

Information theory. We use sans-serif letters to denote random variables. We reserve \mathcal{E} to denote an arbitrary event. All random variables will be assumed to be discrete and we shall adopt the convention $0 \log \frac{1}{0} = 0$. When it is clear from context, we may abbreviate an event $X = x$ as just x . All logarithms are taken with base 2.

► **Definition 3** (Entropy). *The (binary) entropy of X is defined as:*

$$\mathbb{H}(X) = \sum_{x \in \text{supp}(X)} \Pr(x) \cdot \log \frac{1}{\Pr(x)}.$$

The entropy of X conditioned on \mathcal{E} is defined as:

$$\mathbb{H}(X | \mathcal{E}) = \sum_{x \in \text{supp}(X)} \Pr(x | \mathcal{E}) \cdot \log \frac{1}{\Pr(x | \mathcal{E})}.$$

► **Definition 4** (Conditional entropy). *We define the conditional entropy of X given Y and \mathcal{E} as:*

$$\mathbb{H}(X | Y, \mathcal{E}) = \sum_{y \in \text{supp}(Y)} \Pr(y | \mathcal{E}) \cdot \mathbb{H}(X | y, \mathcal{E}).$$

Henceforth, we shall omit writing the $\text{supp}(\cdot)$ when it is clear from context.

► **Lemma 5** (Chain rule for entropy). *It holds for all X, Y, Z and \mathcal{E} that:*

$$\mathbb{H}(X, Y | Z, \mathcal{E}) = \mathbb{H}(X | Z, \mathcal{E}) + \mathbb{H}(Y | X, Z, \mathcal{E}).$$

► **Lemma 6** (Conditioning reduces entropy). *It holds for all X, Y, Z and \mathcal{E} that:*

$$\mathbb{H}(X | Y, Z, \mathcal{E}) \leq \mathbb{H}(X | Z, \mathcal{E}).$$

Equality holds if and only if X and Y are independent conditioned on Z, \mathcal{E} .

As a corollary of Lemmas 5 and 6, we get that entropy is subadditive.

► **Corollary 7** (Subadditivity of entropy). *It holds for all X, Y, Z and \mathcal{E} that:*

$$\mathbb{H}(X, Y | Z, \mathcal{E}) \leq \mathbb{H}(X | Z, \mathcal{E}) + \mathbb{H}(Y | Z, \mathcal{E}).$$

► **Lemma 8.** *It holds for all X and \mathcal{E} that:*

$$0 \leq \mathbb{H}(X | \mathcal{E}) \leq \log(|\text{supp}(X)|).$$

The second inequality is tight if and only if X conditioned on \mathcal{E} is the uniform distribution over $\text{supp}(X)$.

As a corollary of Lemmas 5 and 8, we also have the following bound.

■ **Algorithm 1** The protocol $\pi_{d,L}(G, v_0)$.

-
- 1: If $\frac{L}{d} > 10^3 \log n$, run $\pi_{d,L}^{\text{sparse}}(G, v_0)$.
 - 2: Otherwise, run $\pi_{d,L}^{\text{dense}}(G, v_0)$.
-

► **Corollary 9.** *It holds for all X, Y, Z and \mathcal{E} that:*

$$\mathbb{H}(X | Y, Z, \mathcal{E}) \geq \mathbb{H}(X | Z, \mathcal{E}) - \mathbb{H}(Y | Z, \mathcal{E}).$$

► **Lemma 10** (Shearer's inequality). *Let $n > 0$ and X_1, \dots, X_n be random variables. Let $X = (X_1, \dots, X_n)$ and $Z \subseteq [n]$ be a set-valued random variable that is independent of X . For any event \mathcal{E} determined by X , it holds that:*

$$\mathbb{H}(X | \mathcal{E}) \cdot \min_{i \in [n]} \Pr(i \in Z) \leq \mathbb{H}(X_Z | Z, \mathcal{E}).$$

Proof. We have:

$$\begin{aligned} \mathbb{H}(X_Z | Z, \mathcal{E}) &= \sum_Z \Pr(Z | \mathcal{E}) \cdot \mathbb{H}(X_Z | Z, \mathcal{E}) && \text{(Definition 4)} \\ &= \sum_Z \Pr(Z | \mathcal{E}) \cdot \mathbb{H}(X_Z | Z, \mathcal{E}) \\ &= \sum_Z \Pr(Z) \cdot \mathbb{H}(X_Z | \mathcal{E}) && \text{(as } Z \text{ is independent of } X, \mathcal{E}) \\ &= \sum_Z \Pr(Z) \cdot \sum_{i \in Z} \mathbb{H}(X_i | X_{Z \cap [i-1]}, \mathcal{E}) && \text{(Lemma 5)} \\ &\geq \sum_Z \Pr(Z) \cdot \sum_{i \in Z} \mathbb{H}(X_i | X_{[i-1]}, \mathcal{E}) && \text{(Lemma 6)} \\ &= \sum_{i=1}^n \sum_Z \Pr(Z) \cdot \mathbf{1}(i \in Z) \cdot \mathbb{H}(X_i | X_{[i-1]}, \mathcal{E}) \\ &\geq \min_{i \in [n]} \Pr(i \in Z) \cdot \sum_{i=1}^n \mathbb{H}(X_i | X_{[i-1]}, \mathcal{E}) \\ &= \mathbb{H}(X | \mathcal{E}) \cdot \min_{i \in [n]} \Pr(i \in Z). && \text{(Lemma 5)} \end{aligned}$$

◀

► **Corollary 11.** *Let $n > 0$ and X_1, \dots, X_n be random variables. Let $X = (X_1, \dots, X_n)$ and $Z \subseteq [n]$ be a set-valued random variable that is independent of X . For any random variable Y determined by X , it holds that:*

$$\mathbb{H}(X | Y) \cdot \min_{i \in [n]} \Pr(i \in Z) \leq \mathbb{H}(X_Z | Z, Y).$$

4 Streaming Algorithm

The goal of this section is to prove the following Theorem 12. Algorithms $\pi_{d,L}$, $\pi_{d,L}^{\text{sparse}}$, and $\pi_{d,L}^{\text{dense}}$, are given in Algorithms 1–3, respectively. By “output FAIL,” we mean that the algorithm outputs (v_0, \dots, v_0) . We will elaborate on the details of their streaming implementation later on.

■ **Algorithm 2** The protocol $\pi_{d,L}^{\text{sparse}}(G, v_0)$.

-
- 1: For each vertex $v \in V$, if $d_{\text{out}}(v) \leq \frac{2L}{d}$, store all its outgoing edges, and otherwise, sample $\frac{2L}{d}$ outgoing edges with replacement, uniformly at random, denoted $u_{v,1}, \dots, u_{v, \frac{2L}{d}}$.
 - 2: Output a path (v_0, v_1, \dots, v_L) such that for all $i \in [L]$, if $d_{\text{out}}(v_{i-1}) \leq \frac{2L}{d}$, v_i is an independent and uniformly random out-neighbor of v_{i-1} , and otherwise, $v_i = u_{v_{i-1}, t}$ where $t = |\{i' \in [i] \mid v_{i'-1} = v_{i-1}\}|$. If $t > \frac{2L}{d}$ for some $i \in [L]$, output FAIL.
-

■ **Algorithm 3** The protocol $\pi_{d,L}^{\text{dense}}(G, v_0)$.

-
- 1: Set $V_0 = \{v_0\}$.
 - 2: Independently sample subsets $V_1, \dots, V_L \subseteq V$, each of size $s = \min(\frac{n}{d} \cdot \ln(10Ln), n)$, uniformly at random.
 - 3: For each $i \in [L]$ and vertex $v \in V_{i-1}$, independently sample an outgoing edge (v, u) , uniformly at random, conditioned on $u \in V_i$.
 - 4: Output the (unique) path (v_0, v_1, \dots, v_L) such that for all $i \in [L]$, (v_{i-1}, v_i) is sampled. If there is no edge from v_{i-1} to V_i for some $i \in [L]$, output FAIL.
-

► **Theorem 12.** *Let $d, L \in \mathbb{N}$. $\pi_{d,L}$ solves the task of sampling an L -step random walk with distortion 0.1. It uses $\tilde{O}(\sum_v \min(d_{\text{out}}(v), L/d))$ space.*

As Algorithm 1 is simply a combination of Algorithms 2 and 3, Theorem 12 is a direct corollary of the following Lemmas 13 and 14. Indeed, when $L/d \leq 10^3 \log n$, we know that L/d is larger than $\min(d_{\text{out}}(v), L/d)$ by a multiplicative factor of at most $10^3 \log n$ as $d_{\text{out}}(v) \geq d \geq 1$. Note that $\sum_v \min(d_{\text{out}}(v), L/d)$ can be sublinear in n when $L/d \ll 1$.

We also emphasize that the proof of Lemma 14 does not require $L/d \leq 10^3 \log n$. In other words, Algorithm 3 alone is sufficient to achieve $\tilde{O}(Ln/d)$ space in the worst case. However, as noted in Section 1, Algorithm 2 and Lemma 13 are necessary for achieving universal optimality.

► **Lemma 13.** *Let $d, L \in \mathbb{N}$ such that $L/d > 10^3 \log n$. $\pi_{d,L}^{\text{sparse}}$ solves the task of sampling an L -step random walk with distortion 0.1. It uses $\tilde{O}(\sum_v \min(d_{\text{out}}(v), L/d))$ space.*

► **Lemma 14.** *Let $d, L \in \mathbb{N}$. $\pi_{d,L}^{\text{dense}}$ outputs FAIL with probability at most 0.1 and otherwise samples an (unbiased) L -step random walk. It uses $\tilde{O}(Ln/d)$ space.*

The following two sections constitute proofs of Lemmas 13 and 14, respectively.

4.1 Proof of Lemma 13

The correctness of Algorithm 2 follows from the following claim.

▷ **Claim 15.** *Let $d, L \in \mathbb{N}$ such that $L/d > 10^3 \log n$. For any directed graph $G = (V, E)$ with minimum out-degree at least d and vertex $v_0 \in V$, with probability at least 0.9, an L -step random walk from v_0 in G visits any vertex at most $2L/d$ times.*

Proof. Fix $v \in V$. For $i \in [L]$, let V_i be the random variable for the i -th vertex on an L -step random walk from v_0 in G . For any sub-path (v_0, \dots, v_{i-1}) , $\Pr(V_i = v \mid V_0 = v_0, \dots, V_{i-1} = v_{i-1}) \leq \frac{1}{d}$ because v_{i-1} has at least d outgoing edges, with v_i possibly being one of them.

Note that the events $V_i = v$ may not be independent. Nevertheless, as each step of a random walk is sampled independently, the probability of v being visited more than $2L/d$

times is always upper bounded by the probability that L independent biased coins, each with $1/d$ probability of getting a head, generate more than $2L/d$ heads in total. By Lemma 2 and the assumption $L/d > 10^3 \log n$, this happens with probability at most $1/n^2$. The claim then follows by union bounding over all vertices. \triangleleft

As to the space usage, Algorithm 2 can be implemented using reservoir sampling [23]. Concretely, for each vertex $v \in V$, it takes $\mathcal{O}(\log n)$ bits to keep track of $d_{\text{out}}(v)$. Meanwhile, it stores all its outgoing edges while $d_{\text{out}}(v) \leq 2L/d$. Once $d_{\text{out}}(v) > 2L/d$ during the stream, it can use $2L/d$ independent reservoir samplers, each taking $\mathcal{O}(\log n)$ bits. So each vertex v takes $\tilde{\mathcal{O}}(\min(d_{\text{out}}(v), L/d))$ space. Since the reservoir samplers are unbiased, this concludes the proof of Lemma 13 by applying Claim 15.

4.2 Proof of Lemma 14

To implement Algorithm 3, for each $i \in [L]$ and vertex $v \in V_{i-1}$, one reservoir sampler is enough since all V_i 's can be sampled in advance. The space bound follows as L is at most $\text{poly}(n)$.

We now show that conditioned on no failure, Algorithm 3 samples an (unbiased) L -step random walk. In other words, for every $i \in [L]$, v_i has to be a uniformly random out-neighbor of v_{i-1} in G . To this end, consider any two out-neighbors u, u' of v_{i-1} in G , and subset $U \subseteq V$ of size s . Suppose $V_i = U$. If $u, u' \in U$, each of them has the same probability of being sampled as v_i . If $u, u' \notin U$, none of them would be sampled as v_i . If exactly one of u, u' is in U , assume $u \in U$ without loss of generality. Since V_i is a uniformly random subset of size s , $U' = (U \cup \{u'\}) \setminus \{u\}$ has the same probability as U of being sampled as V_i . Furthermore, the number of out-neighbors of v_{i-1} in V_i is the same in the two cases. As a result, the probability of u' being sampled as v_i in the case $V_i = U'$ is the same as the probability of u being sampled as v_i in the case $V_i = U$. Altogether, this means that the overall probability of being sampled as v_i is the same for u, u' .

It then remains to show that Algorithm 3 outputs FAIL with probability at most 0.1. Note that the only possibility of failure is when there exists $i \in [L]$ and vertex $v \in V_{i-1}$ such that v has no out-neighbor in V_i . This happens only if $s = n/d \cdot \ln(10Ln)$. For any $i \in [L]$ and $v \in V_{i-1}$, the probability of v not having out-neighbors in V_i is at most

$$\frac{\binom{n-d}{s}}{\binom{n}{s}} = \prod_{j=1}^s \frac{n-d-j+1}{n-j+1} \leq \left(\frac{n-d}{n}\right)^s \leq \exp\left(-\frac{ds}{n}\right) = \frac{1}{10Ln}.$$

This concludes the proof of Lemma 14 by union bounding over all $1 + s(L-1) \leq Ls \leq Ln$ options for i, v .

5 Universal Optimality

The goal of this section is to show that the algorithm described in Algorithm 1 is universally optimal upto logarithmic factors.

We define a triple of positive integers $L, d, n \in \mathbb{N}$ to be *significant* if $n > 0$ is large enough for asymptotic inequalities to hold and $L \geq (\log n)^{10}$ and $(\log n)^{10} \leq d \leq \frac{n}{(\log n)^{10}}$. For $d, n > 0$, we define the set $\mathcal{G}_{n,d}$ to be the set of all directed simple graphs where there are n

¹⁰If $d \geq n/(\log n)^{10}$, $\pi_{d,L}$ (or simply $\pi_{d,L}^{\text{dense}}$) uses at most $\tilde{\mathcal{O}}(L)$ space. This must be universally optimal as outputting a length- L path requires $\Omega(L)$ space.

vertices and every vertex has out-degree at least d . We say that a streaming algorithm Alg is a *random-walk algorithm* if given significant integers L, d, n , a graph $G = (V, E) \in \mathcal{G}_{n,d}$ and a start vertex $v \in V$, the algorithm Alg samples an L -step random walk from the start vertex with distortion at most 0.1. The theorem we show in this section is the following Theorem 16. Theorem 1 then follows directly from Theorem 16 and Theorem 12.

► **Theorem 16.** *Let $L, d, n \in \mathbb{N}$ be significant, Alg be a random-walk algorithm, and $G = (V, E) \in \mathcal{G}_{n,d}$ with a designated start vertex. There exists a streaming order for the edges of G on which Alg needs memory at least $\frac{1}{(\log n)^4} \cdot \sum_{v \in V} \min(d_{\text{out}}(v), \frac{L}{d})$.*

The proof of Theorem 16 spans this entire section. Fix L, D, n, Alg , and G be as in the theorem statement. As $G \in \mathcal{G}_{n,d}$, we have that for all $v \in V$, it holds that $d \leq d_{\text{out}}(v) < n$. This means that there exists $\log d < k \leq \log n$ such that:

$$\frac{1}{\log n} \cdot \sum_{v \in V} \min\left(d_{\text{out}}(v), \frac{L}{d}\right) \leq \sum_{\substack{v \in V \\ 2^{k-1} \leq d_{\text{out}}(v) < 2^k}} \min\left(d_{\text{out}}(v), \frac{L}{d}\right). \quad (1)$$

We fix this k for the rest of this proof and define V' to be the set of vertices in the sum on the right. Namely,

$$V' = \{v \in V \mid 2^{k-1} \leq d_{\text{out}}(v) < 2^k\}. \quad (2)$$

Our first lemma shows that the set V' is large.

► **Lemma 17.** *It holds that $|V'| \geq \frac{d}{\log n}$.*

Proof. We first prove this under the assumption that $L \leq d^2$. In this case, we have $\frac{L}{d} \leq d \leq d_{\text{out}}(v)$ for all $v \in V$. Plugging this in Equation (1), we get $\frac{1}{\log n} \cdot \frac{nL}{d} \leq |V'| \cdot \frac{L}{d}$ which gives the lemma as $d \leq n$. Now, assume that $L > d^2$. This means that $\frac{L}{d} > d$ implying that the minimum in the left hand side in Equation (1) can be lower bounded by d while the minimum on the right hand side can be upper bounded by n . Plugging in, we get $\frac{1}{\log n} \cdot nd \leq n \cdot |V'|$. Rearranging gives the result. ◀

We use $\Pr_{U \subseteq V}(\cdot)$ and $\mathbb{E}_{U \subseteq V}[\cdot]$ to denote the probability and expectation when U is a subset of vertices chosen uniformly at random. For any subset $U \subseteq V$, we define $\text{cut}(U) = E \cap ((U \times \bar{U}) \cup (\bar{U} \times U))$ to be the set of edges in the cut formed by U . Also, for all $U \subseteq V$, define:

$$\text{Big}(U) = \{v \in V \mid |N_{\text{out}}(v) \cap \text{cut}(U)| \geq 0.1 \cdot d_{\text{out}}(v)\}. \quad (3)$$

We show that there exists $U \subseteq V$ such that many of the vertices in V' lie in $\text{Big}(U) \cap U$.

► **Lemma 18.** *There exists $U \subseteq V$ satisfying:*

$$|V' \cap \text{Big}(U) \cap U| \geq \frac{1}{5} \cdot |V'|.$$

Proof. For any $v \in V'$ such that $d_{\text{out}}(v) \leq 10$, we have that $v \notin \text{Big}(U)$ (for some $U \subseteq V$) if and only if all its out-neighbors are on the same side of the cut as v . As every vertex has at least one out-neighbor, we get that this event happens with probability at most $\frac{1}{2}$ for such v for a uniformly random $U \subseteq V$. For $v \in V'$ such that $d_{\text{out}}(v) \geq 10$, we have by Lemma 2 that:

$$\Pr_{U \subseteq V}(v \notin \text{Big}(U)) \leq \Pr_{U \subseteq V}(|N_{\text{out}}(v) \cap \text{cut}(U)| \leq 0.1 \cdot d_{\text{out}}(v)) \leq e^{-0.1 \cdot d_{\text{out}}(v)} < \frac{1}{2}.$$

Combining, we get that $\Pr_{U \subseteq V}(v \notin \text{Big}(U)) \leq \frac{1}{2}$ for all $v \in V'$. It follows that:

$$\mathbb{E}_{U \subseteq V} [V' \setminus \text{Big}(U)] \leq \frac{1}{2} \cdot |V'|.$$

In particular, there exists a cut U for which $V' \setminus \text{Big}(U) \leq \frac{1}{2} \cdot |V'|$. From this, we get that at least one of the quantities $V' \cap \text{Big}(U) \cap U$ and $V' \cap \text{Big}(U) \cap \bar{U}$ is at least $\frac{1}{5} \cdot |V'|$. To finish the proof, observe from Equation (3) that $\text{Big}(U) = \text{Big}(\bar{U})$ and thus, at least one of U and \bar{U} satisfies the lemma. \blacktriangleleft

Next, we fix U to be the set promised by Lemma 18 and define $U' = V' \cap \text{Big}(U) \cap U$ and $d' = 0.1 \cdot \max(d, 2^{k-1})$. It follows that $|U'| \geq \frac{1}{5} \cdot |V'| \geq \frac{d}{5 \log n}$ by Lemma 17. Moreover, for all $v \in U'$, we have that

$$\begin{aligned} |N_{\text{out}}(v) \cap (\{v\} \times \bar{U}')| &\geq |N_{\text{out}}(v) \cap (\{v\} \times \bar{U})| \\ &= |N_{\text{out}}(v) \cap \text{cut}(U)| \\ &\geq 0.1 \cdot d_{\text{out}}(v) \\ &\geq 0.1 \cdot \max(d, 2^{k-1}) \\ &= d'. \end{aligned} \tag{4}$$

Next, we isolate a suitably chosen set of vertices $W_1 \subseteq V$.

► **Lemma 19.** *There exists a subset $W_1 \subseteq V$ such that all the following hold:*

$$\begin{aligned} 2d &\leq |W_1| \leq 4d, \\ |U' \setminus W_1| &\geq \frac{1}{2} \cdot |U'|, \\ \forall v \in U' : \quad |N_{\text{out}}(v) \cap (\{v\} \times (\bar{U}' \cap \bar{W}_1))| &\geq \frac{1}{2} \cdot |N_{\text{out}}(v) \cap (\{v\} \times \bar{U}')|. \end{aligned}$$

Proof. Proof by the probabilistic method. Consider the set-valued random variable W_1 obtained by including all $v \in V$ with probability $3d/n$. By Lemma 2 and using the fact that $d > (\log n)^{10}$, we get that the first inequality holds except with probability $\frac{1}{n^2}$. For the other inequalities, consider any set $S \subseteq V$ such that $|S| \geq (\log n)^5$. Note that because of our bound $|U'| \geq \frac{d}{5 \log n}$ and the bound in Equation (4), we have that U' and¹¹ $N_{\text{out}}(v) \cap (\{v\} \times \bar{U}')$ for all $v \in U'$ are such sets. As $d \leq \frac{n}{(\log n)^{10}}$ implies that $3d/n \leq 0.01$, we have from Lemma 2 that $\Pr(|S \cap W_1| \geq \frac{1}{2} \cdot |S|) \leq \frac{1}{n^2}$ for any such set S . A union bound over all sets S considered in the lemma statement finishes the proof. \blacktriangleleft

Henceforth, we let W_1 be as promised by Lemma 19. We also define $W_2 = U' \setminus W_1$ and conclude from Lemma 19 that $|W_2| \geq \frac{1}{2} \cdot |U'| \geq \frac{d}{10 \log n}$. Finally, we define $W_3 = V \setminus (W_1 \cup W_2)$. For all $v \in W_2$, we have:

$$\begin{aligned} |N_{\text{out}}(v) \cap (\{v\} \times W_3)| &= |N_{\text{out}}(v) \cap (\{v\} \times \overline{W_1 \cup W_2})| \\ &= |N_{\text{out}}(v) \cap (\{v\} \times \overline{U' \cup W_1})| \\ &\geq \frac{1}{2} \cdot |N_{\text{out}}(v) \cap (\{v\} \times \bar{U}')| \\ &\geq d'/2. \end{aligned} \tag{5}$$

¹¹Technically, as it needs to be a subset of V , the projection of $N_{\text{out}}(v) \cap (\{v\} \times \bar{U}')$ on its second coordinate is such a set, but the same argument applies.

55:16 Universally Optimal Streaming Algorithm for Random Walks in Dense Graphs

For all $v \in W_2$, let E'_v be an (arbitrary) subset of its out-edges going to vertices in W_3 of size $d'/2$. Observe from Equation (5) that at least one such subset exists and therefore, E'_v is well defined.

Next, for all $v \in W_2$, define the set-valued random variable $X_v \subseteq E'_v$ to be an independent and uniformly random subseteq of E'_v of size $d'/4$. For notational convenience, for any subset $Z \subseteq W_2$, we will abbreviate $(X_v)_{v \in Z}$ to X_Z and omit the subscript when $Z = W_2$. Also, define the set-valued random variable $Z \subseteq W_2$ to be a uniformly random subset of W_2 of size $\frac{d}{20 \log n}$ chosen independently of X . As $|W_2| \geq \frac{d}{10 \log n}$, we have that the random variable Z is well defined. The following lemma is straightforward.

► **Lemma 20.** *Let $v \in W_2$ and $S_v \subseteq E'_v$. It holds that:*

$$|\{X_v \in \text{supp}(X_v) \mid S_v \subseteq X_v\}| \leq \binom{d'/2}{d'/4} \cdot 2^{-|S_v|}.$$

Proof. The size of the set we have to bound equals $\binom{d'/2 - |S_v|}{d'/4 - |S_v|} \leq \binom{d'/2}{d'/4} \cdot 2^{-|S_v|}$. ◀

5.1 A “Search” Lower Bound

Before showing the desired lower bound on the memory used by the algorithm **Alg**, we first define a related distributional search problem and show a lower bound for algorithms solving that. For this, let $X = (X_v)_{v \in W_2}$ and Z be as defined above. We use these random variables to define a graph G' with the vertex set $V = W_1 \cup W_2 \cup W_3$ and the following edges:

- (1) For all $v \in V$ and all $v' \neq v \in W_1$, there is an edge from v to v' . These edges are called the fixed edges.
- (2) For all $v \in W_2$, all the edges in X_v are present in G' . These edges are called the X -edges.
- (3) For all $v \in W_1$ and all $v' \in Z \subseteq W_2$, there is an edge from v to v' . These edges are called the Z -edges.

We let G' be the random variable corresponding to G' . Note that the random variables G' and (X, Z) determine each other, and, more specifically, the X -edges and X determine each other and the Z -edges and Z determine each other. Also, note from the fact that $|W_1| \geq 2d$ (see Lemma 19) that the out-degree of any vertex in any graph G' in the support of G' is at least d . It follows that the support of G' is a subset of $\mathcal{G}_{n,d}$.

Having defined the random variable G' , the search problem we consider receives the edges of a graph $G \sim G'$, which we can equivalently view as a pair (X, Z) , in the stream, with the X -edges arriving first (in lexicographic order), then the other edges (in lexicographic order). The goal of the algorithm is to make one-pass over the stream, and with probability at least 0.8, output at least $\frac{d}{(\log n)^3} \cdot \min(d', \frac{L}{d})$ distinct edges in $\bigcup_{v \in Z} X_v$.

Consider any streaming algorithm \mathcal{A} that solves the above problem. As we study \mathcal{A} on a specific distribution of inputs (specifically, the distribution G'), we can assume without loss of generality that \mathcal{A} is deterministic. Let M the random variable denoting the memory state of the algorithm \mathcal{A} after it has processed the X -edges but before processing the other edges and let O be the random variable denoting the output of \mathcal{A} . As \mathcal{A} is deterministic, observe that M is a function of X and that O is a function of (M, Z) . We abbreviate events of the form $M = M$ to just M , $Z = Z$ to just Z , and so on.

► **Lemma 21.** *It holds that:*

$$\mathbb{H}(X_Z \mid Z, M) \leq \frac{d}{20 \log n} \cdot \log \binom{d'/2}{d'/4} - \frac{1}{2} \cdot \frac{d}{(\log n)^3} \cdot \min\left(d', \frac{L}{d}\right).$$

Proof. Let E be an indicator random variable that is 1 if and only if the algorithm \mathcal{A} is correct on input G' (which determines and is determined by (X, Z)). We have:

$$\begin{aligned} \mathbb{H}(X_Z | Z, M) &\leq \mathbb{H}(X_Z | Z, M, E) + \mathbb{H}(E | Z, M) && \text{(Corollary 9)} \\ &\leq 1 + \mathbb{H}(X_Z | Z, M, E) && \text{(Lemma 8)} \\ &\leq 1 + \sum_{Z, M, E} \Pr(Z, M, E) \cdot \mathbb{H}(X_Z | Z, M, E) && \text{(Definition 4)} \\ &\leq 1 + \sum_{Z, M, E} \Pr(Z, M, E) \cdot \mathbb{H}(X_Z | Z, M, E). \end{aligned}$$

We now use the upper bound for the entropy in Lemma 8. For the terms corresponding to $E = 0$, we use Lemma 20 with $S_v = \emptyset$ to get $|\text{supp}(X_v)| \leq \binom{d'/2}{d'/4}$ for all $v \in Z$. For the terms corresponding to $E = 1$, we use the fact that Z and M together determine the output of the algorithm, which when $E = 1$, determines $\frac{d}{(\log n)^3} \cdot \min(d', \frac{L}{d})$ of the edges in $\bigcup_{v \in Z} X_v$. By Lemma 20, this means that conditioned on $Z, M, E = 1$, the random variable X_Z is supported on a set of size at most $\left(\binom{d'/2}{d'/4}\right)^{\frac{d}{20 \log n}} \cdot 2^{-\frac{d}{(\log n)^3} \cdot \min(d', \frac{L}{d})}$. We get:

$$\begin{aligned} \mathbb{H}(X_Z | Z, M) &\leq 1 + \frac{d}{20 \log n} \cdot \log \binom{d'/2}{d'/4} - \Pr(E = 1) \cdot \frac{d}{(\log n)^3} \cdot \min\left(d', \frac{L}{d}\right) \\ &\leq 1 + \frac{d}{20 \log n} \cdot \log \binom{d'/2}{d'/4} - 0.8 \cdot \frac{d}{(\log n)^3} \cdot \min\left(d', \frac{L}{d}\right) && \text{(As } \Pr(E = 1) \geq 0.8\text{)} \\ &\leq \frac{d}{20 \log n} \cdot \log \binom{d'/2}{d'/4} - \frac{1}{2} \cdot \frac{d}{(\log n)^3} \cdot \min\left(d', \frac{L}{d}\right). && \text{(As } L, d \geq (\log n)^{10} \text{ and } d' \geq 0.1d\text{)} \end{aligned}$$

◀

Combining Lemma 21 with Corollary 11, we get:

$$\mathbb{H}(X | M) \cdot \frac{d}{20 \log n} \cdot \frac{1}{|W_2|} \leq \frac{d}{20 \log n} \cdot \log \binom{d'/2}{d'/4} - \frac{1}{2} \cdot \frac{d}{(\log n)^3} \cdot \min\left(d', \frac{L}{d}\right).$$

Using the fact that $\mathbb{H}(X) = |W_2| \cdot \log \binom{d'/2}{d'/4}$ and simplifying, we get:

$$\mathbb{H}(X | M) \leq \mathbb{H}(X) - 10 \cdot \frac{|W_2|}{(\log n)^2} \cdot \min\left(d', \frac{L}{d}\right).$$

From Corollary 9 and the fact that $d' \geq 2^k/20$, this means that:

$$\mathbb{H}(M) \geq \frac{1}{2} \cdot \frac{|W_2|}{(\log n)^2} \cdot \min\left(2^k, \frac{L}{d}\right).$$

Now, from Equation (1) and as $|W_2| \geq \frac{1}{2} \cdot |U'| \geq \frac{1}{10} \cdot |V'|$, we get:

$$\mathbb{H}(M) \geq \frac{1}{(\log n)^4} \cdot \sum_{v \in V} \min\left(d_{\text{out}}(v), \frac{L}{d}\right).$$

As $\mathbb{H}(M)$ is upper bounded by the logarithm of the size of the support of M , we have that the algorithm \mathcal{A} must use at least $\frac{1}{(\log n)^4} \cdot \sum_{v \in V} \min(d_{\text{out}}(v), \frac{L}{d})$ bits of memory.

5.2 Finishing the Proof

To finish the argument, we prove that the streaming Alg (that is a random-walk algorithm) implies a streaming algorithm that solves our search problem from Section 5.1 with probability at least 0.8. For this, recall that any graph G' in the support of \mathcal{G}' lies in $\mathcal{G}_{n,d}$. This means that we can run the algorithm Alg (with the start vertex being an arbitrary vertex $v_0 \in W_1$ and the order of edges being the same as in the search problem) and it will sample an L -step random walk from v_0 with distortion at most 0.1. Let $P_{G'} = (v_0, v_1, \dots, v_L)$ be the random variable denoting the path, written as a sequence of vertices, taken by an L -step random walk from v_0 .

► **Lemma 22.** *For all graphs G' , which we can equivalently view as a pair (X, Z) , it holds that:*

$$\Pr\left(P_{G'} \text{ has } < \frac{d}{(\log n)^3} \cdot \min\left(d', \frac{L}{d}\right) \text{ distinct edges from } \bigcup_{v \in Z} X_v\right) < 0.1$$

Proof. Note that when the walk is at any vertex $v \in W_1$, the edge from v will either go to another vertex in W_1 (when it is a fixed edge added in Item 1 in Section 5.1) or it will go to a vertex in $Z \subseteq W_2$ (when it is a Z -edge added in Item 3). In the latter case, either the second edge will go to a vertex in W_1 (a fixed edge) or it will go to a vertex in W_3 (when it is an X -edge added in Item 2). If it does go to a vertex in W_3 , the third edge has to be fixed and take it to a vertex in W_1 . Thus, at least one of any three consecutive vertices in $P_{G'}$ must lie in W_1 .

In particular, if we break the random walk into chunks of length 10, then each such chunk $l \in [L/10]$ will contain at least two vertices in W_1 within it. For all $l \in [L/10]$, we define $u_l^{(1)}, u_l^{(2)}$ to be random variables denoting the positions (these are values in the interval $[(l-1) \cdot 10 + 1, (l-1) \cdot 10 + 6]$) of the first two vertices in W_1 in chunk l . As $|W_1| \leq 4d$ by Lemma 19, $|Z| = \frac{d}{20 \log n}$ and $|X_v| = d'/4 \geq d/40$ for all $v \in W_2$, we have by the above discussion that, for all $l \in [L/10]$:

$$\Pr\left(u_l^{(2)} - u_l^{(1)} = 3\right) \geq \frac{1}{1000 \log n}.$$

Because of the fact that $P_{G'}$ is a random walk, the same bound holds even when conditioned on $u_1^{(1)}, u_1^{(2)}, \dots, u_{l-1}^{(1)}, u_{l-1}^{(2)}$. This means that we can find mutually independent indicator random variables $Y_1, \dots, Y_{L/10}$, jointly distributed with $P_{G'}$ such that for all $l \in [L/10]$, we have $\Pr(Y_l = 1) = \frac{1}{1000 \log n}$ and $Y_l = 1 \implies u_l^{(2)} - u_l^{(1)} = 3$ with probability 1. From Lemma 2, it follows that:

$$\Pr\left(\sum_{l=1}^{L/10} \mathbb{1}\left(u_l^{(2)} - u_l^{(1)} = 3\right) \leq \frac{L}{(\log n)^2}\right) \leq \Pr\left(\sum_{l=1}^{L/10} Y_l \leq \frac{L}{(\log n)^2}\right) \leq \frac{1}{n^2},$$

as $L \geq (\log n)^{10}$ and n is large enough.

For the rest of this proof, we condition on values $(u_l^{(1)}, u_l^{(2)})_{l \in [L/10]}$ of $(u_l^{(1)}, u_l^{(2)})_{l \in [L/10]}$ such that the above event does not occur, and also condition on all the vertices visited in the walk except the vertices $(v_{u_l^{(1)}+1}, v_{u_l^{(1)}+2})_{l \in [L/10]: u_l^{(2)} - u_l^{(1)} = 3}$. Observe that, under this conditioning, for all $l \in [L/10]$ such that $u_l^{(2)} - u_l^{(1)} = 3$, the edge $(v_{u_l^{(1)}+1}, v_{u_l^{(1)}+2})$ is an independent and uniformly random edge in $\bigcup_{v \in Z} X_v$. As the total number of edges in

$\bigcup_{v \in Z} X_v$ is $\frac{dd'}{80 \log n}$ and these edges constitute at least $\frac{L}{(\log n)^2}$ independent uniformly random samples from it, the number of distinct samples will be at least $\frac{1}{100} \cdot \min\left(\frac{dd'}{80 \log n}, \frac{L}{(\log n)^2}\right) \geq \frac{d}{(\log n)^3} \cdot \min\left(d', \frac{L}{d}\right)$ except with probability at most 0.05. A union bound finishes the proof. ◀

As for any graph $G' \sim G$, the algorithm Alg samples from $P_{G'}$ with distortion at most 0.1, we have from Lemma 22 that the algorithm Alg solves the search problem with at least probability 0.8. From the result proved in Section 5.1, we get that the memory needed by the Alg after processing the X-edges is at least $\frac{1}{(\log n)^4} \cdot \sum_{v \in V} \min(d_{\text{out}}(v), \frac{L}{d})$. As the X-edges are always a subset of the edges in the graph G , there is a streaming order of the edges of G such that these are the edges that the algorithm Alg sees first when it is run on G . It follows that Alg needs memory at least $\frac{1}{(\log n)^4} \cdot \sum_{v \in V} \min(d_{\text{out}}(v), \frac{L}{d})$ when run on G and the proof of Theorem 16 is complete.

References

- 1 Peyman Afshani, Jérémy Barbay, and Timothy M. Chan. Instance-optimal geometric algorithms. *J. ACM*, 64(1):3:1–3:38, 2017. doi:10.1145/3046673.
- 2 Reid Andersen, Fan Chung, and Kevin Lang. Using pagerank to locally partition a graph. *Internet Mathematics*, 4(1):35–64, 2007. doi:10.1080/15427951.2007.10129139.
- 3 Reid Andersen and Yuval Peres. Finding sparse cuts locally using evolving sets. In *Symposium on Theory of computing (STOC)*, pages 235–244, 2009. doi:10.1145/1536414.1536449.
- 4 Aditya Bhaskara, Samira Daruki, and Suresh Venkatasubramanian. Sublinear algorithms for MAXCUT and correlation clustering. In *International Colloquium on Automata, Languages, and Programming (ICALP)*, volume 107 of *LIPICs*, pages 16:1–16:14, 2018. doi:10.4230/LIPICs.ICALP.2018.16.
- 5 Sergey Brin and Lawrence Page. The anatomy of a large-scale hypertextual web search engine. *Computer networks and ISDN systems*, 30(1-7):107–117, 1998. doi:10.1016/S0169-7552(98)00110-X.
- 6 Moses Charikar, Liadan O’Callaghan, and Rina Panigrahy. Better streaming algorithms for clustering problems. In *Symposium on Theory of computing (STOC)*, pages 30–39, 2003. doi:10.1145/780542.780548.
- 7 Lijie Chen, Gillat Kol, Dmitry Paramonov, Raghuvansh R. Saxena, Zhao Song, and Huacheng Yu. Near-optimal two-pass streaming algorithm for sampling random walks over directed graphs. In *International Colloquium on Automata, Languages, and Programming (ICALP)*, volume 198, pages 52:1–52:19, 2021. doi:10.4230/LIPICs.ICALP.2021.52.
- 8 Ashish Chiplunkar, John Kallaugher, Michael Kapralov, and Eric Price. Factorial lower bounds for (almost) random order streams. In *63rd IEEE Annual Symposium on Foundations of Computer Science, FOCS 2022, Denver, CO, USA, October 31 - November 3, 2022*, pages 486–497. IEEE, 2022. doi:10.1109/FOCS54457.2022.00053.
- 9 Ronald Fagin, Amnon Lotem, and Moni Naor. Optimal aggregation algorithms for middleware. *J. Comput. Syst. Sci.*, 66(4):614–656, 2003. doi:10.1016/S0022-0000(03)00026-6.
- 10 Juan A. Garay, Shay Kutten, and David Peleg. A sublinear time distributed algorithm for minimum-weight spanning trees. *SIAM journal on computing*, 27(1):302–316, 1998. doi:10.1137/S0097539794261118.
- 11 Oded Goldreich and Luca Trevisan. Three theorems regarding testing graph properties. *Random Struct. Algorithms*, 23(1):23–57, 2003. doi:10.1002/rsa.10078.
- 12 Mark Jerrum and Alistair Sinclair. Approximating the permanent. *SIAM journal on computing*, 18(6):1149–1178, 1989. doi:10.1137/0218077.
- 13 Mark R Jerrum, Leslie G Valiant, and Vijay V Vazirani. Random generation of combinatorial structures from a uniform distribution. *Theoretical computer science*, 43:169–188, 1986. doi:10.1016/0304-3975(86)90174-X.

- 14 Ce Jin. Simulating random walks on graphs in the streaming model. In Avrim Blum, editor, *Innovations in Theoretical Computer Science Conference (ITCS)*, volume 124 of *LIPICs*, pages 46:1–46:15, 2019. doi:10.4230/LIPICs.ITCS.2019.46.
- 15 Hossein Jowhari, Mert Saglam, and Gábor Tardos. Tight bounds for lp samplers, finding duplicates in streams, and related problems. In *Symposium on Principles of Database Systems (PODS)*, pages 49–58, 2011. doi:10.1145/1989284.1989289.
- 16 John Kallaugher, Michael Kapralov, and Eric Price. Simulating random walks in random streams. In *Symposium on Discrete Algorithms (SODA)*, pages 3091–3126, 2022. doi:10.1137/1.9781611977073.120.
- 17 Omer Reingold. Undirected connectivity in log-space. *Journal of the ACM (JACM)*, 55(4):1–24, 2008. doi:10.1145/1391289.1391291.
- 18 Tim Roughgarden, editor. *Beyond the Worst-Case Analysis of Algorithms*. Cambridge University Press, 2020. doi:10.1017/9781108637435.
- 19 Atish Das Sarma, Sreenivas Gollapudi, and Rina Panigrahy. Estimating pagerank on graph streams. *J. ACM*, 58(3):13:1–13:19, 2011. doi:10.1145/1970392.1970397.
- 20 Aaron Schild. An almost-linear time algorithm for uniform random spanning tree generation. In *Symposium on Theory of Computing (STOC)*, pages 214–227, 2018. doi:10.1145/3188745.3188852.
- 21 Daniel A Spielman and Shang-Hua Teng. A local clustering algorithm for massive graphs and its application to nearly linear time graph partitioning. *SIAM Journal on computing*, 42(1):1–26, 2013. doi:10.1137/080744888.
- 22 Gregory Valiant and Paul Valiant. An automatic inequality prover and instance optimal identity testing. *SIAM J. Comput.*, 46(1):429–455, 2017. doi:10.1137/151002526.
- 23 Jeffrey Scott Vitter. Random sampling with a reservoir. *ACM Trans. Math. Softw.*, 11(1):37–57, 1985. doi:10.1145/3147.3165.