

Graph Distances in the Streaming Model: The Value of Space *

Joan Feigenbaum ^{†‡} Sampath Kannan ^{§¶} Andrew McGregor ^{§||} Siddharth Suri ^{§**}
Jian Zhang ^{†††}

Abstract

We investigate the importance of space when solving problems based on graph distance in the streaming model. In this model, the input graph is presented as a stream of edges in an arbitrary order. The main computational restriction of the model is that we have limited space and therefore cannot store all the streamed data; we are forced to make space-efficient summaries of the data as we go along. For a graph of n vertices and m edges, we show that testing many graph properties, including connectivity (*ergo* any reasonable decision problem about distances) and bipartiteness, requires $\Omega(n)$ bits of space. Given this, we then investigate how the power of the model increases as we relax our space restriction. Our main result is an efficient randomized algorithm that constructs a $(2t + 1)$ -spanner in one pass. With high probability, it uses $O(t \cdot n^{1+1/t} \log^2 n)$ bits of space and processes each edge in the stream in $O(t^2 \cdot n^{1/t} \log n)$ time. We find approximations to diameter and girth via the constructed spanner. For $t = \Omega(\frac{\log n}{\log \log n})$, the space requirement of the algorithm is $O(n \cdot \text{polylog } n)$, and the per-edge processing time is $O(\text{polylog } n)$. We also show a corresponding lower bound of t for the approximation ratio achievable when the space restriction is $O(t \cdot n^{1+1/t} \log^2 n)$.

We then consider the scenario in which we are allowed multiple passes over the input stream. Here, we investigate whether allowing these extra passes will compensate for a given space restriction. We show that

finding vertices at distance d from a particular vertex will always take d passes, for all $d \in \{1, \dots, t/2\}$, when the space restriction is $o(n^{1+1/t})$. For girth, we show the existence of a direct trade-off between space and passes in the form of a lower bound on the product of the space requirement and number of passes. Finally, we conclude with two general techniques for speeding up the per-edge computation time of streaming algorithms while increasing the space by at most a log factor.

1 Introduction

In recent years, streaming has become an active area of research and an important paradigm for processing massive data sets [2, 17, 12, 18, 15, 13, 14]. Much of the existing work has focused on computing statistics for a stream of data elements, *e.g.*, frequency moments [2], L^p norms [12, 18], histograms [15, 13], and wavelet decompositions [14]. More recently, there have been extensions of the streaming research to the study of traditional graph problems [4, 11, 17]. Solving graph problems in this model has raised new challenges, because many existing approaches to the design of graph algorithms are rendered useless by the sequential-access limitation and the space limitation of the streaming model.

Many real world networks, *e.g.*, telecommunications networks and the WWW, can be modeled as massive graphs. Also, dense massive graphs may appear in applications such as structured data mining, where the relationships among the data items in the data set are represented as graphs. These massive graphs are more readily studied in a streaming, as opposed to random-access, model.

The input to a graph algorithm in the streaming model is a stream presenting the edge set of a graph. That is, for a graph $G(V, E)$ with vertex set $V = \{v_1, v_2, \dots, v_n\}$ and edge set $E = \{e_1, e_2, \dots, e_m\}$, a *graph stream* is the sequence $e_{i_1}, e_{i_2}, \dots, e_{i_m}$, where $e_{i_j} \in E$, and i_1, i_2, \dots, i_m is an arbitrary permutation of $\{1, 2, \dots, m\}$. Note that it is not required that all the edges incident on a vertex be grouped together in the stream.

The salient features of a streaming algorithm are:

*This work was supported by the DoD University Research Initiative (URI) administered by the Office of Naval Research under Grant N00014-01-1-0795.

†Yale University, Email: {feigenbaum-joan, zhangjian}@cs.yale.edu

‡Supported in part by ONR grant N00014-01-1-0795 and NSF grant ITR-0331548.

§University of Pennsylvania, Email: {kannan, andrewm, ssuri}@cis.upenn.edu

¶Supported in part by ARO grant DAAD 19-01-1-0473 and NSF grant CCR-0105337.

||Supported in part by NSF grant ITR-0205456.

**Supported in part by NIH grant T32 HG000046-05.

††Supported by NSF grant ITR-0331548.

1) sequential access to the data stream, 2) fast per-data-item processing time, 3) small space requirements compared to the length of the stream, and 4) access to the stream in a small number of passes. While all the existing streaming algorithms share these features, the quantities that define “small” in features (3) and (4) may vary from algorithm to algorithm. For example, most streaming algorithms use $\text{polylog } n$ space on a stream of length n . The streaming-clustering algorithm in [16], however, uses n^ϵ space. Also, the “streaming”¹ algorithm in [7] uses space \sqrt{n} . Most streaming algorithms access the input stream in one pass, but there are also multiple-pass algorithms [11, 21]. Thus, it is reasonable to have a broader view of the streaming model and treat the quantities in features (3) and (4) as parameters of the model. Note that it is also possible to view the quantity that defines “fast” in feature (2) as a parameter. We show that a simple technique can reduce the worst-case per-data-item complexity if the amortized complexity is small. We focus our attention on the two parameters in features (3) and (4), *i.e.*, the space $S(\cdot)$ and the number of passes $P(\cdot)$ required by the algorithm. The streaming model with these two parameters is a generalization of previous variations such the semi-streaming model [22, 11], in which the space restriction is $O(n \cdot \text{polylog } n)$, and the pass-efficient model [7]. For a better understanding of traditional graph-theory problems in the streaming model, it is necessary to investigate these problems with respect to different values of the parameters.

1.1 Our Contributions

Our main results in one-pass streaming are upper and lower bounds for all space restrictions on the computation of distances in undirected graphs $G = (V, E)$. For two vertices x and y in the graph G , the distance between x and y in G , the distance between x and y is the length of the shortest path between x and y in G . Distances can be approximated by graph spanners. For the graph $G = (V, E)$, a subgraph $H = (V, E')$, where $E' \subseteq E$, is a t -spanner of G if, for any vertices $x, y \in V$, $d_G(x, y) \leq d_H(x, y) \leq t \cdot d_G(x, y)$. The number t is the *stretch factor* of the spanner. Using the spanner, we can compute approximations for girth and diameter of G .

We devise a randomized streaming algorithm that constructs a $(2t + 1)$ -spanner for an unweighted undirected graph in one pass. With probability $1 - \frac{1}{n^{\Omega(1)}}$, the algorithm uses $O(t \cdot n^{1+1/t} \log^2 n)$ bits of space and processes each edge in the stream in $O(t^2 \cdot n^{1/t} \log n)$ time. Using this spanner, the all-pairs distances in the

graph can be $(2t + 1)$ -approximated. Note that, for $t = \log n / \log \log n$, the algorithm uses $O(n \cdot \text{polylog } n)$ bits of space and processes each edge in $O(\text{polylog } n)$ time. A complementary result shows that, with $O(t \cdot n^{1+1/t} \log^2 n)$ bits of space, we can not approximate the distance between u and v better than by a factor t even if we know the vertices u and v .

In [11], a simple semi-streaming spanner construction is briefly presented. This algorithm constructs a $\log n$ -spanner in one pass using $O(n \cdot \text{polylog } n)$ space. However, this algorithm needs $O(n)$ time to process each edge in the input stream. Such a per-edge processing time is prohibitive, especially in the streaming model. The work of [9] studies the construction of additive spanners in the streaming model. A subgraph $H = (V, E')$ is a (ϵ, β) -spanner of G if, for any vertices $x, y \in V$, $d_G(x, y) \leq d_H(x, y) \leq (1 + \epsilon) \cdot d_G(x, y) + \beta$. (Hence, H is also called “additive spanner.”) For large distances, additive spanner provides a better approximation. However, the algorithm of [9] requires multiple passes over the input stream, while our construction only needs one pass. For weighted undirected graphs, we also present a geometric grouping technique by which our algorithm can be extended to construct $((1 + \epsilon) \cdot (2t + 1))$ -spanners for weighted undirected graphs.

In [25], Thorup and Zwick provide a construction of distance oracles for approximating distances in graphs. Although *all-pairs-shortest-path* distances can be approximated using this oracle, their oracle construction requires the computation of shortest-path trees for certain vertices. Indeed, many constructions of spanners and distance oracles have this requirement, *i.e.*, they need to compute *some* distances between certain pairs of vertices in order to build a data structure from which the all-pairs-shortest-path distances can be approximated. Note that, in the streaming model, it is difficult to compute the distance between two particular vertices. Hence, a straightforward adaptation of such approaches will not work in the streaming model, and thus our algorithm takes a different approach.

For $S(n) = o(n)$, we identify the class of *balanced* properties. A graph property \mathcal{P} is *balanced* if, for all n sufficiently large, there exists a graph $G = (V, E)$ with $|V| = n$ and $u \in V$ such that the minimum of the size of the two sets $\{v : G = (V, E \cup \{(u, v)\}) \text{ has } \mathcal{P}\}$ and $\{v : G = (V, E \cup \{(u, v)\}) \text{ has } \neg \mathcal{P}\}$ is $\Omega(n)$. Using an information-theoretic approach it can be easily shown that testing any balanced graph property requires $\Omega(n)$ bits of space. Given that many basic graph properties, including connectivity and bipartiteness, are balanced, this result indicates that this area may not be very interesting for graph problems. Trivially, not being able

¹Instead of “streaming”, the authors of [7] use the term “pass-efficient algorithms.”

to test for connectivity rules out testing any property based on graph distance, including diameter and girth.

For truly massive graphs, in one pass, it may be impractical to expect that we have sufficient space to solve a graph problem. Hence, we also consider the scenario in which we are allowed multiple passes over the input stream with the idea of designing algorithms that work with tighter space constraints. We show that finding vertices at distance d from a particular vertex will always take d passes for all $d \in \{1, \dots, t/2\}$ when the space restriction is $o(n^{1+1/t})$. Thus, the construction of a BFS tree, even for a tree with small (constant) depth, requires multiple passes over the stream. Note that the construction of BFS trees (locally or globally) has been used intensively as an important subroutine in many approximations [1, 3, 8] of the distances in graphs. Such a lower bound pushes into the domain of multiple-pass streaming any algorithm that is based on constructing BFS trees.

For girth, we show the existence of a direct trade-off between space and number of passes when determining whether the girth of a graph is less than g . This takes the form of a lower bound on the product of the space S (in bits) and the number of passes P , i.e., $SP = \Omega\left((n/g)^{1+\frac{1}{g-5}}\right)$

Finally, we present a method for local amortization of per-data item complexity. We also present a technique for adapting existing partially dynamic graph algorithms to the semi-streaming model.

1.2 Structure of the Paper We organize the paper in three main sections. In Section 2, we demonstrate the computational power of our model as a function of the space restriction for one-pass streaming. We first show that balanced graph properties are hard to compute for algorithms with $S(n) = o(n)$. We then proceed to our main results on the estimation of distances in a graph. In Section 3, we describe the results in the multi-pass streaming domain. We show the multiple-pass lower bound for finding all the vertices that are a certain distance away from a particular vertex. We also show the tradeoff between $S(n)$ and $P(n)$ for determining girth. In Section 4, we present the buffering technique that reduces the worst-case per-data-item complexity if the amortized complexity is small.

2 The Importance of Space in One Pass Streaming

In this section we describe the jump in computational power of one pass, graph streaming algorithms that occurs at $\Omega(n)$ space.

2.1 $S(n) = o(n)$ We start by recognizing the entire domain of sub-linear (in n) space restriction as having no interesting computing power for graph problems. We identify a general type of graph property² and show that testing any such graph property requires $\Omega(n)$ space.

DEFINITION 2.1. *We say a graph property \mathcal{P} is balanced if there exists a constant $c > 0$ such that for all sufficiently large n , there exists a graph $G = (V, E)$ with $|V| = n$ and $u \in V$ such that:*

$$\begin{aligned} \min\{ & |\{v : G = (V, E \cup \{(u, v)\}) \text{ has } \mathcal{P}\}|, \\ & |\{v : G = (V, E \cup \{(u, v)\}) \text{ has } \neg\mathcal{P}\}\}| \} \\ & \geq cn \end{aligned}$$

Consider the scenario where the graph $G \setminus (u, v)$ is streamed in before the edge (u, v) , i.e. (u, v) is presented at the very end of the stream. Note that the “indexing” problem [17, 19] (Given a bit string as a stream and an indexing number at the end of the stream, the problem asks for the bit value in the string indexed by the indexing number.) can be reduced to testing a balanced property in such a graph stream. Hence,

LEMMA 2.1. *Testing for any balanced graph property \mathcal{P} requires $S(n) = \Omega(n)$.*

Note that many interesting properties are balanced e.g. connectivity, bipartiteness, does there exist a vertex of a certain degree, etc.

2.2 $S(n) = \Omega(n \cdot \text{polylog } n)$ For all space restrictions in this range, we explore the power of the model in regard to graph distances. We present upper and lower bounds for approximating these distances.

To begin with, we provide a streaming algorithm that constructs a $(2t + 1)$ -spanner for an unweighted undirected graph in one pass. We then extend this algorithm to construct $((1 + \epsilon) \cdot (2t + 1))$ -spanners for weighted undirected graphs, using a geometric grouping technique.

First, we give a high level view of the algorithm. Intuitively a dense graph will contain a small number of subgraphs (compared to n) that are dense and have small diameters. We can group the vertices in each of the subgraphs into a subset, and treat each subset as a super-node. If we remove, from the original graph, those subgraphs that have been grouped into super-nodes, the remaining graph, whose diameter may be large, will not have too many edges. We call the remaining graph the

²A graph property is simply a boolean function whose variables are the elements of the adjacency matrix of the graph but whose value is independent of the labeling of the nodes of the graph.

sparse part of the original graph. The induced graph on the super-nodes, and the sparse part of the original graph then form a spanner.

Note that the vertices in a super-node form a cluster of small diameter. Previous constructions [3, 8, 5] of such clusters all employ an approach similar to BFS, *i.e.*, the clusters are constructed layer by layer. In the streaming model, this would necessitate multiple passes over the input stream. We devise a randomized labeling scheme where the set of vertices with the same label form a cluster. By bypassing the BFS, this randomized construction enables our algorithm to run in one pass through the stream.

We now describe the generation of the labels in detail. A label l used in our construction is a positive integer. Given two parameters n and t , the set of labels L used by our algorithm is generated in the following way. Initially, we have the labels $1, 2, \dots, n$. We denote by L^0 this set of labels and call them the *level 0* labels. Independently, and with probability $\frac{1}{n^{1/t}}$, each label $l \in L^0$ will be selected into the set S^0 and l will be marked as *selected*. From each label l in S^0 , we generate a new label $l' = l + n$. We denote by L^1 the set of newly generated labels and call them level 1 labels. We then apply the above selection and new label generation procedure on L^1 to get the set of level 2 labels L^2 . We continue this until the level $\lfloor \frac{t}{2} \rfloor$ labels $L^{\lfloor \frac{t}{2} \rfloor}$ are generated. If a level $i+1$ label l is generated from a level i label l' , we call l the *successor* of l' and denote by $Succ(l') = l$. The set of labels we will use in our algorithm is the union of labels of level $1, 2, \dots, \lfloor \frac{t}{2} \rfloor$, *i.e.*, $L = \cup_0^t L^i$. Note that L can be generated before the algorithm sees the edges in the stream. But, in order to generate the labels in L , except in the case $t = O(\log n)$, the algorithm needs to know n , the number of vertices in the graph, before seeing the edges in the input stream. For $t = O(\log n)$, a simple modification of the above method can be used to generate L without knowing n , because the probability for a label to be selected is just $\frac{1}{2 \cdot O(1)}$.

At first glance it might appear that our labeling scheme resembles the sequences of sets initially constructed by the algorithm Thorup and Zwick [25]. In our construction the generation of the clusters of vertices depends on both the labels and the edge set. Thorup and Zwick, however, generate their sequence of sets independent of the edges. But, this is just the first step in their algorithm. Subsequently, their algorithm computes the exact distance between a fixed vertex and each of the sets of vertices in the sequence. Computing the exact distances between a pair of vertices is difficult in the streaming model. Our algorithm takes a very different approach from [25] to avoid this difficulty. \square

While going through the stream, our algorithm will label each vertex with labels chosen from L . The algorithm may label a vertex v with multiple labels, however, v will be labeled by at most one label from L^i , for $i = 1, 2, \dots, \lfloor \frac{t}{2} \rfloor$. Moreover, if v is labeled by a label l , and l is selected, the algorithm will also label v with the label $Succ(l)$.

Denote by l^i a label of level i , *i.e.*, $l^i \in L^i$. Let $L(v) = \{l^0, l^{k_1}, l^{k_2}, \dots, l^{k_j}\}$, $0 < k_1 < k_2 < \dots < k_j < t$ be the collection of labels that has been assigned to the vertex v . Let $Height(v) = \max\{j | l^j \in L(v)\}$ and $Top(v) = l^k \in L(v)$ s.t. $k = Height(v)$. Let $C(l)$ be the collection of vertices that are labeled with the label l .

The sets $L(v)$ and $C(l)$ will grow while the algorithm goes through the stream and labels the vertices. For each $C(l)$, our algorithm stores a BFS tree, $Tree(l)$, on the vertices of $C(l)$. We say an edge (u, v) connects $C(l)$ and $C(l')$ if u is labeled with l and v is labeled with l' . For some pairs of labels $l, l' \in L^{\lfloor \frac{t}{2} \rfloor}$, our algorithm will store edges that connects $C(l)$ and $C(l')$. We denote by H the set of such edges stored by our algorithm. In addition, for each vertex v , we denote by $M(v)$ the other edges incident to v that are stored by our algorithm. Intuitively, the subgraph induced by $\cup_{v \in V} M(v)$ is the sparse part of the graph G . The spanner constructed by the algorithm is the union of the BFS-trees for all the labels, $M(v)$ for all the vertices and the set H . The detailed algorithm is given in Figure 1, and we proceed with its analysis.

LEMMA 2.2. *At the end of the stream, for all $v \in V$, $|M(v)| \leq O(t \cdot n^{1/t} \log n)$ with probability $1 - \frac{1}{n^{\Omega(1)}}$.*

Proof. Let $M^{(i)}(v) \subseteq M(v)$ be the set of edges added to $M(v)$ during the period when $Height(v) = i$. Let $L(M(v)) = \cup_{(u,v) \in M(v)} L(u)$ be the set of labels that have been assigned to the vertices in $M(v)$. An edge (u, v) is added to $M(v)$ only in step 2(b)(ii). Note that in this case, $\lfloor \frac{t}{2} \rfloor \geq Height(u) \geq Height(v)$. Hence, the set $L_v(u)$ is not empty. Also by the condition in step 2(b)(ii), none of the labels in $L_v(u)$ appears in $L(M(v))$. Thus, by adding the edge (u, v) to $M(v)$, we introduce at least one new label to $L(M(v))$. $M^{(i)}(v)$ will then introduce a set B of distinct labels to $L(M(v))$. Furthermore, the size of B is at least $|M^{(i)}(v)|$. Note that the labels in B are not marked as selected. Otherwise, the algorithm would have taken step 2(b)(i) instead of step 2(b)(ii). Hence, the size of B follows a geometric distribution, *i.e.*, $\mathbb{P}(|B| = k) \leq (1 - \frac{1}{n^{1/t}})^k$. Thus, with probability $1 - \frac{1}{n^{\Omega(1)}}$, $|M^{(i)}(v)| \leq O(n^{1/t} \log n)$. Because i can take at most $O(t)$ values, $|M(v)| = \sum_i |M^{(i)}(v)|$ is at most $O(t \cdot n^{1/t} \log n)$ with probability $1 - \frac{1}{n^{\Omega(1)}}$. \square

Algorithm 1 (Efficient One Pass Spanner Construction).

The input to the algorithm is a unweighted undirected graph $G = (V, E)$, presented as a stream of edges, and two positive integer parameters n and t . (See the descriptions before Lemma 2.2 for the definition of H , $M(v)$, $L(v)$, $C(l)$, $Tree(l)$, $Succ(l)$, $Top(v)$ and $Height(v)$.)

1. Generate the set L of labels as described. $\forall v_i \in V$, label vertex v_i with label $i \in L^0$, and set $M(v) \leftarrow \emptyset$, $H \leftarrow \emptyset$.
2. Upon seeing an edge (u, v) in the stream, if $L(v) \cap L(u) = \emptyset$, consider the following cases:
 - (a) If $Height(v) = Height(u) = \lfloor \frac{t}{2} \rfloor$, and there is no edge in H that connects $C(Height(v))$ and $C(Height(u))$, set $H \leftarrow H \cup \{(u, v)\}$.
 - (b) Otherwise, without loss of generality, $\lfloor \frac{t}{2} \rfloor \geq Height(u) \geq Height(v)$. Consider the collection of labels $L_v(u) = \{l^{k_1}, l^{k_2}, \dots, l^{Height(u)}\} \subseteq L(u)$, where $k_1 \geq Height(v)$ and $k_1 < k_2 < \dots < Height(u)$. Let $l = l^i \in L_v(u)$ such that l^i is marked as selected and there is no $l^j \in L_v(u)$ with $j < i$ that is marked as selected.
 - i. If such a label l exists, label the vertex v with the successor $l' = Succ(l)$ of l , i.e., $L(v) \leftarrow L(v) \cup \{l'\}$. Incorporate the edge in the BFS-tree $Tree(l')$. If l' is selected, label v with $l'' = Succ(l')$ and incorporate the edge in the tree $Tree(l'')$. Continue this until we see an label that is not marked as selected.
 - ii. If no such label l exists and there is no edge (u', v) in $M(v)$ such that u, u' are labeled with the same label $l \in L_v(u)$, add (u, v) to $M(v)$, i.e., set $M(v) \leftarrow M(v) \cup \{(u, v)\}$.
3. After seeing all the edges in the stream, output the union of the BFS-trees for all the labels, $M(v)$ for all the vertices and the set H as the spanner.

Figure 1: An Efficient One Pass Algorithm for Computing Sparse Spanners

LEMMA 2.3. *The algorithm stores at most $O(t \cdot n^{1+1/t} \log n)$ edges, with probability $1 - \frac{1}{n^{\Omega(1)}}$.*

Proof. The algorithm stores edges in the set H , in the BFS trees for each cluster $C(l)$, and in the sets $M(v)$, $\forall v \in V$. By the Chernoff bound and the union bound, with probability $1 - \frac{1}{n^{\Omega(1)}}$, the number of clusters at level $t/2$ is $O(\sqrt{n})$ and the size of the set H is $O(n)$.

For each label l , the algorithm stores a BFS tree for the set of vertices $C(l)$. Note that for $i = 1, 2, \dots, \lfloor \frac{t}{2} \rfloor$, a vertex is labeled with at most one label in L^i . Hence, $\cup_{l \in L^i} C(l) \subseteq V$. Thus, the overall number of edges in the BFS trees is at most $O(t \cdot n)$.

By Lemma 2.2, with high probability, $|M(v)| \leq O(t \cdot n^{1/t} \log n)$. By the union bound, with high probability $\sum_{v \in V} |M(v)| \leq O(t \cdot n^{1+1/t} \log n)$. \square

THEOREM 2.1. *For an unweighted undirected graph of n vertices, presented as a stream of edges, and a positive number t , there is a randomized streaming algorithm that constructs a $(2t + 1)$ -spanner of the graph in one pass. With probability $1 - \frac{1}{n^{\Omega(1)}}$, the algorithm uses $O(t \cdot n^{1+1/t} \log^2 n)$ bits of space and processes each edge in the stream in $O(t^2 \cdot n^{1/t} \log n)$ time.*

Proof. Consider Algorithm 1 in Figure 1. At the beginning of the algorithm, for all the labels $l \in L^0$, $C(l)$ is a singleton set and the depth of the BFS tree for $C(l)$ is zero. We now bound, for label l^i , where $i > 0$, the depth of the BFS tree T^i on the vertices in $C(l^i)$. A tree grows when an edge (u, v) is incorporated into the tree in step 2(b)(i). Note that in this case, l^i is a successor of some label l^{i-1} of level $i - 1$. Assume that the depth $d_{T^i}(v)$ of the vertex v in the tree is one more than the depth $d_{T^i}(u)$ of the vertex u . Then $u \in C(l^{i-1})$, and the depth $d_{T^{i-1}}(u)$ of u in the BFS tree T^{i-1} of $C(l^{i-1})$ is the same as $d_{T^i}(u)$. Hence, $d_{T^i}(v) = d_{T^{i-1}}(u) + 1$ where T^i is a tree of level i and T^{i-1} is a tree of level $i - 1$. Given that $d_{T^0}(x) = 0$ for all $x \in V$, the depth of a BFS tree for $C(l)$, where l is a label of level i , is at most i .

We proceed to show that for any edge which the algorithm does not store, there is a path of length at most $2t + 1$ that connects the two endpoints of the edge. The algorithm ignores three types of edges. Firstly, if $L(u) \cap L(v) \neq \emptyset$, the edge (u, v) is ignored. In this case, let l be one of the label(s) in $L(u) \cap L(v)$, u and v are both on the BFS tree for $C(l)$. Hence, there is a path of length at most t connecting u and v . Secondly,

(u, v) will be ignored if $\text{Height}(v) = \text{Height}(u) = \lfloor \frac{t}{2} \rfloor$ and there is already an edge connecting $C(\text{Top}(u))$ and $C(\text{Top}(v))$. In this case, the path connecting u and v has a length at most $2t + 1$. Finally, in step 2(b)(ii), (u, v) will be ignored if there is already another edge in $M(v)$ that connects v to some $u' \in C(l)$ where $l \in L(u)$. Note that u and u' are both on the BFS tree of $C(l)$. Hence, there is a path of length at most $t + 1$ connecting u and v .

Hence, the stretch factor of the spanner constructed by Algorithm 1 is $2t + 1$. By Lemma 2.3, with high probability, the algorithm stores $O(t \cdot n^{1+1/t} \log n)$ edges and requires $O(t \cdot n^{1+1/t} \log^2 n)$ bits of space. Also note that the bottleneck in the processing of each edge lies on step 2(b)(ii), where for each label in $L_v(u)$, we need to examine the whole set of $M(v)$. This takes at most $O(t^2 \cdot n^{1/t} \log n)$ time. \square

Once the spanner is built, all-pairs-shortest distances of the graph can be computed from the spanner. The computation does not need to access the input stream and thus can be viewed as post-processing. Although Algorithm 1 constructs spanners for an unweighted undirected graph, we can use it, along with a geometric grouping technique to construct spanners for weighted undirected graphs.

Consider a weighted undirected graph $G = (V, E, \omega)$. Let ω_{\min} be the minimum weight and let ω_{\max} be the maximum weight. We divide the range $[\omega_{\min}, \omega_{\max}]$ into $\log_{(1+\epsilon)} \frac{\omega_{\max}}{\omega_{\min}}$ intervals of the form $[(1 + \epsilon)^i \omega_{\min}, (1 + \epsilon)^{i+1} \omega_{\min})$. Any weight in interval $[(1 + \epsilon)^i \omega_{\min}, (1 + \epsilon)^{i+1} \omega_{\min})$ we round down to $(1 + \epsilon)^i \omega_{\min}$. Let $G^i = (V, E^i)$ be the induced graph of G where E^i is the set of edges in E whose weight, after rounding, is $(1 + \epsilon)^i$. For each G^i we run a copy of Algorithm 1 in parallel, and the union of the spanners constructed for G^i forms a spanner for G . Note that this can be done without prior knowledge of ω_{\min} and ω_{\max} .

THEOREM 2.2. *For a weighted undirected graph of n vertices, whose maximum edge weight is ω_{\max} , and whose minimum edge weight is ω_{\min} , there is a randomized streaming algorithm that constructs a $((1 + \epsilon) \cdot (2t + 1))$ -spanner of the graph in one pass. With probability $1 - \frac{1}{n^{\Omega(1)}}$, the algorithm uses $O(\log_{1+\epsilon} \frac{\omega_{\max}}{\omega_{\min}} \cdot t \cdot n^{1+1/t} \log^2 n)$ bits of space and processes each edge in the stream in $O(t^2 \log_{1+\epsilon} \frac{\omega_{\max}}{\omega_{\min}} \cdot n^{1/t} \log n)$ time.*

In the case where $t = \frac{\log n}{\log \log n}$, Algorithm 1 computes a $(2 \frac{\log n}{\log \log n} + 1)$ -spanner in one pass. It uses $O(n \log^4 n)$ bits of space and processes each edge in $\log^4 n$ time. Once we have the spanner, the all-pair-shortest-path distances, as well as the diameter of the

graph, can be approximated. Also, it indirectly $\frac{2t+2}{3}$ approximates the girth of the graph since, if the constructed spanner is a strict subgraph of G we know the girth of G must have been between 3 and $2t + 2$.

2.2.1 A Lower Bound for One Pass Distance and Diameter Estimation If we are solely trying to estimate the distance between two nodes u and v , it may appear that constructing a graph spanner that gives no special attention to u and v , but rather approximates *all* distances, is an unnecessarily crude approach. In this section, however, we demonstrate that the spanner construction approach yields an approximation at most a factor 2 away from optimal.

A general technique to prove such results is showing the existence of *zombie* edges, edges whose presence or absence in the input graph can not be decided from the memory configuration. Their existence usually follows from an application of the pigeon hole principle: that for any algorithm, there exists two graphs G and G' that map to the same memory configuration and therefore there exists a zombie edge $e \in E(G) \setminus E(G')$. If we can make a zombie edge important in some way to the answer to the problem, we prove a limit on our ability to solve the problem.

For distance estimation we consider an edge $e = (u, v)$ important if it is *k-critical*, i.e. if the shortest path from u to v in $G \setminus (u, v)$ is length $\geq k$. In Lemma 2.4 we show the existence of graphs with many *k-critical* edges. Furthermore, we show these graphs may have relatively small diameter. By arguing about the existence of zombie *k-critical* edges, we then prove, in Theorem 2.3, lower bounds for the approximation ratio possible for both short and long distances.

LEMMA 2.4. *For $\epsilon > 0$ and sufficiently large n , there exists a graph $G = (V, E)$ with $|V| = n$, $|E| = (n^{1+\gamma})/4$ and unit edge weights such that the majority of edges are $\frac{1-\epsilon}{\gamma}$ -critical and the majority of the subgraphs in the set:*

$\{G' : G'$ formed from G by deleting at most half the *k-critical* edges $\}$.

have diameter less or equal to $\frac{1-\epsilon}{\gamma} + 1$. We call this family of subgraphs $\mathcal{F}(G)$.

Existence of such a family of graphs can be demonstrated by proving that almost all graphs in $\mathcal{G}_{n,n^{\gamma-1}}$ have the desired properties. We omit the details here but refer the reader to a similar proof in [11] for a slightly less general result. The following theorem generalizes a result in [11] to all space restrictions (as opposed to just the semi-streaming space restriction) and represents a tightening of the previous lower bound for the semi-

streaming restriction.

THEOREM 2.3. *In one pass and $S(n) = o(n^{1+\gamma})$, it is impossible to approximate the distance between two nodes $s, t \in V$ better than by a $\frac{1-\epsilon}{\gamma}$ ratio. Furthermore, this lower bound holds even if we are promised that $d_G(s, t)$ is equal to the diameter of G .*

Proof. Let $k = \frac{1-\epsilon}{\gamma}$, $D = k + 1$ and r , some arbitrarily large constant. Consider a graph G on n vertices with the properties in Lemma 2.4. Each graph in $\mathcal{F}(G)$ has $> n^{1+\gamma}/8$ edges and hence $|\mathcal{F}(G)| \geq 2^{n^{1+\gamma}/8-1} = 2^{\omega(S(n))}$. Therefore, for any algorithm, there exists a zombie edge which is k -critical. If, after streaming in $G_1 \in \mathcal{F}(G)$ and identifying zombie edge (t_1, s_1) we add edges (s, t_1) and (s_1, t) of zero weight, we can at best claim $d_G(s, t) \leq k$ when in fact $d_G(s, t)$ could be 1. This gives us the first result.

To prove a lower bound for approximating the diameter, consider a graph G on $\eta = \frac{n-2D}{rD}$ with the properties in Lemma 2.4. Our instance consists of rD disjoint graphs G_1, \dots, G_{rD} from $\mathcal{F}(G)$. For any algorithm, there exists a zombie edge (t_i, s_i) which is k -critical in each G_i . Finally we stream in edges (s_i, t_{i+1}) for $i = 2, \dots, rD-1$ of zero weight and two disjoint length D paths, one with end points s and t_1 and the other with end points s_D and t . Because of these two paths, the diameter is realized by the shortest path between s and t . Our construction gives a graph that has diameter $2D + rD$ while the algorithm can not guarantee that the diameter is less than $2D + rDk$. Hence the best approximation ratio is $> \frac{2D+rDk}{2D+rD} \rightarrow k$ for increasingly large r . \square

3 Trade-offs between Passes and Storage Space

3.1 Distances In this subsection we look at the hardness of finding the nodes at a specified distance from a particular node. This has obvious implications for the construction of BFS-trees. Let $\Gamma_i(s)$ be all the nodes up to distance i away from s . When $i = 1$ we omit the subscript. Let $S_i(s) = \Gamma_i(s) \setminus \Gamma_{i-1}(s)$.

Consider the following construction of a family of random graphs, $\mathcal{H}_{n,\epsilon,d}$, on $n = d\eta + 1$ vertices. Let $\epsilon \leq \frac{1}{2d}$. Partition the vertices into $d + 1$ layers, $A_0 = \{s\}, A_1, \dots, A_d$ where $|A_i| = \eta$. The graph has edges only between successive layers: s has one randomly chosen neighbor in A_1 and for $1 \leq i < d$, each node in A_i has exactly η^ϵ neighbors in A_{i+1} chosen uniformly at random. Let B_i be the subset of A_i that is exactly distance i away from s . The following lemma about the properties of this random construction follow from some reasonably straightforward applications of the Chernoff bound.

LEMMA 3.1. *For $\epsilon = \frac{1}{2d}$, $G \in \mathcal{H}_{n,\epsilon,d}$ has the following properties:*

1. *For $1 \leq i \leq d - 1$, $|\Gamma(B_i) \cup A_{i+1}| = \eta^\epsilon |B_i|$ with probability $q > 1 - \frac{1}{\Omega(n)}$.*
2. *For $1 \leq i \leq d - 1$, every subset S of size $\eta/4$ of A_i has $\Gamma(S) \cup A_{i+1} = A_{i+1}$ with probability $p > 1 - \frac{1}{\Omega(n)}$.*

LEMMA 3.2. *Let d satisfy $(\frac{n}{d})^{1+\frac{1}{2d}} = \omega(S(n))$ and let $\epsilon = 1/(2d)$. Partitioning A_d into B_d and $A_d \setminus B_d$ in a graph $G \in \mathcal{H}_{n,\epsilon,d}$, even given the layer number of each node, requires d passes in the streaming model with space restriction $S(n)$.*

Proof. [Sketch] Consider a randomly chosen instance constructed as above that has the properties of Lemma 3.1. Let E_i be the set of edges between A_{i-1} and A_i and observe, for $i > 1$, $|E_i| = \eta^{1+\epsilon}$. We stream the edges in the following order E_d, E_{d-1}, \dots, E_1 , and we allow unlimited computation per input edge. Now, $v \in B_i$ if and only if there exists $u \in B_{i-1}$ such that $(u, v) \in E_i$. We call u a *witness* for v . To certify that $v \in A_i \setminus B_i$ we need to know that v has no witnesses, *i.e.* $\Gamma(v) \cap B_{i-1} = \emptyset$. The result follows from the following claim.

Claim: For $1 \leq i \leq d$, classifying more than half the elements of A_i (as elements of B_i or of $A_i \setminus B_i$) requires at least $i - 1$ passes and the start $(E_d, E_{d-1}, \dots, E_i)$ of the i th pass.

Proof of Claim. The proof is by induction on i . Clearly the claim is true for $i = 1$. Assume it is true for $j - 1$.

Consider the first $j - 1$ passes of the stream. In each pass, when seeing edges from E_j there are still at least $\eta/2$ potential witnesses for vertices in B_j by the induction hypothesis. Hence, when seeing E_j we are unable to conclude any of the nodes in A_j are not in B_j by property 2 of G .

Consider the $j - 1$ th pass of the stream when we are now seeing edges from E_{j-1} . At this stage we may successfully partition A_{j-1} . An algorithm may have remembered some information about E_j , but, we argue, because of the memory restriction, that no algorithm could have remembered a sufficient amount to use our new found partition of A_{j-1} to deduce more than half of the elements of $A_j \setminus B_j$.

Denote by T the time when we stop seeing edges from E_j in the $j - 1$ th pass. Let $UP_{(i,T)}$ be the vertices in A_i that are unpartitioned at time T - these are all potential witnesses for vertices in A_{i+1} at this time. (NB. when we say a vertex in A_i is unpartitioned we

mean that the current memory state is insufficient to deduce whether it is in B_i or $A_i \setminus B_i$.)

Let $S_1 = UP_{(j-1,T)} \setminus \Gamma(B_j)$ and $S_2 = UP_{(j,T)}$. Let x'_i be the characteristic vector of the subset of A_j to which $v_i \in A_{j-1}$ is connected. Consider only v_i which are in S_1 and let $x_i = x'_i|_{S_2}$, the substring of x'_i on support set S_2 . This gives rise to an $|S_2| \times |S_1|$ matrix such that each column, x_i , has η^ϵ 1's (since every node in A_{j-1} has n^ϵ neighbors in A_j and if $v \in A_{j-1}$ has a neighbor in $A_j \setminus S_2$ then $v \in \Gamma(B_j)$ as $A_j \setminus S_2 \subset B_j$). We call this matrix M .

Consider the situation that, while seeing E_{j-1} we fully partition A_{j-1} and learn a new vertex v_l in B_{j-1} . Note there must be at least one such vertex since we know $|B_{j-1}| = \eta^{\epsilon(j-2)}$ and if, before seeing E_{j-1} , we had already found $\eta^{\epsilon(j-2)}$ nodes in B_{j-1} we would have concluded that nodes not in this set must be in $A_{j-1} \setminus B_{j-1}$ and we would have totally partitioned A_{j-1} in contradiction to the induction hypothesis. Now, v_l could have been any vertex in S_1 . To deduce more than a quarter of the vertices in S_2 are in $A_j \setminus B_j$, requires the information that some $|S_2|/4$ indices in the column x_l of M correspond to 0 entries. Since any algorithm could not have known v_l in advance, to do this for all possible v_l requires

$$\log p \left(\frac{|S_2|}{\eta^\epsilon} \right)^{|S_1|} - \log \left(\frac{|S_2|/4}{\eta^\epsilon} \right)^{|S_1|}$$

bits of information about E_j . By the induction hypothesis, this is $\Omega \left(\left(\frac{n}{d} \right)^{1+\frac{1}{2d}} \right)$. Thus, any algorithm that uses $o \left(\left(\frac{n}{d} \right)^{1+\frac{1}{2d}} \right)$ storage, cannot store enough information about E_j , and therefore it cannot guarantee that for any v_l , it learns some $|S_2|/4$ indices in the column x_l of M correspond to 0 entries. That is, any algorithm cannot guarantee that it partitions more than $|S_2|/4$ vertices in A_j into $A_j \setminus B_j$. Therefore the total number of vertices partitioned in A_j is $< |B_j| + |S_2|/4 \leq \eta/4 + \eta/4 = \eta/2$ since at time T less than $|B_j|$ vertices were partitioned. \square

Our main result here follows directly from the above lemma.

COROLLARY 3.1. *For all $d \in \{1, \dots, t/2\}$, computing $S_d(s)$ takes d passes when the space restriction is $o(n^{1+1/t})$.*

We note that this bound does not come directly from the multi-round communication complexity lower bound for the pointer chasing problem [24, 23, 6]. This problem implicitly defines a sparse, directed graph over n vertices. Thus, the communication complexity lower

bound for this problem can be no higher than $n \log n$, whereas in our model, we prove a higher lower bound depending on the value of t .

In regards to the semi-streaming space restriction of [11], we can not find $S_d(s)$ for $d < \frac{\log n}{(\log \log n)^2}$ in less than d passes. Conversely, note that with d passes we can easily find $S_d(s)$ even in the semi-streaming model - in the i th pass we can construct $S_i(s)$ from $S_{i-1}(s)$.

3.2 Girth In this section we prove a lower bound on the product $S(n)P(n)$ necessary to check whether a graph has girth $\leq g$. We shall make use of the following result from [20],

LEMMA 3.3. *Let $k \geq 1$ be an odd integer, $t = \lceil \frac{k+2}{4} \rceil$ and q be a prime power. There exists a bipartite, q -regular graph on $q \leq n \leq 2q^{k-t+1}$ nodes, with number of edges $\Omega(n^{1+1/(k-t+2)})$ and with girth $\geq k+5$.*

THEOREM 3.1. *To check whether a graph has girth $\leq g$, it is required $P(n)S(n) = \Omega \left((n/g)^{1+\frac{1}{g-5}} \right)$*

Proof. Let $G = (L \cup R, E)$ be a q regular bipartite graph with $|L| = |R| = n/g$ and $q = \Omega(n^{1/(g-5)})$. Let D_i be the set of neighbors of the i th node in L . We consider a family of graphs \mathcal{G} indexed by $(S_1, \dots, S_{n/g}, T_1, \dots, T_{n/g})$ where T_i, S_i are arbitrary subsets of D_i . Each graph consists of g , size n/g sets of nodes V_1, \dots, V_g . There is a perfect matching between V_i and V_{i+1} for $i = 1, \dots, \lceil g/2 \rceil, \lceil g/2 \rceil + 2, \dots, g$. Let $V_{\lceil g/2 \rceil+1} = \{u_1, \dots, u_{n/g}\}$, $V_{\lceil g/2 \rceil} = \{v_1, \dots, v_{n/g}\}$ and $V_1 = \{w_1, \dots, w_{n/g}\}$. Then we connect v_i to $\{u_j : j \in S_i\}$ and $\{w_j : j \in T_j\}$. The girth of G is g if $\sum_{i \in [n/g]} |T_i \cap S_i| \geq 1$ and at least $g+1$ otherwise.

Consider the nq/g bit string, A , induced by concatenating the characteristic vectors of the S_i . Likewise B for the T_i . Then the girth of the graph G is g if and only if A and B are not disjoint. The bound follows from the communication complexity of the set disjointness problem [17, 19]. \square

4 Towards Fast per-data item Processing

In section 2.2 we showed a spanner construction that processes each edge much faster than previous spanner construction algorithms. In this section we explore two general methods for decreasing the per-edge computation time of a streaming algorithm. As a consequence, we will show how some results from [10] give rise to efficient graph streaming algorithms.

Our first observation is that we can locally amortize per-edge processing by using some of our storage space as a *buffer* for incoming edges. While the algorithm processes a time intensive edge, subsequent edges can

be buffered subject to the availability of space. This potentially yields a decrease in the minimal allowable time between the arrival of incoming edges.

THEOREM 4.1. *Consider a streaming algorithm that runs in space $S(n)$ and uses computation time $\tau(m, n)$ to process the entire stream. This streaming algorithm can be simulated by a 1 pass streaming algorithm that uses $O(S(n) \log n)$ storage space, and has worst case time per-edge $\tau(m, n)/S(n)$.*

Next we turn to capitalizing on work done to speed up dynamic graph algorithms. *Dynamic graph algorithms* allow edges to be inserted and deleted in any order, and the current graph to be queried for \mathcal{P} at any point. *Partially dynamic algorithms*, on the other hand, are those that only allow edge insertions and querying. In [10] the authors describe a technique called sparsification, which they use to speed up existing dynamic graph algorithms that decide if a graph has property \mathcal{P} or not. Sparsification is based on maintaining strong certificates throughout the updates to the graph.

DEFINITION 4.1. *For any graph property \mathcal{P} , and graph G , a strong certificate for G is a graph G' on the same vertex set such that, for any H , $G \cup H$ has the property \mathcal{P} if and only if $G' \cup H$ has the property.*

It is easy to see that strong certificates obey a transitivity property. If G' is a strong certificate of property \mathcal{P} for graph G , and if G'' is a strong certificate for G' , then G'' is a strong certificate for G . Strong certificates also obey a compositional property. If G' and H' are strong certificates of \mathcal{P} for G and H , then $G' \cup H'$ is a strong certificate for $G \cup H$.

In order to achieve their speedup, the authors of [10], ensure that the certificates they maintain are not only strong but also sparse. A property is said to have *sparse certificates* if there is some constant c such that for every graph G on an n -vertex set, we can find a strong certificate for G with at most cn edges. Maintaining \mathcal{P} over a sparse certificate allows an algorithm to run over a dense graph, while only using the computational time of a sparse graph.

Now, in the streaming model, we are only concerned with edges being inserted and querying the property at the end of the stream. Moreover, observe that a sparse certificate fits in space, $S(n) = O(n \cdot \text{polylog } n)$. The following theorem states that any algorithm that could be sped up via the three major techniques described in [10], yields a one-pass, $O(n \cdot \text{polylog } n)$ space, streaming algorithm with the improved running time per input edge.

Problem	Time/Edge
Bipartiteness	$\alpha(n)$
Connected comps.	$\alpha(n)$
2-vertex connected comps.	$\alpha(n)$
3-vertex connected comps.	$\alpha(n)$
4-vertex connected comps.	$\log n$
MST	$\log n$
2-edge connected comps.	$\alpha(n)$
3-edge connected comps.	$\alpha(n)$
4-edge connected comps.	$n\alpha(n)$
Constant edge connected comps.	$n \log n$

Table 1: One-pass, $O(n \cdot \text{polylog } n)$ space streaming algorithms given by theorem 4.2.

THEOREM 4.2. *Let \mathcal{P} be a property for which we can find a sparse certificate in time $f(n, m)$. Then there exists a one-pass, semi-streaming algorithm that maintains a sparse certificate for \mathcal{P} using $f(n, O(n))/n$ time per edge.*

Proof. Let the edges in the stream be denoted by e_1, e_2, \dots, e_m . Let G_i denote the subgraph given by e_1, e_2, \dots, e_i . Inductively assume we have a sparse certificate C_{jn} for G_{jn} , where $1 \leq j \leq \lfloor m/n \rfloor$, constructed in time $f(n, O(n))/n$ per edge. Also, inductively assume that we have buffered the next n edges, $e_{jn+1}, e_{jn+2}, \dots, e_{(j+1)n}$. Let $T = C_{jn} \cup \{e_{jn+1}, e_{jn+2}, \dots, e_{(j+1)n}\}$. By the composition of strong certificates, T is a strong certificate for $G_{(j+1)n}$. Let $C_{(j+1)n}$ be the sparse certificate of T . By the transitivity of strong certificates, $C_{(j+1)n}$ is a sparse certificate of $G_{(j+1)n}$. Since C_{jn} is sparse, $|T| = (c+1)n$. Thus, computing $C_{(j+1)n}$ takes time $f(n, O(n))$. By theorem 4.1 this results in $f(n, O(n))/n$ time per edge, charged over $e_{jn+1}, e_{jn+2}, \dots, e_{(j+1)n}$. This computation can be done while the next n edges are being buffered. Let $k = \lfloor m/n \rfloor$, then the final sparse certificate will be $C_k \cup \{e_{k+1}, e_{k+2}, \dots, e_m\}$. \square

We note that for $f(n, m)$ which is linear or sub-linear in m , a better speed up may be achieved by buffering more than n edges, a possibility when we have more space. In [10] the authors provide many algorithms for computing various graph properties which they speed up using sparsification. Applying theorem 4.2 to these algorithms yields the list of streaming algorithms outlined in Table 1. For $l \geq 2$, the l -vertex and l -edge connectivity problems have either not been explicitly considering in the streaming model, or have algorithms with significantly slower time per edge (see [11]).

References

- [1] D. Aingworth, C. Chekuri, P. Indyk, and R. Motwani. Fast estimation of diameter and shortest paths(without matrix multiplication). *SIAM Journal on Computing*, 28(4):1167–1181, 1999.
- [2] N. Alon, Y. Matias, and M. Szegedy. The space complexity of approximating the frequency moments. *Journal of Computer and System Sciences*, 58(1):137–147, 1999.
- [3] B. Awerbuch, B. Berger, L. Cowen, and D. Peleg. Near-linear time construction of sparse neighborhood covers. *SIAM Journal on Computing*, 28(1):263–277, 1998.
- [4] Z. Bar-Yossef, R. Kumar, and D. Sivakumar. Reductions in streaming algorithms, with an application to counting triangles in graphs. In *Proc. 13th ACM-SIAM Symposium on Discrete Algorithms*, pages 623–632, 2002.
- [5] S. Baswana and S. Sen. A simple linear time algorithm for computing a $(2k - 1)$ -spanner of $o(n^{1+1/k})$ size in weighted graphs. In *Proc. 30th International Colloquium on Automata, Languages and Computing*, pages 284–296, 2003.
- [6] C. Damm, S. Jukna and J. Sgall. Some Bounds on Multiparty Communication Complexity of Pointer Jumping. *Computational Complexity*, 7(2):109–127, 1998.
- [7] P. Drineas and R. Kannan. Pass efficient algorithms for approximating large matrices. In *Proc. 14th ACM-SIAM Symposium on Discrete Algorithms*, pages 223–232, 2003.
- [8] M. Elkin. Computing almost shortest paths. In *Proc. 20th ACM Symposium on Principles of Distributed Computing*, pages 53–62, 2001.
- [9] M. Elkin and J. Zhang. Efficient algorithms for constructing $(1 + \epsilon, \beta)$ -spanners in the distributed and streaming models. in *Proc. 23th ACM Symposium on Principles of Distributed Computing*, pages 160–168, 2004.
- [10] D. Eppstein, Z. Galil, G. F. Italiano, and A. Nissenzweig. Sparsification – a technique for speeding up dynamic graph algorithms. *J. ACM*, 44(5):669–696, 1997.
- [11] J. Feigenbaum, S. Kannan, A. McGregor, S. Suri, and J. Zhang. On graph problems in a semi-streaming model. In *Proc. of 31st International Colloquium on Automata, Languages and Programming*, LNCS 3142, pages 531–543, 2004.
- [12] J. Feigenbaum, S. Kannan, M. Strauss, and M. Viswanathan. An approximate L^1 difference algorithm for massive data streams. *SIAM Journal on Computing*, 32(1):131–151, 2002.
- [13] A. C. Gilbert, S. Guha, P. Indyk, Y. Kotidis, S. Muthukrishnan, and M. Strauss. Fast, small-space algorithms for approximate histogram maintenance. In *Proc. 34th ACM Symposium on Theory of Computing*, pages 389–398, 2002.
- [14] A. C. Gilbert, Y. Kotidis, S. Muthukrishnan, and M. Strauss. Surfing wavelets on streams: One-pass summaries for approximate aggregate queries. In *Proc. 27th International Conference on Very Large Data Bases*, pages 79–88, 2001.
- [15] S. Guha, N. Koudas, and K. Shim. Data-streams and histograms. In *Proc. 33th ACM Symposium on Theory of Computing*, pages 471–475, 2001.
- [16] S. Guha, N. Mishra, R. Motwani, and L. O’Callaghan. Clustering data streams. In *Proc. 41th IEEE Symposium on Foundations of Computer Science*, pages 359–366, 2000.
- [17] M. R. Henzinger, P. Raghavan, and S. Rajagopalan. Computing on data streams. *Technical Report 1998-001, DEC Systems Research Center*, 1998.
- [18] P. Indyk. Stable distributions, pseudorandom generators, embeddings and data stream computation. In *Proc. 41th IEEE Symposium on Foundations of Computer Science*, pages 189–197, 2000.
- [19] E. Kushilevitz and N. Nisan. *Communication Complexity*. Cambridge University Press, 1997.
- [20] F. Lazebnik, V. Ustimenko, and A. Woldar. A new series of dense graphs of high girth. *Bulletin of the AMS*, 32(1):73–79, 1995.
- [21] J. Munro and M. Paterson. Selection and sorting with limited storage. *Theoretical Computer Science*, 12:315–323, 1980.
- [22] S. Muthukrishnan. Data streams: Algorithms and applications. 2003. <http://athos.rutgers.edu/~muthu/stream-1-1.ps>.
- [23] N. Nisan and A. Widgerson. Rounds in communication complexity revisited. In *Proceedings of the twenty-third annual ACM symposium on Theory of computing*, pages 419–429, 1991.
- [24] S. Ponzio, J. Radhakrishnan, and S. Venkatesh. The communication complexity of pointer chasing. *Journal of Computer and System Sciences*, 62(2):323–355, 2001.
- [25] M. Thorup and U. Zwick. Approximate distance oracles. In *Proc. 33th ACM Symposium on Theory of Computing*, pages 183–192, 2001.