

# Deterministic Independent Sets in the Semi-Streaming Model

Daniel Ye\*

February 8, 2025

## Abstract

We consider the independent set problem in the semi-streaming model. For any input graph  $G = (V, E)$  with  $n$  vertices, an independent set is a set of vertices with no edges between any two elements. In the semi-streaming model,  $G$  is presented as a stream of edges and any algorithm must use  $\tilde{O}(n)$ <sup>1</sup> bits of memory to output a large independent set at the end of the stream.

Prior work has designed various semi-streaming algorithms for finding independent sets. Due to the hardness of finding maximum and maximal independent sets in the semi-streaming model, the focus has primarily been on finding independent sets in terms of certain parameters, such as the maximum degree  $\Delta$ . In particular, there is a simple randomized algorithm that obtains independent sets of size  $\frac{n}{\Delta+1}$  in expectation, which can also be achieved with high probability using more complicated algorithms. For deterministic algorithms, the best bounds are significantly weaker. In fact, the best we currently know is a straightforward algorithm that finds an  $\tilde{\Omega}\left(\frac{n}{\Delta^2}\right)$  size independent set.

We show that this straightforward algorithm is nearly optimal by proving that any deterministic semi-streaming algorithm can only output an  $\tilde{O}\left(\frac{n}{\Delta^2}\right)$  size independent set. Our result proves a strong separation between the power of deterministic and randomized semi-streaming algorithms for the independent set problem.

---

\*(d29ye@uwaterloo.ca) School of Computer Science, University of Waterloo. Supported in part by Sepehr Assadi's Sloan Research Fellowship and NSERC Discovery grant.

<sup>1</sup> $\tilde{O}(\cdot)$  is used to hide polylog factors.

# 1 Introduction

Finding an independent set of a graph is a classical problem in graph theory with wide-ranging applications. An independent set of  $G = (V, E)$  is a subset  $I \subseteq V$  such that none of the vertices in  $I$  have an edge between them. Being able to find large independent sets plays a crucial role in network scheduling, transportation management, and much more. Especially in recent years, there has been an increased need for handling massive graphs to solve various tasks. Consequently, there has been an increased interest in studying the independent set problem in modern models of computation, such as the semi-streaming model introduced in [FKM<sup>+</sup>05]. In this model,  $G$  is presented as a stream of edges and the storage space of the algorithm is bounded to  $\tilde{O}(n)$ , where  $n = |V|$ . Our result pertains to the independent set problem under the semi-streaming model.

It is known that the Maximum Independent Set problem is both NP-hard [Kar10] and hard-to-approximate to within a factor of  $n^{1-\delta}$  for any  $\delta > 0$  [Zuc06]. Not only that, it is hard-to-approximate in the semi-streaming model, regardless of the time complexity of the algorithm [HSSW12]. As such, one line of research has focused on obtaining *Maximal* Independent Set. In this respect, [ACG<sup>+</sup>15] yields an  $O(\log \log n)$ -pass semi-streaming algorithm for this problem, and very recently, [AKNS24] proved this is the optimal number of passes.

Hence, to study algorithms in even fewer passes, the problem must be further relaxed to finding “combinatorially optimal” independent sets. On one hand, we can try finding an optimal bound with respect to degree sequences — the Caro-Wei theorem guarantees an independent set of size  $\sum_v \frac{1}{1+\deg(v)}$ , which is optimal in the sense that there are many graphs that do not admit larger independent sets. To achieve this bound in the semi-streaming model, [HHLS10] devises an algorithm for hypergraphs, which when applied to graphs, yields an expected independent set of size  $\Omega\left(\sum_v \frac{1}{\deg(v)}\right)$  with  $O(n)$  memory and  $O(1)$  time per edge.

On the other hand, if we consider bounds with respect to  $n$  and  $\Delta$  (the maximum degree in a graph), the Caro-Wei theorem guarantees independent sets of size  $\frac{n}{\Delta+1}$ , which is also tight. There is a very easy randomized semi-streaming algorithm for achieving this bound in expectation: randomly permute the vertices, and when an edge arrives in the stream, mark the endpoint that appears later. At the end, output all unmarked vertices. It is a standard result in graph theory that this algorithm outputs an independent set of size  $\frac{n}{\Delta+1}$ . We can even obtain such an independent set with high probability using the  $(\Delta+1)$ -coloring semi-streaming algorithm presented in [ACK19].

Interestingly, these results all make heavy use of randomization. Thus, as with a plethora of other problems in the semi-streaming model, there is a particular interest in derandomizing algorithms to achieve similar performance. Some problems admit single-pass deterministic algorithms matching their randomized counterparts (e.g. connectivity and bipartiteness [AGM12], maximum matching [PS18]). However, others have a proven discrepancy between the theoretical space complexities of deterministic algorithms and randomized ones under this model (e.g. triangle counting [BOV13] and vertex coloring [ACS22]).

For the independent set problem, the state-of-the-art has been a fairly simple deterministic algorithm for finding an independent set of size  $\frac{n}{\Delta^2}$ . We begin by choosing  $\frac{n}{\Delta}$  vertices arbitrarily, then store all the edges between them (resulting in  $O\left(\frac{n}{\Delta} \cdot \Delta\right) = O(n)$  edges). We can then calculate a maximal independent set of this subgraph, resulting in an independent set of size  $O\left(\frac{n}{\Delta^2}\right)$ . There has been very little progress in constructing a deterministic algorithm that does better than this, which has led to the following open question:

*Is there a single-pass semi-streaming deterministic algorithm that can match the performance of randomized ones? In particular, is there a deterministic algorithm that can find an independent set of size better than  $\tilde{O}\left(\frac{n}{\Delta^2}\right)$  in general graphs?*

## 1.1 Our Contributions

Our main result is a negative answer to the open question: the deterministic algorithm stated above is optimal up to polylog factors.

**Result 1.** *Any deterministic single-pass semi-streaming algorithm can only find an independent set of size at most  $\tilde{O}\left(\frac{n}{\Delta^2}\right)$  in a graph of maximum degree  $\Delta$  (even when  $\Delta$  is known).*

To the best of our knowledge, no space lower bounds have been devised for deterministic algorithms solving the independent set problem in the semi-streaming model. This result provides a lower bound that is tight up to logarithmic factors, illustrating another significant separation between deterministic and randomized algorithms in the semi-streaming model.

## 1.2 Our Techniques

We give a summary of our techniques here.

We begin our proof of the space lower bound in [Result 1](#) with a similar setup to [\[ACS22\]](#), who derived a space lower bound for deterministic algorithms solving the vertex coloring problem. In our case, we consider the multi-party communication complexity of the Independent Set problem, where the edges of a graph are partitioned among some number of players. In a predefined order, each player may speak once (outputting  $\tilde{O}(n)$  bits) and is heard by all future players. The goal is for the last player to output a large independent set.

We will design an adversary that adaptively constructs a graph that forces the algorithm to output an independent set of size  $\tilde{O}\left(\frac{n}{\Delta^2}\right)$ . To do so, it is useful to consider the graph of non-edges (later referred to as the missing graph): the graph consisting of edges each player *knows* has not been sent to them or any previous player. This model is useful because the independent set the last player outputs *must* be a clique in their missing graph (otherwise, there would be some input that makes this algorithm incorrect).

Additionally, we make use of the compression lemma of [\[ACS22\]](#). On any arbitrary graph, if we sample each edge with probability  $p$ , for any algorithm that compresses the result into  $s$  bits, the compression lemma finds a summary such that at most  $O\left(\frac{s}{p}\right)$  edges are not in any graph mapped to that summary. This is useful for [\[ACS22\]](#) as their adversary can narrow its search to a small set of vertices that a deterministic algorithm cannot summarize well. In fact, by focusing on vertices with low non-edge degree, it can sample each remaining edge with a higher probability to improve the bound given by the compression lemma. However, for the independent set problem, deterministic algorithms are free to choose vertices from *any* section of the graph (and do not need to deal with *every* vertex as in coloring). As such, our adversary does not have the luxury of searching for some useful set of vertices nor working with vertices with low non-edge degree. Instead, it must remove large independent sets globally from the graph and account for vertices that might not be easy to work with. To achieve this, our adversary generates graphs that are similar in structure to a Turán graph. In particular, as a deterministic algorithm receives its edges, the overall graph will seem more like many densely-connected vertex-disjoint subgraphs. The key to this strategy is a new lemma in our paper that allows for destroying large independent sets (equivalently, cliques) by adding “few” edges (equivalently, removing in the case of cliques).

## 1.3 Related Works

A similar-in-spirit result for vertex coloring is proven in [\[ACS22\]](#), which shows that deterministic algorithms cannot color a graph with  $\exp(\Delta^{o(1)})$  colors in a single pass. Since vertex coloring is

fairly hard, designing an adversary entails finding *some* set of vertices that is a clique from the perspective of a deterministic algorithm. With the independent set problem, however, we must design an adversary that removes *all* large independent sets from the perspective of a deterministic algorithm. Due to this difference, deterministic algorithms can easily find independent sets of size  $\frac{n}{\text{poly}(\Delta)}$  in our setting (whereas no deterministic algorithm can find a coloring using at most  $\text{poly}(\Delta)$  colors). We show that they cannot do better than quadratic, up to a polylog factor.

More generally, there has also been a significant interest in finding independent sets in graph streams [ACG<sup>+</sup>15, AKNS24, HHLS10, CCEW23, CDK17, BKO22, BCW20]. Independent sets in the *online streaming* model are studied in [HHLS16]. Under this model, they devise a deterministic algorithm with performance ratio  $O(2^\Delta)$ , which they prove is also tight. We provide an adjacent result in the *semi-streaming* model, which does not require an algorithm to maintain a feasible solution at all times. Additionally, [CDK19] studied independent sets in vertex-arrival streams, where each element in the input stream is a vertex along with its incident edges to earlier vertices. They show that the maximum independent set problem in the vertex-arrival model is not much easier than the problem in the edge-streaming model.

Independent sets have also been studied in more tangential settings. [CDK18] studied the problem of finding the Caro-Wei bound itself that other algorithms (such as the one in [HHLS10]) achieve, and [BKO22] studied the geometric independent set problem in the streaming model.

## 2 Preliminaries

**Notation.** For an integer  $k \geq 1$ , we denote  $[k] := \{1, 2, \dots, k\}$ . For a tuple  $(X_1, \dots, X_k)$  and any  $i \in [k]$ , we denote  $X_{<i}$  as  $(X_1, \dots, X_{i-1})$ . For any distribution  $\mu$ , we will denote  $\text{supp}(\mu)$  as the support of  $\mu$ .

For a graph  $G = (V, E)$ , we denote  $\Delta(G)$  as the maximum degree of  $G$ , and for any  $v \in V$ ,  $\deg(v)$  as the degree of  $v$  in  $G$ . For any vertex set  $T \subseteq V$ , we will denote the induced subgraph of  $G$  on  $T$  as  $G[T]$ . Often, we will also partition a set  $S$  into a collection of subsets  $\mathcal{P}$  and a subset  $Q$ . When we use this language, we are saying that  $\mathcal{P} \cup \{Q\}$  is a partition of  $S$ .

One part of our strategy also involves partitioning the vertices into many small subsets. For any integer  $g \geq 1$ , we will denote  $\text{PARTITION}(S, g)$  as an arbitrary partition of  $S$  into subsets of size  $g$ , *except* potentially for the last set, which has size  $< g$ .

**Fact 2.1.** For any set  $S$  and  $g \geq 1$ ,  $|\text{PARTITION}(S, g)| \leq \left\lceil \frac{|S|}{g} \right\rceil$ .

Finally, we use the following standard Chernoff bound:

**Proposition 2.2** (Chernoff bound; c.f. [DP09]). Suppose  $X_1, \dots, X_m$  are  $m$  independent random variables in the range  $[0, 1]$ . Let  $X := \sum_{i=1}^m X_i$  and  $\mu_L \leq \mathbb{E}[X] \leq \mu_H$ . Then, for any  $\epsilon > 0$ ,

$$\Pr(X > (1 + \epsilon) \cdot \mu_H) \leq \exp\left(-\frac{\epsilon^2 \cdot \mu_H}{3 + \epsilon}\right) \text{ and } \Pr(X < (1 - \epsilon) \cdot \mu_L) \leq \exp\left(-\frac{\epsilon^2 \cdot \mu_L}{2 + \epsilon}\right).$$

### 2.1 The Communication Complexity of Independent Sets

We prove our space lower bound in [Result 1](#) through a communication complexity argument in the following communication game, which as stated in [Section 1.2](#), is defined similarly to [\[ACS22\]](#).

For integers  $n, \Delta, k, s \geq 1$ , the  $\text{INDEPENDENT-SET}(n, \Delta, k, s)$  game is defined as:

1. There are  $k$  players  $P_1, \dots, P_k$ . Each Player  $P_i$  knows the vertex set  $V$  and receives a set  $E_i$  of edges. Let  $G = (V, E)$ , where  $E = E_1 \cup \dots \cup E_k$  and players are guaranteed  $E_1, \dots, E_k$  are disjoint. Players are guaranteed  $\Delta(G) \leq \Delta$ , and their goal is to output an Independent set of  $G$ .
2. In order from  $i = 1$  to  $i = k$ , each player  $P_i$  writes a public message  $M_i$  based on  $E_i$  and  $M_{<i}$  (all the messages from the previous players) of length at most  $s$ .
3. The goal of the players is to output an independent set of  $G$  by  $P_k$  outputting it as the message  $M_k$ .

The following is standard:

**Lemma 2.3.** *Suppose there is a deterministic streaming algorithm that, on any  $n$ -vertex graph  $G$  with known maximum degree  $\Delta$ , outputs an independent set of  $G$  with size  $r$  using  $s$  bits of space. Then, there also exists a deterministic protocol for  $\text{INDEPENDENT-SET}(n, \Delta, k, s)$  that outputs an independent set of size  $r$ .*

*Proof.* TOPROVE 0 ■

## 2.2 The Missing Graph and Compression Lemma

As stated in Section 1.2, we use the compression lemma in [ACS22]. We begin with two definitions that are similar to [ACS22].

**Definition 2.4.** *For a base graph  $G_{\text{Base}} = (V, E_{\text{Base}})$  and parameters  $p \in (0, 1]$ ,  $d \geq 1$ , we define the random graph distribution  $\mathbb{G} = \mathbb{G}(G_{\text{Base}}, p, d)$  as follows:*

1. *Sample a graph  $G$  on vertices  $V$  and edges  $E$  by picking each edge of  $E_{\text{Base}}$  independently with probability  $p$ .*
2. *Return  $G$  if  $\Delta(G) < 2p \cdot d$ . Otherwise, repeat the process.*

To analyze arbitrary deterministic algorithms, we will often consider the compression algorithm associated with it. We will represent the “information” available to the algorithm as a “missing graph”, which we define here:

**Definition 2.5.** *Consider  $\mathbb{G}(G_{\text{Base}}, p, d)$  for a base graph  $G_{\text{Base}} = (V, E_{\text{Base}})$  and parameters  $p \in (0, 1]$ ,  $d \geq 1$ , and an integer  $s \geq 1$ . A compression algorithm with size  $s$  is any function  $\Phi : \text{supp}(\mathbb{G}) \rightarrow \{0, 1\}^s$  that maps graphs sampled from  $\mathbb{G}$  into  $s$ -bit strings. For any graph  $G \in \text{supp}(\mathbb{G})$ , we refer to  $\Phi(G)$  as the summary of  $G$ . For any summary  $\Phi \in \{0, 1\}^s$ , we define:*

1.  $\mathbb{G}_\Phi$  *as the distribution of graphs mapped to  $\Phi$  by  $\Phi$ .*
2.  $G_{\text{Miss}}(\Phi) = (V, E_{\text{Miss}}(\Phi))$ , *called the missing graph of  $\Phi$ , as the graph on vertices  $V$  and edges missed by all graphs in  $\mathbb{G}_\Phi$ .*

The previous definitions are used extensively in our work. The following lemma also plays a crucial role in our communication lower bound, bounding the number of conclusive missing edges that can be recovered from a compression algorithm of a given size. It is proven in [ACS22].

**Lemma 2.6** (Compression Lemma). *Let  $G_{Base} = (V, E_{Base})$  be an  $n$ -vertex graph,  $s \geq 1$  be an integer, and  $p \in (0, 1)$  and  $d \geq 1$  be parameters such that  $d \geq \max\{\Delta(G_{Base}), \frac{4 \ln(2n)}{p}\}$ . Consider the distribution  $\mathbb{G} := \mathbb{G}(G_{Base}, p, d)$  and suppose  $\Phi : \text{supp}(\mathbb{G}) \rightarrow \{0, 1\}^s$  is a compression algorithm of size  $s$  for  $\mathbb{G}$ . Then, there exists a summary  $\phi^* \in \{0, 1\}^s$  such that in the missing graph of  $\phi^*$ ,*

$$|E_{Miss}(\phi^*)| \leq \frac{\ln 2 \cdot (s + 1)}{p}.$$

### 3 Removing Cliques in the Missing Graph

In this section, we introduce a key tool used by our adversary. More specifically, we develop a procedure that removes large cliques in the missing graph: the graph consisting of edges that an algorithm *knows* are not in the input graph. This is useful because a deterministic algorithm that finds independent sets must be *certain* that none of the vertices in its output have an edge between them, which creates a clique of the same size in its missing graph. To do so, we will design our adversary to generate a Turán-type graph, as mentioned in [Section 1.2](#).

#### 3.1 Removing Cliques in Low-Degree Graphs

The following lemma provides a method to bound the largest size of a clique when the degree is bounded by removing a small number of edges.

**Lemma 3.1.** *Let  $G$  be a graph with  $n$  vertices such that  $\Delta(G) \leq \Delta$ . For any positive integer  $d$ , there is some subgraph  $H$  of  $G$  such that:*

1. *The degree of  $H$  is  $\leq d$ .*
2. *The largest clique in  $G - H$  has size  $\leq 16 \ln(n) \cdot \frac{\Delta}{d} + 10$ .*

*Proof.* [TOPROVE 1](#) ■

#### 3.2 Removing Cliques in General Graphs

To use [Lemma 3.1](#), we need graphs of low maximum degree. However, the [Compression Lemma](#) can only bound the total number of edges. Hence, to employ [Lemma 3.1](#), we will partition a graph with few edges into many subgraphs with low degree and a small “remainder” subgraph. We will start with the following definition:

**Definition 3.2.** *For a positive integer  $b$  and a graph  $G = (V, E)$  with  $m$  edges, we define  $\text{SPLIT}(G, b)$  as a partition of  $V$  into a pair of vertex sets  $(P, Q)$  such that  $\Delta(G[P]) \leq b$  and  $|Q| \leq \frac{2m}{b}$ .*

The following proposition ensures that a split exists:

**Proposition 3.3.** *Suppose we are given a graph  $G = (V, E)$  with  $n$  vertices and  $m$  edges. For every positive integer  $b$ ,  $\text{SPLIT}(G, b)$  exists.*

*Proof.* [TOPROVE 2](#) ■

To “remove” many large cliques in an arbitrary graph, we will run a two-step subalgorithm.

The first step of the subalgorithm involves finding a subgraph of our input graph  $G_{Base}$ , which we denote as  $H_{Base}$ . The graph  $H_{Base}$  is chosen such that, for any algorithm compressing it, we

can easily bound the number of edges in  $H_{Miss}(\phi)$  (the missing graph of  $H_{Base}$  for some message  $\phi$ ) using the **Compression Lemma**. We will then prove that we can partition the vertices of  $H_{Miss}(\phi)$  (which are the same as the vertices of  $G_{Base}$ ) into a collection of vertex sets  $\mathcal{P}$  and a vertex set  $Q$  such that the maximum degree in  $(H_{Miss}(\phi))[P]$  is small for all  $P \in \mathcal{P}$  and the size of  $Q$  is small.

**Lemma 3.4.** *Let  $G_{Base} = (V, E_{Base})$  be a graph with  $n$  vertices. Let  $g \geq 1$  (group size) and  $s \geq 1$  (message size) be arbitrary integers. Let  $d_{comp} \geq 4 \ln(2n)$  (compression degree) and  $d_{filter} \geq 1$  (filter degree) be arbitrary real numbers.*

*Then, there is a subgraph  $H_{Base} \subseteq G_{Base}$  and a distribution  $\mathbb{H} := \mathbb{G}(H_{Base}, p, d)$  (for some  $p$  and  $d$ ) such that, for every compression algorithm  $\Phi : \text{supp}(\mathbb{H}) \rightarrow \{0, 1\}^s$ , we can find a message  $\phi$  and a partition of  $V$  into a collection of vertex sets  $\mathcal{P}$  and a vertex set  $Q$  satisfying:*

1. *For all  $H \in \text{supp}(\mathbb{H})$ ,  $\Delta(H) \leq 2 \cdot d_{comp}$ .*
2. *For all  $P \in \mathcal{P}$ ,  $G_{Base}[P] = H_{Base}[P]$ .*
3. *For any vertex set  $P \in \mathcal{P}$ ,  $\Delta(H_{Miss}(\phi)[P]) \leq d_{filter}$ .*
4. *The size of  $\mathcal{P}$  is  $\leq \lceil \frac{n}{g} \rceil$ .*
5. *The size of  $Q$  is  $\leq \frac{2 \ln(2) \cdot (s+1) \cdot g}{d_{comp} \cdot d_{filter}}$ .*

*Proof.* **TOPROVE 3** ■

In the second step of the subalgorithm, we will apply **Lemma 3.1** on  $H_{Miss}(\phi)[P]$  for all  $P \in \mathcal{P}$ , storing the removed edges in a subgraph  $R$  (which will have low degree). In the end, the maximum clique size in  $H_{Miss}(\phi) - R$  will be small, since it cannot exceed the *sum* of the maximum clique sizes over the subgraphs induced by the vertex sets in  $\mathcal{P}$ .

**Lemma 3.5.** *Suppose we have a graph  $H$  with  $n$  vertices and a partition of its vertices into a collection of vertex sets  $\mathcal{P}$  and a vertex set  $Q$ . Additionally, for all  $P \in \mathcal{P}$ , suppose that the degree of  $G[P]$  does not exceed  $d_{filter}$ .*

*Then, for any integer  $d_{remove} > 0$  (removal degree) there is a subgraph  $R \subseteq H$  such that:*

- *For all  $P \in \mathcal{P}$ , the largest clique in  $(H_{Miss}(\phi) - R)[P]$  has size  $\leq 16 \ln(n) \cdot \frac{d_{filter}}{d_{remove}} + 10$ .*
- *None of the edges in  $R$  is incident to a vertex in  $Q$ .*
- *The degree of  $R$  does not exceed  $d_{remove}$ .*

*Proof.* **TOPROVE 4** ■

## 4 A Communication Lower Bound for Independent Set

We will prove our lower bound in **Result 1** by showing that it is impossible for any set of players to output a large independent set, which is sufficient by **Lemma 2.3**. To do so, we will design an adversary that, for any large independent set, can find an *invalid* graph and set of edges to send to each player such that each player outputs the same message. This will ensure that no deterministic algorithm can confidently output any large independent set.



## 4.1 The Adversary

**Theorem 1.** *There exist constants  $\eta > 0$  and  $\eta_0 > 0$  such that: if  $n$ ,  $s$ , and  $\Delta$  are integers satisfying*

$$\eta < n \leq s, \quad \max \left\{ \eta_0 \cdot \frac{s \ln(n)}{n}, \eta_0 \cdot (\ln n)^2 \right\} < \Delta < \sqrt{n},$$

*then the size of the largest independent set a deterministic algorithm using  $s$  bits of memory can output for all graphs of size  $n$  and maximum degree  $\Delta$  does not exceed  $\tilde{O}\left(\frac{s}{\Delta^2} \cdot \frac{s}{n}\right)$ .*

Under the semi-streaming model,  $s = \tilde{O}(n)$ . Consequently, the largest independent set that a deterministic algorithm can find under this model has size  $\tilde{O}\left(\frac{n}{\Delta^2}\right)$ , which formalizes [Result 1](#).

To begin with the proof of [Theorem 1](#), we will let  $\eta = e^2$  and  $\eta_0 = 128$ . We let  $\ell = \max \left\{ \left\lceil \frac{2e \ln(2)(s+1)}{n} \right\rceil, \lceil 8 \ln n \rceil \right\}$ . By our choice of  $\eta$  and  $\eta_0$ , it is easy to prove that  $\ell < \frac{\Delta}{4 \ln(2n)}$ . We let  $k = \lceil \ln n \rceil + 1$ . For our adversary,  $k$  denotes the number of players, and we will generally send graphs of degree  $\leq \frac{\Delta}{\ell}$  to each player.

At a high level, our adversary adheres to the following structure:

- For all  $i = 1 \dots k$ ,  $G_{Base}(i)$  is the graph where most edges for Player  $i$  will be chosen from.
- Similarly,  $R_i$  will be provided as an additional set of edges to send to Player  $i$  to “manually” remove large cliques in the missing graph for the previous players’ messages.
- For all  $i \in [k]$ , Player  $i$  receives  $R_i$  and a subgraph  $H_i \subseteq G_{Base}(i)$ .
- [Lemma 3.4](#) is used to find an adversarial distribution of subgraphs  $\mathbb{H}_i$ , which is used to obtain  $H_i$ . [Lemma 3.5](#) is used to determine  $R_{i+1}$ , which is a subgraph of  $G_{Base}(i)$ .
- All graph parameters are derived adversarially based on what each player communicates. The compression algorithm associated with each player also affects the input for future players by determining  $G_{Base}(i+1)$  and  $R_{i+1}$ .
- The adversarially generated input graph,  $G_{input}$ , is the union of all graphs  $H_i$  and  $R_i$ .

An **adversary** that generates a “hard” input.

1. We start with  $G_{Base}(1)$  as the clique on  $n$  vertices and  $R_1$  as an empty graph.
2. For  $i = 1 \dots k$ ,
  - (a) Let  $n_i$  be the number of vertices in  $G_{Base}(i)$ , where  $G_{Base}(i) = (V_{Base}(i), E_{Base}(i))$ .
  - (b) If  $n_i < \frac{n}{\Delta^2}$ , then terminate the algorithm, letting  $\mathbb{H} = \emptyset$  and  $R_{i+1} = \emptyset$ . For the sake of our analysis, we let  $G_{Base}(i+1) = G_{Base}(i)$  and run the adversary until  $i = k$ .
  - (c) Otherwise, apply [Lemma 3.4](#) with  $G_{Base}(i)$ , group size  $\lfloor \frac{n_i \Delta^2}{n} \rfloor$ , message size  $s$ , compression degree  $\frac{\Delta}{\ell}$ , and filter degree  $\ell^2 \Delta$ .  
We have a distribution  $\mathbb{H}_i$  with a base graph  $H_{Base}(i)$ , and we define  $\Phi_i = \Phi(\mathbb{H}_i, M_{<i})$  as the compression algorithm for Player  $i$  after receiving  $R_i$ .  
By [Lemma 3.4](#), there is a message  $M_i := \phi$  and some partition of  $V_{Base}(i)$  into a collection of vertex sets  $\mathcal{P}_i$  and a vertex set  $Q_i$  satisfying the conclusions of [Lemma 3.4](#).



- (d) Apply **Lemma 3.5** with  $H_{Miss}(M_i)$  and removal degree  $\frac{\Delta}{2}$  to find some  $R_{i+1}$  such that the conclusions of **Lemma 3.5** holds.
  - (e) Let  $G_{Base}(i+1) = (G_{Base}(i) - (H_{Base}(i) - H_{Miss}(M_i)))[Q_i]$ .
  - 3. Finally, we generate the edges we send each player. For  $i = 1, \dots, k$ ,
    - (a) If  $\mathbb{H} = \emptyset$ , let  $H_i = \emptyset$ . Otherwise, choose  $H_i \in \text{supp}(\mathbb{H}_i)$  such that  $\Phi_i(H_i) = M_i$ .
    - (b) We send Player  $i$  the graph  $H_i \cup R_i$ .
- Then, the input graph is the union of these graphs. In particular,  $G_{input} := \bigcup_{i=1}^k H_i \cup R_i$ .

To begin, we must ensure that the input graph  $G_{input}$  is valid. The following lemma ensures that we are not sending a multigraph to the players.

**Lemma 4.1.** *No two players will receive the same edge (i.e.  $G_{input}$  is not a multigraph).*

*Proof.* **TOPROVE 5** ■

Next, we must also ensure that the maximum degree of the input graph agrees with the maximum degree given to each player.

**Lemma 4.2.** *For each vertex  $v \in G_{input}$ ,  $\deg(v) \leq \Delta$ .*

*Proof.* **TOPROVE 6** ■

**Lemma 4.1** and **Lemma 4.2** prove that the input graph is valid. The next step is to prove that the adversary terminates — for the sake of our analysis, we will instead show that the final base graph  $G_{Base}(k)$  is small.

**Lemma 4.3.**  $n_k \leq \frac{n}{\Delta^2}$ .

*Proof.* **TOPROVE 7** ■

Having proven that the input graph is valid and the algorithm terminates, it remains to bound the largest clique at each iteration. **Lemma 3.5** allows us to bound the largest clique in  $H_{Miss}(\phi)[P]$  for all  $P \in \mathcal{P}_i$ . The following lemma bounds the size of  $\mathcal{P}_i$ , which will allow us to bound the size of the largest clique in the missing graph over the entirety of  $\mathcal{P}_i$ .

**Lemma 4.4.** *For all  $i \in [k]$ , the size of  $\mathcal{P}_i$  is  $\leq \frac{3n}{\Delta^2}$ .*

*Proof.* **TOPROVE 8** ■

Finally, we prove that Player  $k$  cannot conclusively find a large independent set by proving that the adversary can find a “breaking” graph if the output is ever too large.

**Lemma 4.5.** *Suppose Player  $k$  outputs an Independent set  $A$  of size greater than*

$$\frac{n}{\Delta^2} + \frac{n}{\Delta^2} \cdot k \cdot (96\ell^2 \ln(n) + 30).$$

*Then, there is another graph  $G'_{input}$  and set of edges to send to each player such that each player outputs the same message but an edge exists between some  $u, v$  in the output of Player  $k$ .*

*Proof.* **TOPROVE 9** ■

**Lemma 4.5** leads naturally to a bound on the largest independent set a deterministic algorithm can find, which the following formalizes.

**Lemma 4.6.** *There does not exist a deterministic algorithm using at most  $s$  bits of memory that outputs an independent set of size greater than  $\frac{n}{\Delta^2} + \frac{n}{\Delta^2} \cdot k \cdot (96\ell^2 \ln(n) + 30)$  for all graphs with  $n$  vertices and maximum degree  $\Delta$ .*

*Proof.* **TOPROVE 10** ■

To conclude, we substitute  $\ell = \max \left\{ \left\lceil \frac{2e \ln(2)(s+1)}{n} \right\rceil, \lceil 8 \ln n \rceil \right\}$  and  $k = \lceil \ln n \rceil + 1$  into the bounds we showed in **Lemma 4.6**.

*Proof of Theorem 1.* This proof is simply a matter of loosening the bound in **Lemma 4.6** until we get a “simpler” form. We will start with the choices of  $k$  and  $\ell$ . For  $k$ , we have

$$k = \lceil \ln n \rceil + 1 \leq 2 \ln n + 1 \leq 3 \ln n.$$

For the first possible value of  $\ell$ , we have

$$\left\lceil \frac{2e \ln(2)(s+1)}{n} \right\rceil \leq \left\lceil \frac{4e \ln(2) \cdot s}{n} \right\rceil \leq 8e \ln(2) \cdot \frac{s}{n}.$$

For the other potential value of  $\ell$ , we have

$$\lceil 8 \ln n \rceil \leq 16 \ln n.$$

Hence,  $\ell \leq \max\{8e \ln 2 \cdot \frac{s}{n}, 16 \ln n\}$ . Since both terms are positive and greater than 1, we can multiply both upper bounds to get

$$\ell \leq 8e \ln 2 \cdot \frac{s}{n} \cdot 16 \ln n = 128e \ln 2 \cdot \frac{s \ln n}{n}.$$

To simplify calculations, we will let  $r = 128e \ln(2)$ , which is a constant.

By **Lemma 4.6**, no deterministic algorithm can find an independent set of size greater than

$$\begin{aligned} \frac{n}{\Delta^2} + \frac{n}{\Delta^2} \cdot k \cdot (96\ell^2 \ln(n) + 30) &\leq \frac{n}{\Delta^2} \left( 1 + 3 \ln(n) \cdot \left[ 96 \cdot \left( \frac{rs \ln n}{n} \right)^2 \cdot \ln(n) + 30 \right] \right) \\ &\leq \frac{n}{\Delta^2} \cdot 2 \cdot 6(\ln(n))^2 \left( 96 \cdot \left( \frac{rs \ln n}{n} \right)^2 \right) \\ &= 1152r^2 \ln^4(n) \cdot \frac{s^2}{n\Delta^2} \\ &= \tilde{O} \left( \frac{s}{\Delta^2} \cdot \frac{s}{n} \right). \end{aligned}$$

Hence, the size of the largest independent set that a deterministic algorithm can find for *all* graphs with  $n$  vertices and maximum degree  $\Delta$  does not exceed  $\tilde{O} \left( \frac{s}{\Delta^2} \cdot \frac{s}{n} \right)$ . ■

## Acknowledgements

I would like to thank Sepehr Assadi for introducing me to this problem and for his numerous discussions, comments, and feedback throughout the development of this work. Without him, this would not have been possible. I would also like to thank Vihan Shah, Janani Sundaresan, and Christopher Trevisan for their insightful comments and meticulous remarks on earlier versions of this work.

## References

- [ACG<sup>+</sup>15] KookJin Ahn, Graham Cormode, Sudipto Guha, Andrew McGregor, and Anthony Wirth, *Correlation clustering in data streams*, International Conference on Machine Learning, PMLR, 2015, pp. 2237–2246. [2](#), [4](#)
- [ACK19] Sepehr Assadi, Yu Chen, and Sanjeev Khanna, *Sublinear algorithms for  $(\Delta + 1)$  vertex coloring*, Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms, SIAM, 2019, pp. 767–786. [2](#)
- [ACS22] Sepehr Assadi, Andrew Chen, and Glenn Sun, *Deterministic graph coloring in the streaming model*, Proceedings of the 54th Annual ACM SIGACT Symposium on Theory of Computing, 2022, pp. 261–274. [2](#), [3](#), [4](#), [5](#)
- [AGM12] Kook Jin Ahn, Sudipto Guha, and Andrew McGregor, *Analyzing graph structure via linear measurements*, Proceedings of the twenty-third annual ACM-SIAM symposium on Discrete Algorithms, SIAM, 2012, pp. 459–467. [2](#)
- [AKNS24] Sepehr Assadi, Christian Konrad, Kheeran K Naidu, and Janani Sundaresan,  *$O(\log \log n)$  passes is optimal for semi-streaming maximal independent set*, Proceedings of the 56th Annual ACM Symposium on Theory of Computing, 2024, pp. 847–858. [2](#), [4](#)
- [BCW20] Ainesh Bakshi, Nadiia Chepurko, and David P. Woodruff, *Weighted maximum independent set of geometric objects in turnstile streams*, Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques, APPROX/RANDOM 2020, August 17-19, 2020, Virtual Conference (Jaroslav Byrka and Raghu Meka, eds.), LIPIcs, vol. 176, Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2020, pp. 64:1–64:22. [4](#)
- [BKO22] Sujoy Bhore, Fabian Klute, and Jelle J Oostveen, *On streaming algorithms for geometric independent set and clique*, International Workshop on Approximation and Online Algorithms, Springer, 2022, pp. 211–224. [4](#)
- [BOV13] Vladimir Braverman, Rafail Ostrovsky, and Dan Vilenchik, *How hard is counting triangles in the streaming model?*, Automata, Languages, and Programming: 40th International Colloquium, ICALP 2013, Riga, Latvia, July 8-12, 2013, Proceedings, Part I 40, Springer, 2013, pp. 244–254. [2](#)
- [CCEW23] Xiuge Chen, Rajesh Chitnis, Patrick Eades, and Anthony Wirth, *Sublinear-space streaming algorithms for estimating graph parameters on sparse graphs*, Algorithms and Data Structures Symposium, Springer, 2023, pp. 247–261. [4](#)

- [CDK17] Graham Cormode, Jacques Dark, and Christian Konrad, *Independent set size approximation in graph streams*, arXiv preprint arXiv:1702.08299 (2017). 4
- [CDK18] ———, *Approximating the caro-wei bound for independent sets in graph streams*, International Symposium on Combinatorial Optimization, Springer, 2018, pp. 101–114. 4
- [CDK19] Graham Cormode, Jacques Dark, and Christian Konrad, *Independent sets in vertex-arrival streams*, 46th International Colloquium on Automata, Languages, and Programming, ICALP 2019, July 9–12, 2019, Patras, Greece (Christel Baier, Ioannis Chatzigiannakis, Paola Flocchini, and Stefano Leonardi, eds.), LIPIcs, vol. 132, Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2019, pp. 45:1–45:14. 4
- [DP09] Devdatt P Dubhashi and Alessandro Panconesi, *Concentration of measure for the analysis of randomized algorithms*, Cambridge University Press, 2009. 4
- [FKM<sup>+</sup>05] Joan Feigenbaum, Sampath Kannan, Andrew McGregor, Siddharth Suri, and Jian Zhang, *On graph problems in a semi-streaming model*, Theoretical Computer Science **348** (2005), no. 2-3, 207–216. 2
- [HHLS10] Bjarni V Halldórsson, Magnús M Halldórsson, Elena Losievskaja, and Mario Szegedy, *Streaming algorithms for independent sets*, Automata, Languages and Programming: 37th International Colloquium, ICALP 2010, Bordeaux, France, July 6–10, 2010, Proceedings, Part I 37, Springer, 2010, pp. 641–652. 2, 4
- [HHLS16] ———, *Streaming algorithms for independent sets in sparse hypergraphs*, Algorithmica **76** (2016), 490–501. 4
- [HSSW12] Magnús M Halldórsson, Xiaoming Sun, Mario Szegedy, and Chengu Wang, *Streaming and communication complexity of clique approximation*, International Colloquium on Automata, Languages, and Programming, Springer, 2012, pp. 449–460. 2
- [Kar10] Richard M Karp, *Reducibility among combinatorial problems*, Springer, 2010. 2
- [PS18] Ami Paz and Gregory Schwartzman, *A  $(2 + \epsilon)$ -approximation for maximum weight matching in the semi-streaming model*, ACM Transactions on Algorithms (TALG) **15** (2018), no. 2, 1–15. 2
- [Zuc06] David Zuckerman, *Linear degree extractors and the inapproximability of max clique and chromatic number*, Proceedings of the thirty-eighth annual ACM symposium on Theory of computing, 2006, pp. 681–690. 2