# Randomized Composable Coresets for Matching and Vertex Cover

Sepehr Assadi[*]        Sanjeev Khanna[*]

## Abstract

A common approach for designing scalable algorithms for massive data sets is to distribute the computation across, say $k$, machines and process the data using limited communication between them. A particularly appealing framework here is the simultaneous communication model whereby each machine constructs a small representative summary of its own data and one obtains an approximate/exact solution from the union of the representative summaries. If the representative summaries needed for a problem are small, then this results in a *communication-efficient* and *round-optimal* (requiring essentially no interaction between the machines) protocol. Some well-known examples of techniques for creating summaries include sampling, linear sketching, and composable coresets. These techniques have been successfully used to design communication efficient solutions for many fundamental graph problems. However, two prominent problems are notably absent from the list of successes, namely, the *maximum matching* problem and the *minimum vertex cover* problem. Indeed, it was shown recently that for both these problems, even achieving a modest approximation factor of polylog($n$) requires using representative summaries of size $\widetilde{\Omega}(n^2)$ i.e. essentially no better summary exists than each machine simply sending its entire input graph.

The main insight of our work is that the intractability of matching and vertex cover in the simultaneous communication model is inherently connected to an *adversarial* partitioning of the underlying graph across machines. We show that when the underlying graph is randomly partitioned across machines, both these problems admit *randomized composable coresets* of size $\widetilde{O}(n)$ that yield an $\widetilde{O}(1)$-approximate solution[1]. In other words, a small *subgraph* of the input graph at each machine can be identified as its representative summary and the final answer then is obtained by simply running any maximum matching or minimum vertex cover algorithm on these combined subgraphs. This results in an $\widetilde{O}(1)$-approximation *simultaneous* protocol for these problems with $\widetilde{O}(nk)$ total communication when the input is randomly partitioned across $k$ machines. We also prove our results are optimal in a very strong sense: we not only rule out existence of smaller randomized composable coresets for these problems but in fact show that our $\widetilde{O}(nk)$ bound for total communication is optimal for *any* simultaneous communication protocol (i.e. not only for randomized coresets) for these two problems. Finally, by a standard application of composable coresets, our results also imply MapReduce algorithms with the same approximation guarantee in one or two rounds of communication, improving the previous best known round complexity for these problems.

[1]Here and throughout the paper, we use $\widetilde{O}(\cdot)$ notation to suppress polylog($n$) factors, where $n$ is the number of vertices in the graph.

# 1  Introduction

Recent years have witnessed tremendous algorithmic advances for efficient processing of massive data sets. A common approach for designing scalable algorithms for massive data sets is to distribute the computation across machines that are interconnected via a communication network. These machines can then jointly compute a function on the union of their inputs by exchanging messages. Two main measures of efficiency in this setting are the *communication cost* and the *round complexity*; we shall formally define these terms in details later in the paper but for the purpose of this section, communication cost measures the total number of bits exchanged by all machines and round complexity measures the number of rounds of interaction between them.

An important and widely studied framework here is the *simultaneous* communication model whereby each machine constructs a small representative summary of its own data and one obtains a solution for the desired problem from the union of the representative summary of combined pieces. The appeal of this framework lies in the simple fact that the *simultaneous protocols* are inherently *round-optimal*; they perform in only one round of interaction. The only measure that remains to be optimized is the communication cost – this is now determined by the size of the summary created by each machine. An understanding of the communication cost for a problem in the simultaneous model turns out to have value in other models of computation as well. For instance, a lower bound on the maximum communication needed by any machine implies a matching lower bound on the space complexity of the same problem in dynamic streams [7, 47].

Two particularly successful techniques for designing small summaries for simultaneous protocols are *linear sketches* and *composable coresets*. Linear sketching technique corresponds to taking a linear projection of the input data as its representative summary. The "linearity" of the sketches is then used to obtain a sketch of the combined pieces from which the final solution can be extracted. There has been a considerable amount of work in designing linear sketches for graph problems in recent years [5, 6, 10, 17, 18, 20, 40, 41, 50]. Coresets are subgraphs (in general, subsets of the input) that suitably preserve properties of a given graph, and they are said to be composable if the union of coresets for a collection of graphs yields a coreset for the union of the graphs. Composable coresets have also been studied extensively recently [11, 12, 15, 36, 52, 53], and indeed several graph problems admit natural composable coresets; for instance, connectivity, cut sparsifiers, and spanners (see [49], Section 2.2; the "merge and reduce" approach). Successful applications of these two techniques has yielded $\widetilde{O}(n)$ size summaries for many graph problems (see further related work in Section 1.3). However, two prominent problems are notably absent from the list of successes, namely, the *maximum matching* problem and the *minimum vertex cover* problem. Indeed, it was shown recently [10] that both matching and vertex cover require summaries of size $n^{2-o(1)}$ for even computing a polylog$(n)$-approximate solution[2].

This state-of-affairs is the starting point for our work, namely, intractability of matching and vertex cover in the simultaneous communication model. Our main insight is that a natural *data oblivious partitioning scheme* completely alters this landscape: both problems admit $\widetilde{O}(1)$-approximate composable coresets of size $\widetilde{O}(n)$ provided the edges of the graph are randomly partitioned across the machines. The idea that random partitioning of data can help in distributed computation was nicely illustrated in the recent work of [52] on maximizing submodular functions. Our work can be seen as the first illustration of this idea in the domain of graph algorithms. The applicability of this idea to graph theoretic problems has been cast as an open problem in [52].

---

[2]The authors in [10] only showed the inapproximability result for the matching problem. However, a simple modification of their result proves an identical lower bound for the vertex cover problem as well.

**Randomized Composable Coresets** We follow the notation of [52] with a slight modification to adapt to our application in graphs. Let $E$ be an edge-set of a graph $G(V, E)$; we say that a partition $\{E^{(1)}, \ldots, E^{(k)}\}$ of the edges $E$ is a *random k-partitioning* iff the sets are constructed by assigning each edge in $E$ independently to a set $E^{(i)}$ chosen uniformly at random. A random partitioning of the edges naturally defines partitioning the graph $G(V, E)$ into $k$ graphs $G^{(1)}, \ldots, G^{(k)}$ whereby $G^{(i)} := G(V, E^{(i)})$ for any $i \in [k]$, and hence we use random partitioning for both the edge-set and the input graph interchangeably.

**Definition** (Randomized Composable Coresets [52]). *For a graph-theoretic problem $P$, consider an algorithm* ALG *that given any graph $G(V, E)$, outputs a subgraph* ALG$(G) \subseteq G$ *with at most $s$ edges. Let $G^{(1)}, \ldots, G^{(k)}$ be a* random *k-partitioning of a graph $G$. We say that* ALG *outputs an $\alpha$-approximation randomized composable core-set of size $s$ for $P$ if $P\left(\mathsf{ALG}(G^{(1)}) \cup \ldots \cup \mathsf{ALG}(G^{(k)})\right)$ is an $\alpha$-approximation for $P(G)$ w.h.p., where the probability is taken over the random choice of the $k$-partitioning. For brevity, we use randomized coresets to refer to randomized composable coresets.*

We further augment this definition by allowing the coresets to also contain a *fixed solution* to be *directly* added to the final solution of the composed coresets. In this case, size of the coreset is measured both in the number of edges in the output subgraph plus the number of vertices and edges picked by the fixed solution (this is mostly relevant for our coreset for the vertex cover problem).

## 1.1 Our Results

We show existence of randomized composable coresets for matching and vertex cover.

**Result 1.** *There exist randomized coresets of size $\widetilde{O}(n)$ that w.h.p. (over the random partitioning of the input) give an $O(1)$-approximation for maximum matching, and an $O(\log n)$-approximation for minimum vertex cover.*

In contrast to the above result, when the graph is *adversarially* partitioned, the results of [10] show that the best approximation ratio conceivable for these problems in $\widetilde{O}(n)$ space is only $\Theta(n^{1/3})$. We further remark that Result 1 can also be extended to the weighted version of the problems. Using the Crouch-Stubbs technique [22] one can extend our result to achieve a coreset for weighted matching (with a factor 2 loss in approximation and extra $O(\log n)$ term in the space). Similar ideas of "grouping by weight" of edges can also be used to extend our coreset for weighted vertex cover with an $O(\log n)$ factor loss in approximation and space; we omit the details.

The $\widetilde{O}(n)$ space bound achieved by our coresets above is considered a "sweet spot" for graph streaming algorithms [30, 54] as many fundamental problems are provably intractable in $o(n)$ space (sometimes not enough to even store the answer) while admit efficient solutions in $\widetilde{O}(n)$ space. However, in the simultaneous model, these considerations imply only that the total size of all $k$ coresets must be $\Omega(n)$, leaving open the possibility that coreset output by each machine may be as small as $\widetilde{O}(n/k)$ in size (similar in spirit to coresets of [52]). Our next result rules out this possibility and proves the optimality of our coresets size.

**Result 2.** *Any $\alpha$-approximation randomized coreset for the matching problem must have size $\Omega(n/\alpha^2)$, and any $\alpha$-approximation randomized coreset for the vertex cover problem must have size $\Omega(n/\alpha)$.*

We now elaborate on some applications of our results.

**Distributed Computation**  We use the following distributed computation model in this paper, referred to as the *coordinator model* (see, e.g., [62]). The input is distributed across $k$ machines. There is also an additional party called the *coordinator* who receives no input. The machines are allowed to only communicate with the coordinator, not with each other. A protocol in this model is called a *simultaneous* protocol iff the machines simultaneously send a message to the coordinator and the coordinator then outputs the answer with no further interaction. *Communication cost* of a protocol in this model is the total number of bits communicated by all parties.

Result 1 can also be used to design simultaneous protocols for matching and vertex cover with $\widetilde{O}(nk)$ total communication and the same approximation guarantee stated in Result 1 in the case the input is partitioned randomly across $k$ machines. Indeed, each machine only needs to compute a coreset of its input, sends it to the coordinator, and coordinator computes an exact maximum matching or a 2-approximate minimum vertex cover on the union of the coresets. We further prove that the communication cost of theses protocols are essentially optimal.

> **Result 3.** *Any $\alpha$-approximation simultaneous protocol for the maximum matching problem, resp. the vertex cover problem, requires total communication of $\Omega(nk/\alpha^2)$ bits, resp. $\Omega(nk/\alpha)$ bits,* even *when the input is* partitioned randomly *across the machines.*

Result 3 is a strengthening of Result 2; it rules out *any* representative summary (not necessarily a randomized coreset) of size $o(n/\alpha^2)$ (resp. $o(n/\alpha)$) that can be used for $\alpha$-approximation of matching (resp. vertex cover) when the input is partitioned randomly.

For the matching problem, it was shown previously in [35] that when the input is adversarially partitioned in the coordinator model, any protocol (not necessarily simultaneous) requires $\Omega(nk/\alpha^2)$ bits of communication to achieve an $\alpha$-approximation of the maximum matching. Result 3 extends this to the case of *randomly partitioned* inputs albeit only for simultaneous protocols.

**MapReduce Framework**  We show how to use our randomized coresets to obtain improved MapReduce algorithms for matching and vertex cover in the MapReduce computation model formally introduced in [42,46]. Let $k = \sqrt{n}$ be the number of machines, each with a memory of $\widetilde{O}(n\sqrt{n})$; we show that *two* rounds of MapReduce suffice to obtain an $O(1)$-approximation for matching and $O(\log n)$-approximation for vertex cover. In the first round, each machine randomly partitions the edges assigned to it across the $k$ machines; this results in a random $k$-partitioning of the graph across the machines. In the second round, each machine sends a randomized composable coreset of its input to a designated central machine $M$; as there are $k = \sqrt{n}$ machines and each machine is sending $\widetilde{O}(n)$ size coreset, the input received by $M$ is of size $\widetilde{O}(n\sqrt{n})$ and hence can be stored entirely on that machine. Finally, $M$ computes the answer by combining the coresets (similar to the case in the coordinator model). Note that if the input was distributed randomly in the first place, we could have implemented this algorithm in only one round of MapReduce (see [52] for details on when this assumption applies).

Our MapReduce algorithm outperforms the previous algorithms of [46] for matching and vertex cover in terms of the number of rounds it uses, albeit with a larger approximation guarantee. In particular, [46] achieved a 2-approximation to both matching and vertex cover in 6 rounds of MapReduce when using similar space as ours on each machine (the number of rounds of this algorithm is always at least 3 even if we allow $\widetilde{O}(n^{5/3})$ space per each machine). The improvement on the number of rounds is significant in this context; the transition between different rounds in a MapReduce computation is usually the dominant cost of the computation [46] and hence, minimizing the number of rounds is an important goal in the MapReduce framework.

## 1.2 Our Techniques

**Randomized Coreset for Matching**   Greedy and Local search algorithms are the typical choices for composable coresets (see, e.g., [36,52]). It is then natural to consider the greedy algorithm for the maximum matching problem as a randomized coreset: the one that computes a *maximal matching*. However, one can easily show that this choice of coreset performs poorly in general; there are simple instances in which choosing arbitrary maximal matching in the graph $G^{(i)}$ results only in an $\Omega(k)$-approximation.

Somewhat surprisingly, we show that a simple change in strategy results in an efficient randomized coreset: *any maximum matching* of the graph $G^{(i)}$ can be used as an $O(1)$-approximate randomized coreset for the maximum matching problem. Unlike the previous work in [36,52] that relied on analyzing a specific algorithm (or a specific family of algorithms) for constructing a coreset, we prove this result by exploiting structural properties of the maximum matching (i.e., the optimal solution) directly, independent of the algorithm that computes it. As a consequence, our coreset construction requires no prior coordination (such as consistent tie-breaking rules used in [52]) between the machines and in fact each machine can use a different algorithm for computing the maximum matching required by the coreset.

**Randomized Coreset for Vertex Cover**   In the light of our coreset for the matching problem, one might wonder whether a minimum vertex cover of a graph can also be used as its randomized coreset. However, it is easy to show that the answer is negative here – there are simple instances (e.g., a star on $k$ vertices) on which this leads to an $\Omega(k)$ approximation ratio. Indeed, the *feasibility constraint* in the vertex cover problem depends heavily on the input graph as a whole and not only the coreset computed by each machine, unlike the case for matching and in fact most problems that admit a composable coreset [12,36,52]. This suggests the necessity of using edges in the coreset to *certify* the feasibility of the answer. On the other hand, only sending edges seems too restrictive: a vertex of degree $n-1$ can safely be assumed to be in an optimal vertex cover, but to certify this, one needs to essentially communicate $\Omega(n)$ edges. This naturally motivates a slightly more general notion of coresets – the coreset contains both subsets of vertices (to be always included in the final vertex cover) and edges (to guide the choice of additional vertices in the vertex cover).

To obtain a randomized coreset for vertex cover, we employ an iterative "peeling" process where we remove the vertices with the highest residual degree in each iteration (and add them to the final vertex cover) and continue until the residual graph is sufficiently sparse, in which case we can return this subgraph as the coreset. The process itself is a modification of the algorithm by Parnas and Ron [59]; we point out that other modifications of this algorithm has also been used previously for matching and vertex cover [16,38,58].

However, to employ this algorithm as a coreset we need to argue that the set of vertices peeled across different machines is not too large as these vertices are added directly to the final vertex cover. The intuition behind this is that random partitioning of edges in the graph should result in vertices to have essentially the same degree across the machines and hence each machine should peel the same set of vertices in each iteration. But this intuition runs into a technical difficulty: the peeling process is quite sensitive to the exact degree of vertices and even slight changes in degree results in moving vertices between different iterations that potentially leads to a cascading effect. To address this, we design a *hypothetical* peeling process (which is aware of the actual minimum vertex cover in $G$) and show that the our actual peeling process is in fact "sandwiched" between two application of this peeling process with different degree threshold for peeling vertices. We then use this to argue that the set of all vertices peeled across the machines are always contained in the solution of the hypothetical peeling process which in turn can be shown to be a relatively small set.

4

**Lower Bounds for Randomized Coresets.** Our lower bound results for randomized coresets for matching are based on the following simple distribution: the input graph consists of union of two bipartite graphs, one of which is a random $k$-regular graph $G_1$ with $n/2\alpha$ vertices on each side while the other graph $G_2$ is a perfect matching of size $n - n/2\alpha$. Thus the input graph almost certainly contains a matching of size $n - o(n)$ and any $\alpha$-approximate solution must collect $\Omega(n/\alpha)$ edges from $G_2$ overall i.e. $\Omega(n/\alpha k)$ edges from $G_2$ from each machine on average. After random partitioning, the input given to each machine is essentially a matching of size $n/2\alpha$ from $G_1$ and a matching of size roughly $n/k$ from $G_2$. The local information at each machine is not sufficient to differentiate between edges of $G_1$ and $G_2$, and thus any coreset that aims to include $\Omega(n/\alpha k)$ edges from $G_2$, can not reduce the input size by more than a factor of $\alpha$. Somewhat similar ideas can also be shown to work for the vertex cover problem.

**Communication Complexity Lower Bounds** We briefly highlight the ideas used in obtaining the lower bounds described in Result 3. We will focus on the vertex cover problem to describe our techniques. Our lower bound result is based on analyzing (a variant of) the following distribution: the input graph $G(L, R, E)$ consists of a bipartite graph $G_1$ plus a single edge $e^\star$. $G_1$ is a graph on $n/2\alpha$ vertices $L_1 \subseteq L$, each connected to $k$ random neighbors in $R$, and $e^\star$ is an edge chosen uniformly at random between $L \setminus L_1$ and $R$. This way $G$ admits a minimum vertex cover of size at most $n/2\alpha + 1$. However, when this graph is randomly partitioned, the input to each machine is essentially a matching of size $n/2\alpha$ chosen from the graph $G_1$ with possibly one more edge $e^\star$ (in exactly one machine chosen uniformly at random). The local information at the machine receiving the edge $e^\star$ is not sufficient to differentiate between the edges of $G_1$ and $e^\star$ and thus if the message sent by this machine is much smaller than its input size (i.e., $o(n/\alpha)$ bits), it most likely does not "convey enough information" to the coordinator about the identity of $e^\star$. This in turn forces the coordinator to use more than $n/2$ vertices in order to cover $e^\star$, resulting in an approximation factor larger than $\alpha$.

Making this intuition precise is complicated by the fact that the input across the players are highly correlated, and hence the message sent by one player, can also reveal extra information about the input of another (e.g. a relatively small communication from the players is enough for the coordinator to know the identity of entire $L_1$). To overcome this, we show that by conditioning on proper parts of the input, we can limit the correlation in the input of players and then use the *symmetrization* technique of [62] to reduce the simultaneous $k$-player vertex cover problem to a one-way two-player problem named the *hidden vertex problem* (HVP). Loosely speaking, in HVP, Alice and Bob are given two sets $S, T \subseteq [n]$, each of size $n/\alpha$, with the promise that $|S \setminus T| = 1$ and their goal is to find a set $C$ of size $o(n)$ which contains the single element in $S \setminus T$. We prove a lower bound of $\Omega(n/\alpha)$ bits for this problem using a subtle reduction from the well-known set disjointness problem. In this reduction, Alice and Bob use the protocol for HVP on "non-legal" instances (i.e., the ones for which HVP is not well-defined) to reduce the original disjointness instance between sets $A, B$ on a universe $[N]$ to a lopsided disjointness instance $(A, B')$ whereby $|B'| = o(N)$, and then solve this new instance in $o(N)$ communication (using the Håstad-Wigderson protocol [34]), contradicting the $\Omega(N)$ lower bound on the communication complexity of disjointness.

The lower bound for the matching problem is also proven along similar lines (over the hard distribution mentioned earlier for this problem) using a careful combinatorial argument instead of the reduction from the disjointness problem.

## 1.3 Further Related Work

Maximum matching and minimum vertex cover are among the most studied problems in the context of massive graphs including, in dynamic graphs [14,55,58,64], sub-linear algorithms [33,56,57,59,66], streaming algorithms [3–6, 9, 10, 20–22, 26–32, 37, 38, 43, 44, 48, 49, 51, 61], MapReduce computation [5,46], and different distributed computation models [8,24,32,35]. Most relevant to our work are the linear sketches of [20] for computing an *exact* minimum vertex cover or maximum matching in $O(\text{opt}^2)$ space (opt is the size of the solution), and linear sketches of [10,20] for $\alpha$-approximating maximum matching in $\widetilde{O}(n^2/\alpha^3)$ space. These results are proven to be tight by [21], and [10], respectively. Finally, [10] also studied the simultaneous communication complexity of bipartite matching in the vertex-partition model and proved that obtaining better than an $O(\sqrt{k})$-approximation in this model requires strictly more than $\widetilde{O}(n)$ communication from each player (see [10] for more details on this model).

Coresets, composable coresets, and randomized composable coresets are respectively introduced in [2], [36], and [52]. Composable coresets have been used previously in the context of nearest neighbor search [1], diversity maximization [36], clustering [12, 15], and submodular maximization [11, 36, 52]. Moreover, while not particularly termed a composable coreset, the "merge and reduce" technique in the graph streaming literature (see [49], Section 2.2) is identical to composable coresets. Similar ideas as randomized coreset for optimization problems has also been used in random arrival streams [38, 44]. Moreover, communication complexity lower bounds have also been studied previously under the random partitioning of the input [19, 39].

# 2 Preliminaries

**Notation.** For any integer $m$, $[m] := \{1, \ldots, m\}$. Let $G(V, E)$ be a graph; $\mathsf{MM}(G)$ denotes the maximum matching size in $G$ and $\mathsf{VC}(G)$ denotes the minimum vertex cover size. We assume that these quantities are $\omega(k \log n)$[3]. For a set $S \subseteq V$ and $v \in V$, $N_S(v) \subseteq S$ denotes the neighbors of $v$ in the set $S$. For an edge set $E' \subseteq E$, we use $V(E')$ to refer to vertices incident on $E'$.

**Useful Concentration of Measure Inequalities.** We use the following standard version of Chernoff bound (see, e.g., [25]) throughout.

**Proposition 2.1** (Chernoff bound). *Let $X_1, \ldots, X_n$ be independent random variables taking values in $[0, 1]$ and let $X := \sum_{i=1}^{n} X_i$. Then,*

$$\Pr\left(|X - \mathbb{E}\left[X\right]| > t\right) \leq 2 \cdot \exp\left(-\frac{2t^2}{n}\right)$$

We also need the method of bounded differences in our proofs. A function $f(x_1, \ldots, x_n)$ satisfies the *Lipschitz property* with constant $d$, iff for all $i \in [n]$, $|f(a) - f(a')| \leq d$, whenever $a$ and $a'$ differ only in the $i$-th coordinate.

**Proposition 2.2** (Method of bounded differences). *If $f$ satisfies the Lipschitz property with constant $d$ and $X_1, \ldots, X_n$ are independent random variables, then,*

$$\Pr\left(|f(X) - \mathbb{E}\left[f(X)\right]| > t\right) \leq 2 \cdot \exp\left(-\frac{2t^2}{n \cdot d^2}\right)$$

A proof of this proposition can be found in [25] (see Section 5).

---

[3]Otherwise, we can use the algorithm of [20] to obtain *exact* coresets of size $\widetilde{O}(k^2)$ as mentioned in Section 1.3.

**Communication Complexity** We prove our lower bounds for distributed protocols using the framework of communication complexity, and in particular in the *multi-party simultaneous communication model* and the *two-player one-way communication model*.

Formally, in the multi-party simultaneous communication model, the input is partitioned across $k$ players $P^{(1)}, \ldots, P^{(k)}$. All players have access to an infinite shared string of random bits, referred to as *public randomness* (or *public coins*). The goal is for the players to compute a specific function of the input by simultaneously sending a message to a central party called the coordinator (or the referee). The coordinator then needs to output the answer using the messages received by the players. We refer to the case when the input is partitioned randomly as the *random partition* model.

In the two-player one-way communication model, the input is partitioned across two players, namely Alice and Bob. The players again have access to public randomness, and the goal is for Alice to send a single message to Bob, so that Bob can compute a function of the joint input. The *communication cost* of a protocol in both models is the total length of the messages sent by the players. In Section 5.3.1, we also consider general two-player communication model which allows a *two-way* communication, i.e., both Alice and Bob can send messages to each other. We refer the reader to an excellent text by Kushilevitz and Nisan [45] for more details.

# 3   Randomized Coresets for Matching and Vertex Cover

We present our randomized composable coresets for matching and vertex cover in this section.

## 3.1   An $O(1)$-Approximation Randomized Coreset for Matching

The following theorem formalizes Result 1 for matching.

**Theorem 1.** *Any* maximum matching *of a graph $G(V, E)$ is an $O(1)$-approximation randomized composable coreset of size $O(n)$ for the maximum matching problem.*

We remark that our main interest in Theorem 1 is to achieve *some* constant approximation factor for randomized composable coresets of the matching problem and as such we did not optimize the constant in the approximation ratio. Nevertheless, our result already shows that the approximation ratio of this coreset is *at most* 9 (in fact, with a bit more care, we can reduce this factor down to 8; however, as this is not the main contribution of this paper, we omit the details).

Let $G(V, E)$ be any graph and $G^{(1)}, \ldots, G^{(k)}$ be a random $k$-partitioning of $G$. To prove Theorem 1, we describe a simple process for combining the maximum matchings (i.e., the coresets) of $G^{(i)}$'s, and prove that this process results in a constant factor approximation of the maximum matching of $G$. We remark that this process is only required for the analysis, i.e., to show that there exists a large matching in the union of coresets; in principle, any (approximation) algorithm for computing a maximum matching can be applied to obtain a large matching from the coresets.

Consider the following greedy process for computing an approximate matching in $G(V, E)$:

---

GreedyMatch($G$)**:**

1. Let $M^{(0)} := \emptyset$. For $i = 1$ to $k$:

2. Let $M^{(i)}$ be a *maximal matching* obtained by adding to $M^{(i-1)}$ the edges in an *arbitrary maximum matching* of $G^{(i)}$ that do not violate the matching property.

3. return $M := M^{(k)}$.

---

**Lemma 3.1.** GreedyMatch *is an $O(1)$-approximation algorithm for the maximum matching problem w.h.p (over the randomness of the edge partitioning).*

Before proving Lemma 3.1, we show that Theorem 1 easily follows from this lemma.

*Proof of Theorem 1.* Let ALG be any algorithm that given a graph $G(V, E)$, ALG$(G)$ outputs an arbitrary maximum matching of $G$. It is immediate to see that to implement GreedyMatch, we only need to compute a maximal matching on the output of ALG on each graph $G^{(i)}$ where $G^{(i)}$'s form a random $k$-partitioning of $G$. Consequently, since GreedyMatch outputs an $O(1)$-approximate matching (by Lemma 3.1), the graph $H := G^{(1)} \cup \ldots \cup G^{(k)}$ should contain an $O(1)$-approximate matching as well. We emphasize here that the use of GreedyMatch for finding a large matching in $H$ is *only* for the purpose of analysis. ∎

In the rest of this section, we prove Lemma 3.1. Recall that MM$(G)$ denotes the maximum matching size in the input graph $G$. Let $c > 0$ be a small constant to be determined later. To prove Lemma 3.1, we will show that $\left|M^{(k)}\right| \geq c \cdot$ MM$(G)$ w.h.p, where $M^{(k)}$ is the output of GreedyMatch. Notice that the matchings $M^{(i)}$ (for $i \in [k]$) constructed by GreedyMatch are random variables depending on the random $k$-partitioning.

Our general approach for the proof of Lemma 3.1 is as follows. Suppose at the beginning of the $i$-th step of GreedyMatch, the matching $M^{(i-1)}$ is of size $o($MM$(G))$. It is easy to see that in this case, there is a matching of size $\Omega($MM$(G))$ in $G$ that is entirely incident on vertices of $G$ that are not matched by $M^{(i-1)}$. We can further show that in fact $\Omega($MM$(G)/k)$ edges of this matching are appearing in $G^{(i)}$, *even* when we condition on the assignment of the edges in the first $(i-1)$ graphs. The next step is then to argue that the existence of these edges forces *any* maximum matching of $G^{(i)}$ to match $\Omega($MM$(G)/k)$ edges in $G^{(i)}$ between the vertices that are not matched by $M^{(i-1)}$; these edges can always be added to the matching $M^{(i-1)}$ to form $M^{(i)}$. This ensures that while the maximal matching in GreedyMatch is of size $o($MM$(G))$, we can increase its size by $\Omega($MM$(G)/k)$ edges in each of the first $k/3$ steps, hence obtaining a matching of size $\Omega($MM$(G))$ at the end. The following key lemma formalizes this argument.

**Lemma 3.2.** *For any $i \in [k/3]$, if $\left|M^{(i-1)}\right| \leq c \cdot$ MM$(G)$, then, w.p. $1 - O(1/n)$,*

$$\left|M^{(i)}\right| \geq \left|M^{(i-1)}\right| + \left(\frac{1 - 6c - o(1)}{k}\right) \cdot \mathsf{MM}(G)$$

To continue we define some notation. Let $M^\star$ be an arbitrary maximum matching of $G$. For any $i \in [k]$, we define $M^{\star < i}$ as the part of $M^\star$ assigned to the first $i - 1$ graphs in the random $k$-partitioning, i.e., the graphs $G^{(1)}, \ldots, G^{(i-1)}$. We have the following simple concentration result.

**Claim 3.3.** *W.p. $1 - O(1/n)$, for any $i \in [k]$,*

$$\left|M^{\star < i}\right| \leq \left(\frac{i - 1 + o(i)}{k}\right) \cdot \mathsf{MM}(G).$$

*Proof.* Fix an $i \in [k]$; each edge in $M^\star$ is assigned to $G^{(1)}, \ldots, G^{(i-1)}$, w.p. $(i - 1)/k$, hence in expectation, size of $M^{\star < i}$ is $\frac{i-1}{k} \cdot$ MM$(G)$. The claim now follows from a standard application of Chernoff bound (recall that, throughout the paper, we assume MM$(G) = \omega(k \log n)$). ∎

We now prove Lemma 3.2.

*Proof of Lemma 3.2.* Fix an $i \in [k/3]$ and the set of edges for $E^{(1)}, \ldots, E^{(i-1)}$; this also fixes the matching $M^{(i-1)}$ while the set of edges in $E^{(i)}, \ldots, E^{(k)}$ together with the matching $M^{(i)}$ are still random variables. We further assume that after fixing the edges in $E^{(1)}, \ldots, E^{(i-1)}$, $\left| M^{\star < i} \right| \leq \frac{i-1+o(i)}{k} \cdot \mathsf{MM}(G)$ which happens w.p. $1 - O(1/n)$ by Claim 3.3.

We first define some notation. Let $V_{\mathsf{old}}$ be the set of vertices incident on $M^{(i-1)}$ and $V_{\mathsf{new}}$ be the remaining vertices. Let $E^{\geq i}$ be the set of edges in $E \setminus \left( E^{(1)} \cup \ldots \cup E^{(i-1)} \right)$. We partition $E^{\geq i}$ into two parts: (i) $E_{\mathsf{old}}$: the set of edges with *at least one endpoint* in $V_{\mathsf{old}}$, and (ii) $E_{\mathsf{new}}$: the set of edges *incident entirely* on $V_{\mathsf{new}}$. Our goal is to show that w.h.p. *any* maximum matching of $G^{(i)}$ matches $\Omega(\mathsf{MM}(G)/k)$ vertices in $V_{\mathsf{new}}$ to each other by using the edges in $E_{\mathsf{new}}$; the lemma then follows easily from this.

Notice that the edges in the graph $G^{(i)}$ are chosen by independently assigning each edge in $E^{\geq i}$ to $G^{(i)}$ w.p. $1/(k-i+1)$.[4] This independence allows us to treat the edges in $E_{\mathsf{old}}$ and $E_{\mathsf{new}}$ separately; we can fix the set of sampled edges of $G^{(i)}$ in $E_{\mathsf{old}}$ denoted by $E_{\mathsf{old}}^i$ without changing the distribution of edges in $G^{(i)}$ chosen from $E_{\mathsf{new}}$. Let $\mu_{\mathsf{old}} := \mathsf{MM}(G(V, E_{\mathsf{old}}^i))$, i.e., the maximum number of edges that can be matched in $G^{(i)}$ using only the edges in $E_{\mathsf{old}}^i$. In the following, we show that w.h.p., there exists a matching of size $\mu_{\mathsf{old}} + \Omega(\mathsf{MM}(G)/k)$ in $G^{(i)}$; by the definition of $\mu_{\mathsf{old}}$, this implies that *any* maximum matching of $G^{(i)}$ has to use at least $\Omega(\mathsf{MM}(G)/k)$ edges in $E_{\mathsf{new}}$, proving the lemma.

Let $M_{\mathsf{old}}$ be any arbitrary maximum matching of size $\mu_{\mathsf{old}}$ in $G(V, E_{\mathsf{old}}^i)$. Let $V_{\mathsf{new}}(M_{\mathsf{old}})$ be the set of vertices in $V_{\mathsf{new}}$ that are incident on $M_{\mathsf{old}}$. We show that there is a large matching in $G(V, E_{\mathsf{new}})$ that avoids $V_{\mathsf{new}}(M_{\mathsf{old}})$.

**Claim 3.4.** *There exists a matching in $G(V, E_{\mathsf{new}})$ of size $\left( \frac{k-i+1-o(i)}{k} - 4c \right) \cdot \mathsf{MM}(G)$ that avoids the vertices of $V_{\mathsf{new}}(M_{\mathsf{old}})$.*

*Proof.* We first bound the size of $V_{\mathsf{new}}(M_{\mathsf{old}})$. Since any edge in $M_{\mathsf{old}}$ has at least one endpoint in $V_{\mathsf{old}}$, we have $|V_{\mathsf{new}}(M_{\mathsf{old}})| \leq |M_{\mathsf{old}}| \leq |V_{\mathsf{old}}|$. By the assertion of the lemma, $\left| M^{(i-1)} \right| < c \cdot \mathsf{MM}(G)$, and hence $|V_{\mathsf{new}}(M_{\mathsf{old}})| \leq |V_{\mathsf{old}}| < 2c \cdot \mathsf{MM}(G)$.

Moreover, by the assumption that $\left| M^{\star < i} \right| \leq \frac{i-1+o(i)}{k} \cdot \mathsf{MM}(G)$, there is a matching $M$ of size $\frac{k-i+1-o(i)}{k} \cdot \mathsf{MM}(G)$ in the graph $G(V, E^{\geq i})$. By removing the edges in $M$ that are either incident on $V_{\mathsf{old}}$ or $V_{\mathsf{new}}(M_{\mathsf{old}})$, at most $4c \cdot \mathsf{MM}(G)$ edges are removed from $M$. Now the remaining matching is entirely contained in $E_{\mathsf{new}}$ and also avoids $V_{\mathsf{new}}(M_{\mathsf{old}})$, hence proving the claim. ∎

We are now ready to finalize the proof. Let $M_{\mathsf{new}}$ be the matching guaranteed by Claim 3.4. Each edge in this matching is chosen in $G^{(i)}$ w.p. $1/(k-i+1)$ independent of the other edges; hence, by Chernoff bound (and the assumption that $\mathsf{MM}(G) = \omega(k \log n)$), there is a matching of size

$$(1 - o(1)) \cdot \left( \frac{1}{k} - \frac{o(i)}{k(k-i+1)} - \frac{4c}{k-i+1} \right) \cdot \mathsf{MM}(G)$$

$$\geq \left( \frac{1 - 6c - o(1)}{k} \right) \cdot \mathsf{MM}(G) \qquad (i \leq k/3)$$

in the edges of $M_{\mathsf{new}}$ that appear in $G^{(i)}$. This matching can be directly added to the matching $M_{\mathsf{old}}$, implying the existence of a matching of size $\mu_{\mathsf{old}} + \left( \frac{1-6c-o(1)}{k} \right) \cdot \mathsf{MM}(G)$ in $G^{(i)}$. As argued

---

[4]This is true even when we condition on the size of $\left| M^{\star < i} \right|$ since this event does not depend on the choice of edges in $E^{\geq i}$.

before, this ensures that any maximum matching of $G^{(i)}$ contains at least $\left(\frac{1-6c-o(1)}{k}\right) \cdot \mathsf{MM}(G)$ edges in $E_{\mathsf{new}}$. These edges can always be added to $M^{(i-1)}$ to form $M^{(i)}$, hence proving the lemma. ∎

Having proved Lemma 3.2, we can easily conclude Lemma 3.1.

*Proof of Lemma 3.1.* Recall that $M := M^{(k)}$ is the output matching of GreedyMatch. For the first $k/3$ steps of GreedyMatch, if at any step we obtained a matching of size $c \cdot \mathsf{MM}(G)$, then we are already done. Otherwise, at each step, by Lemma 3.2, w.p. $1 - O(1/n)$, we increase the size of the maximal matching by $\left(\frac{1-6c-o(1)}{k}\right) \cdot \mathsf{MM}(G)$ edges; consequently, by taking a union bound on the $k/3$ steps, w.p. $1 - o(1)$, the size of the maximal matching would be $\left(\frac{1-6c-o(1)}{3}\right) \cdot \mathsf{MM}(G)$. By picking $c = 1/9$, we ensure that in either case, the matching computed by GreedyMatch is of size at least $\mathsf{MM}(G)/9 - o(\mathsf{MM}(G))$, proving the lemma. ∎

## 3.2 An $O(\log n)$-Approximation Randomized Coreset For Vertex Cover

The following theorem formalizes Result 1 for vertex cover.

**Theorem 2.** *There exists an $O(\log n)$-approximation randomized composable coreset of size $O(n \log n)$ for the vertex cover problem.*

Let $G(V, E)$ be a graph and $G^{(1)}, \ldots, G^{(k)}$ be a random $k$-partitioning of $G$; we propose the following coreset for computing an approximate vertex cover of $G$. This coreset construction is a modification of the algorithm for vertex cover first proposed by [59].

---

VC-Coreset($G^{(i)}$). An algorithm for computing a composable coreset of each $G^{(i)}$.

1. Let $\Delta$ be the smallest integer such that $n/(k \cdot 2^\Delta) \leq 4 \log n$ and define $G_1^{(i)} := G^{(i)}$.

2. For $j = 1$ to $\Delta - 1$, let:

$$V_j^{(i)} := \left\{ \text{vertices of degree} \geq n/(k \cdot 2^{j+1}) \text{ in } G_j^{(i)} \right\}$$
$$G_{j+1}^{(i)} := G_j^{(i)} \setminus V_j^{(i)}.$$

3. Return $V_{\mathsf{cs}}^{(i)} := \bigcup_{j=1}^{\Delta-1} V_j^{(i)}$ as a *fixed solution* plus the graph $G_\Delta^{(i)}$ as the coreset.

---

In VC-Coreset we allow the coreset to, in addition to returning a subgraph, identify a set of vertices (i.e., $V_{\mathsf{cs}}^{(i)}$) to be added directly to the final vertex cover. In other words, to compute a vertex cover of the graph $G$, we compute a vertex cover of the graph $\bigcup_{i=1}^{k} G_\Delta^{(i)}$ and return it together with the vertices $\bigcup_{i=1}^{k} V_{\mathsf{cs}}^{(i)}$. It is easy to see that this set of vertices indeed forms a vertex cover of $G$: any edge in $G$ that belongs to $G^{(i)}$ is either incident on some $V_j^{(i)}$, and hence is covered by $V_j^{(i)}$, or is present in $G_\Delta^{(i)}$, and hence is covered by the vertex cover of $G_\Delta^{(i)}$.

In the remainder of this section, we bound the approximation ratio of this coreset. To do this, we need to prove that $\left| \bigcup_{i=1}^{k} V_{\mathsf{cs}}^{(i)} \right| = O(\log n) \cdot \mathsf{VC}(G)$. The bound on the approximation ratio then follows as the vertex cover of $\bigcup_{i=1}^{k} G_\Delta^{(i)}$ can be computed to within a factor of 2.

It is easy to prove (and follows from [59]) that the set of vertices $V_{\text{cs}}^{(i)}$ is of size $O(\log n) \cdot \text{VC}(G)$; however, using this fact directly to bound the size of $\bigcup_{i=1}^{k} V_{\text{cs}}^{(i)}$ only implies an approximation ratio of $O(k \log n)$ which is far worse than our goal of achieving an $O(\log n)$-approximation. In order to obtain the $O(\log n)$ bound, we need to argue that not only each set $V_{\text{cs}}^{(i)}$ is relatively small, but also that these sets are all intersecting in many vertices. In order to do so, we introduce a hypothetical algorithm (similar to VC-Coreset) on the graph $G$ and argue that the set $V_{\text{cs}}^{(i)}$ output by VC-Coreset($G^{(i)}$) is, with high probability, a subset of the output of this hypothetical algorithm. This allows us to then bound the size of the union of the sets $V_{\text{cs}}^{(i)}$ for $i \in [k]$.

Let $O^\star$ denote the set of vertices in an arbitrary optimum vertex cover of $G$ and $\overline{O^\star} := V \setminus O^\star$. Consider the following process on the original graph $G$ (defined only for analysis):

---

1. Let $G_1$ be the bipartite graph obtained from $G$ by removing edges between vertices in $O^\star$.

2. For $j = 1$ to $t := \lceil \log n \rceil$, let:

$$O_j := \big\{ \text{vertices in } O^\star \text{ of degree} \geq n/2^j \text{ in } G_j \big\}$$
$$\overline{O}_j := \big\{ \text{vertices in } \overline{O^\star} \text{ of degree} \geq n/2^{j+2} \text{ in } G_j \big\}$$
$$G_{j+1} := G_j \setminus (O_j \cup \overline{O}_j).$$

---

We first prove that the sets $O_j$'s and $\overline{O}_j$'s in this process form an $O(\log n)$ approximation of the minimum vertex cover of $G$ and then show that VC-Coreset($G^{(i)}$) (for any $i \in [k]$) is *mimicking* this hypothetical process in a sense that the set $V_{\text{cs}}^{(i)}$ is essentially *contained* in the union of the sets $O_j$'s and $\overline{O}_j$'s.

**Lemma 3.5.** $\left| \bigcup_{j=1}^{t} O_j \cup \overline{O}_j \right| = O(\log n) \cdot \text{VC}(G)$.

*Proof.* Fix any $j \in [t]$; we prove that $\overline{O}_j \leq 8 \cdot \text{VC}(G)$. The lemma follows from this since there are at most $O(\log n)$ different sets $\overline{O}_j$ and the union of the sets $O_j$'s is a subset of $O^\star$ (with size $\text{VC}(G)$).

Consider the graph $G_j$. The maximum degree in this graph is at most $n/2^{j-1}$ by the definition of the process. Since all the edges in the graph are incident on at least one vertex of $O^\star$, there can be at most $|O^\star| \cdot n/2^{j-1}$ edges between the remaining vertices in $O^\star$ and $\overline{O^\star}$ in $G_j$. Moreover, any vertex in $\overline{O}_j$ has degree at least $n/2^{j+2}$ by definition and hence there can be at most $\left( |O^\star| \cdot n/2^{j-1} \right) / \left( n/2^{j+2} \right) \leq 8 |O^\star| = 8 \cdot \text{VC}(G)$ vertices in $\overline{O}_j$, proving the claim. ∎

We now prove the main relation between the sets $O_j$'s and $\overline{O}_j$'s defined above and the intermediate sets $V_j^{(i)}$'s computed by VC-Coreset($G^{(i)}$). The following lemma is the heart of the proof.

**Lemma 3.6.** *Fix an* $i \in [k]$, *and let* $A_j = V_j^{(i)} \cap O^\star$ *and* $B_j = V_j^{(i)} \cap \overline{O^\star}$. *With probability* $1 - O(1/n)$, *for any* $t \in [\Delta]$:

1. $\bigcup_{j=1}^{t} A_j \supseteq \bigcup_{j=1}^{t} O_j$.

2. $\bigcup_{j=1}^{t} B_j \subseteq \bigcup_{j=1}^{t} \overline{O}_j$.

*Proof.* To simplify the notation, for any $t \in [\Delta]$, we let $A_{<t} = \bigcup_{j=1}^{t-1} A_j$ and $A_{\geq t} = \bigcup_{j=t}^{\Delta} A_j$ (and similarly for $B_j$'s, $O_j$'s, and $\overline{O}_j$'s). We also use $N_S(v)$ to denote the neighbor-set of the vertex $v$ in the set $S \subseteq V$.

Note that the vertex-sets of the graphs $G$ and $G^{(i)}$ are the same and we can "project" the sets $O_j$'s and $\overline{O}_j$'s on graph $G^{(i)}$ as well. In other words, we can say a vertex $v$ in $G^{(i)}$ belongs to $O_j$ iff $v \in O_j$ in the original graph $G$. In the following claim, we crucially use the fact that the graph $G^{(i)}$ is obtained from $G$ by sampling each edge w.p. $1/k$ to prove that the degree of vertices across different sets $O_j$'s (and $\overline{O}_j$'s) in $G^{(i)}$ are essentially the same as in $G$ (up to the scaling factor of $1/k$).

**Claim 3.7.** *For any $j \in [\Delta]$:*

- *For any vertex $v \in O_j$, $\left| N_{\overline{O}_{\geq j}}(v) \right| \geq n/(k \cdot 2^{j+1})$ in the graph $G^{(i)}$ w.p. $1 - O(1/n^2)$.*

- *For any vertex $v \in \overline{O}_{\geq j+1}$, $\left| N_{O_{\geq j}}(v) \right| < n/(k \cdot 2^{j+1})$ in the graph $G^{(i)}$ w.p. $1 - O(1/n^2)$.*

*Proof.* Fix any $j \in [\Delta]$ and $v \in O_j$. By definition of $O_j$, degree of $v$ is at least $n/2^j$ in $G_j$; in other words, $\left| N_{\overline{O}^{\geq j}}(v) \right| \geq n/2^j$ in the graph $G$. Since each edge in $G$ is sampled w.p. $1/k$ in $G^{(i)}$, $\left| N_{\overline{O}^{\geq j}}(v) \right| \geq n/(k \cdot 2^j)$ in $G^{(i)}$ in expectation. Moreover, by the choice of $\Delta$, $n/(k \cdot 2^j) \geq 4 \log n$, and hence by Chernoff bound, w.p. $1 - O(1/n^2)$, $\left| N_{\overline{O}^{\geq j}}(v) \right| \geq n/(k \cdot 2^{j+1})$ in $G^{(i)}$.

Similarly for a vertex $v \in \overline{O}^{\geq j+1}$, degree of $v$ is less than $n/2^{j+2}$ in $G_j$ by definition of $\overline{O}_j$; hence, $|N_{O^{\geq j}}(v)| < n/2^{j+2}$ in the graph $G$. Using a similar argument as before, by Chernoff bound, w.p. $1 - O(1/n^2)$, $|N_{O_{\geq j}}(v)| < n/(k \cdot 2^{j+1})$ in $G^{(i)}$. ∎

By using a union bound on the $n$ vertices in $G$, the statements in Claim 3.7 hold simultaneously for all vertices of $G$ w.p. $1 - O(1/n)$; in the following we condition on this event. We now prove Lemma 3.6 by induction.

Let $v$ be a vertex that belongs to $O_1$; we prove that $v$ belongs to the set $V_1^{(i)}$ of VC-Coreset, i.e., $v \in A_1$. By Claim 3.7 (for $j = 1$), the degree of $v$ in $G_1^{(i)}$ is at least $n/4k$. Note that in $G_1^{(i)}$, $v$ may also have edges to other vertices in $O^\star$ but this can only increase the degree of $v$. This implies that $v$ also belongs to $A_1$ by the threshold chosen in VC-Coreset. Similarly, let $u$ be a vertex in $\overline{O}_{\geq 2}$ (i.e., *not* in $\overline{O}_1$); we show that $u$ is not chosen in $V_1^{(i)}$, implying that $B_1$ can only contain vertices in $\overline{O}_1$. By Claim 3.7, degree of $u$ in $G_1^{(i)}$ is less than $n/4k$. This implies that $u$ does not belong to $B_1$. In summary, we have $O_1 \subseteq A_1$ and $B_1 \subseteq \overline{O}_1$.

Now consider some $t > 1$ and let $v$ be a vertex in $O_t$. By induction, $B_{<t} \subseteq \overline{O}_{<t}$. This implies that the degree of $v$ to $B_{\geq t}$ is at least as large as its degree to $O_{\geq t}$. Consequently, by Claim 3.7 (for $j = t$), degree of $v$ in the graph $G_t^{(i)}$ is at least $n/(k \cdot 2^{t+1})$ and hence $v$ also belongs to $A_t$. Similarly, fix a vertex $u$ in $\overline{O}_{\geq t+1}$. By induction, $A_{<t} \supseteq O_{<t}$ and hence the degree of $u$ to $A_{\geq t}$ is at most as large as its degree to $O_{\geq t}$; note that since $O^\star$ is a vertex cover, $u$ does not have any other edge in $G_t^{(i)}$ except for the ones to $A_{\geq t}$. We can now argue as before that $u$ does not belong to $B_t$. ∎

We are now ready to prove Theorem 2.

*Proof of Theorem 2.* The bound on the coreset size follows immediately from the fact that the graph $G_\Delta^{(i)}$ contains at most $O(n \log n)$ edges and size of $V_{\mathsf{cs}}^{(i)}$ is at most $n$. As argued before, to prove the bound on the approximation ratio, we only need to show that $\bigcup_{i=1}^{k} V_{\mathsf{cs}}^{(i)}$ is of size $O(\log n) \cdot \mathsf{VC}(G)$. Let $A^{(i)} = V_{\mathsf{cs}}^{(i)} \cap O^\star$ and $B^{(i)} = V_{\mathsf{cs}}^{(i)} \cap \overline{O^\star}$; clearly, each $A^{(i)} \subseteq O^\star$ and moreover, by Lemma 3.6 (for $t = \Delta$), each $B^{(i)} \subseteq \cup_{j=1}^{\Delta} \overline{O}_j$. Consequently, $\left| \bigcup_{i=1}^{k} V_{\mathsf{cs}}^{(i)} \right| \leq |O^\star| + \left| \bigcup_{j=1}^{\Delta} \overline{O}_j \right| \leq O(\log n) \cdot \mathsf{VC}(G)$, where the last inequality is by Lemma 3.5. ∎

# 4 Lower Bounds for Randomized Coresets

We formalize Result 2 in this section. As argued earlier, Result 2 is a special case of Result 3 and hence follows from that result; however, as the proof of Result 3 is rather technical and complicated, we also provide a self-contained proof of Result 2 as a warm-up to Result 3.

## 4.1 A Lower Bound for Randomized Composable Coresets of Matching

We prove a lower bound on the size of any randomized composable coreset for the matching problem, formalizing Result 2 for matching.

**Theorem 3.** *For any $k = o(n/\log n)$ and $\alpha = o(\min\{n/k, k\})$, any $\alpha$-approximation randomized composable coreset of the maximum matching problem is of size $\Omega(n/\alpha^2)$.*

By Yao's minimax principle [65], to prove the lower bound in Theorem 3, it suffices to analyze the performance of deterministic algorithms over a fixed (hard) distribution. We propose the following distribution for this task. For simplicity of exposition, in the following, we prove a lower bound for $(\alpha/4)$-approximation algorithms; a straightforward scaling of the parameters proves the lower bound for $\alpha$-approximation.

---

**Distribution $\mathcal{D}_{\mathsf{Matching}}$.** A hard input distribution for the matching problem.

- Let $G(L, R, E)$ (with $|L| = |R| = n$) be constructed as follows:

    1. Pick $A \subseteq L$ and $B \subseteq R$, each of size $n/\alpha$, uniformly at random.
    2. Define $E_{AB}$ as a set of edges between $A$ and $B$, chosen by picking each edge in $A \times B$ w.p. $k \cdot \alpha/n$.
    3. Define $E_{\overline{AB}}$ as a *random* perfect matching between $\overline{A}$ and $\overline{B}$.
    4. Let $E := E_{AB} \cup E_{\overline{AB}}$.

- Let $E^{(1)}, \ldots, E^{(k)}$ be a *random k-partitioning* of $E$ and let the input to player $P^{(i)}$ be the graph $G^{(i)}(L, R, E^{(i)})$.

---

Let $G$ be a graph sampled from the distribution $\mathcal{D}_{\mathsf{Matching}}$. Notice first that the graph $G$ always has a matching of size at least $n - n/\alpha \geq n/2$, i.e., the matching $E_{\overline{AB}}$. Additionally, it is easy to see that any matching of size more than $2n/\alpha$ in $G$ uses at least $n/\alpha$ edges from $E_{\overline{AB}}$: the edges in $E_{AB}$ can only form a matching of size $n/\alpha$ by construction. This implies that any $(\alpha/4)$-approximate solution requires recovering at least $n/\alpha$ edges from $E_{\overline{AB}}$. In the following, we prove that this is only possible if the coresets of the players are sufficiently large.

For any $i \in [k]$, define the *induced matching* $M^{(i)}$ as the unique matching in $G^{(i)}$ that is incident on *vertices of degree exactly one*, i.e., both end-points of each edge in $M^{(i)}$ have degree one in $G^{(i)}$. We emphasize that the notion of induced matching is with respect to the entire graph and not only with respect to the vertices included in the induced matching. We have the following crucial lemma on the size of $M^{(i)}$. The proof is technical and is deferred to Appendix A.

**Lemma 4.1.** *W.p. $1 - O(1/n)$, for all $i \in [k]$, $\left|M^{(i)}\right| = \Theta(n/\alpha)$.*

We are now ready to prove Theorem 3.

*Proof of Theorem 3.* Fix any randomized composable coreset (algorithm) for the matching problem that has size $o(n/\alpha^2)$. We show that such a coreset cannot achieve a better than $(\alpha/4)$-approximation over the distribution $\mathcal{D}_{\mathsf{Matching}}$. As argued earlier, to prove this, we need to show that this coreset only contains $o(n/\alpha)$ edges from $E_{\overline{AB}}$ in expectation.

Fix any player $i \in [k]$, and let $M^{\star(i)}$ be the subset of the matching $E_{\overline{AB}}$ assigned to $P^{(i)}$. It is clear that $M^{\star(i)} \subseteq M^{(i)}$ by the definition of $M^{(i)}$. Moreover, define $X_i$ as the random variable denoting the number of edges from $M^{\star(i)}$ that belong to the coreset sent by player $P^{(i)}$. Notice that $X_i$ is clearly an upper bound on the number of edges of $E_{\overline{AB}}$ that are in the final matching of coordinator and also belong to the input graph of player $P^{(i)}$. In the following, we show that

$$\mathbb{E}\left[X_i\right] = o\left(\frac{n}{k \cdot \alpha}\right) \tag{1}$$

Having proved this, we have that the expected size of the output matching by the coordinator is at most $n/\alpha + \sum_{i=1}^{k} \mathbb{E}\left[X_i\right] = n/\alpha + o(n/\alpha) < (\alpha/4) \cdot \mathsf{MM}(G)$, a contradiction.

We now prove Eq (1). In the following, we condition on the event that $\left|M^{\star(i)}\right| = \Theta(n/k)$ and $\left|M^{(i)}\right| = \Theta(n/\alpha)$; by Chernoff bound (for the first part, since $n/k = \omega(\log n)$) and Lemma 4.1 (for the second part), this event happens with probability $1 - O(1/n)$. As such, this conditioning can only change $\mathbb{E}\left[X_i\right]$ by an additive factor of $O(1)$ which we ignore in the following.

A crucial property of the distribution $\mathcal{D}_{\mathsf{Matching}}$ is that the edges in $M^{\star(i)}$ and the remaining edges in $M^{(i)}$ are indistinguishable in $G^{(i)}$. More formally, for any edge $e \in G^{(i)}$,

$$\Pr\left(e \in M^{\star(i)} \mid e \in M^{(i)}\right) = \frac{\left|M^{\star(i)}\right|}{\left|M^{(i)}\right|} = \Theta(\alpha/k)$$

On the other hand, for a fixed input $M^{(i)}$ to player $P^{(i)}$, the computed coreset $C_i$ is always the same (as the coreset is a deterministic function of the player input). Hence,

$$\mathbb{E}\left[X_i\right] = \sum_{e \in C_i} \Pr\left(e \in M_i^{\star} \mid e \in M^{(i)}\right) = |C_i| \cdot \Theta(\alpha/k) = o(n/\alpha^2) \cdot \Theta(\alpha/k) = o\left(n/(\alpha \cdot k)\right)$$

where the second last equality is by the assumption that the size of the coreset, i.e., $|C_i|$, is $o(n/\alpha^2)$. This finalizes the proof. ∎

## 4.2 A Lower Bound for Randomized Composable Coresets of Vertex Cover

In this section, we prove that the size of the corset for the vertex cover problem in Theorem 2 is indeed optimal. The following is a formal statement of Result 2 for the vertex cover problem.

**Theorem 4.** *For any $k = o(n/\log n)$ and $\alpha = o(\min\{n/k, k\})$, any $\alpha$-approximation randomized composable coreset of the minimum vertex cover problem is of size $\Omega(n/\alpha)$.*

By Yao's minimax principle [65], to prove the lower bound in Theorem 4, it suffices to analyze the performance of deterministic algorithms over a fixed (hard) distribution. We propose the following distribution for this task[5]. For simplicity of exposition, in the following, we prove a lower bound for $(c \cdot \alpha)$-approximation algorithms (for some constant $c > 0$); a straightforward scaling of the parameters proves the lower bound for $\alpha$-approximation as well.

---

[5]We point out that simpler versions of this distribution suffice for proving the lower bound in this section. However, as we would like this proof to also act as a warm-up to the proof of Theorem 6, we use the same distribution that is used to prove that theorem.

**Distribution $\mathcal{D}_{\mathsf{VC}}$.** A hard input distribution for the vertex cover problem.

- Construct $G(L, R, E)$ (with $|L| = |R| = n$) as follows:

  1. Pick $A \subseteq L$ of size $n/\alpha$ uniformly at random.
  2. Let $E_A$ be a set of edges chosen by picking each edge in $A \times R$ w.p. $k/2n$.
  3. Pick a single vertex $v^\star$ uniformly at random from $\overline{A}$ and let $e^\star$ be an edge incident on $v^\star$ chosen uniformly at random.
  4. Let $E := E_A \cup \{e^\star\}$.

- Let $E^{(1)}, \ldots, E^{(k)}$ be a *random $k$-partitioning* of $E$ and let the input to player $P^{(i)}$ be the graph $G^{(i)}(L, R, E^{(i)})$.

For any $i \in [k]$, we define $L_i^1$ as the set of vertices in $L$ with degree *exactly one* in $G^{(i)}$. We further define $R_i^1$ as the set of neighbors of vertices in $L_i^1$ (note that vertices in $R_i^1$ do not *not* necessarily have degree exactly one). We start by proving a simple property of this distribution.

**Lemma 4.2.** *For any $i \in [k]$, $\left|L_i^1\right| = \Theta(n/\alpha)$ and $\left|R_i^1\right| = \Theta(n/\alpha)$ w.p. $1 - o(1)$.*

*Proof.* Fix any player $i \in [k]$ and any vertex $v \in A$. The distribution of neighborhood of $v$ in the graph $G^{(i)}$ is as follows: pick each vertex in $R$ w.p. $1/2n$ independently; this is because each vertex in $R$ is chosen w.p. $k/2n$ to be a neighbor of $v$ in $G$ and then each of these vertices are assigned to the graph $G^{(i)}$ w.p. $1/k$ by the random $k$-partitioning. As such,

$$\Pr\left(d(v) = 1 \text{ in } G^{(i)}\right) = \binom{n}{1} \cdot \frac{1}{2n} \cdot \left(1 - \frac{1}{2n}\right)^{n-1} \approx \frac{1}{2\sqrt{e}} = \Theta(1)$$

Consequently, we have $\mathbb{E}\left[\left|L_i^1\right|\right] = |A| \cdot \Theta(1) = \Theta(n/\alpha)$ and by Chernoff bound, $\left|L_i^1\right| = \Theta(n/\alpha)$ (note that for one player $v^\star$ would also belong to $O_i$ but that only changes the size of $|O_i|$ by one vertex).

We now bound the size of $R_i^1$. Each vertex in $L_i^1$ is choosing one vertex uniformly at random from $R$ and hence we can model this distribution by a simple balls and bins experiment (throwing $\left|L_i^1\right|$ balls into $n$ bins, each independently and uniformly at random), and hence by a standard fact about balls and bins experiments argue that $\left|R_i^1\right| = \Theta(n/\alpha)$ w.p. $1 - o(1)$ as well (see Proposition A.1 in Appendix A for a proof of this fact about balls and bins experiments). ∎

We can now prove Theorem 4.

*Proof of Theorem 4.* Let $i$ be the index of the player $P^{(i)}$ that the edge $e^\star$ is given to. We argue that if the coreset sent by player $P^{(i)}$ is of size $o(n/\alpha)$, then the coordinator cannot obtain a vertex cover of size $o(n)$. As the graph $G$ admits a vertex cover of size $(n/\alpha + 1)$ (pick $A$ and $v^\star$), this proves the theorem.

By Lemma 4.2, the set of vertices in $L$ with degree exactly one in $G^{(i)}$ and the set of their neighbors in $R$, i.e., the sets $L_i^1$ and $R_i^1$, are of size $\Theta(n/\alpha)$ w.p. $1 - o(1)$. In the following, we condition on this event. As the algorithm used by $P^{(i)}$ to create the coreset is deterministic, given a fixed input, it always creates the same coreset. However, a crucial property of the distribution $\mathcal{D}_{\mathsf{VC}}$ is that, conditioned on a fixed assignment to $L_i^1$, the vertex $v^\star$ is chosen uniformly at random from $L_i^1$. This implies that if the coreset of player $P^{(i)}$ contains $o(n/\alpha)$ edges, then w.p. $1 - o(1)$, $e^\star$ is not part of the coreset ($e^\star$ is chosen uniformly at random from the set of all edges incident on $L_i^1$). Similarly, if the coreset fixes $o(n/\alpha)$ vertices to be added to the final solution, w.p. $1 - o(1)$, no end

15

point of $e^\star$ is added to this fixed set ($v^\star$ is chosen uniformly at random from $L_i^1$ of size $\Theta(n/\alpha)$, and the other end point of $e^\star$ is chosen uniformly at random from $R_i^1$ of size $\Theta(n/\alpha)$). Finally, the coresets of other players are all independent of the edge $e^\star$ and hence as long as the total number of fixed vertices sent by the players is $o(n)$, w.p. $1 - o(1)$, no end points of $e^\star$ are present in the fixed solution. Conditioned on these three events, w.p. $1 - o(1)$, the output of the algorithm does not cover the edge $e^\star$ and hence is not a feasible vertex cover.

We remark that this argument holds even if we are allowed to add extra vertices to the final vertex cover (other than the ones fixed by the players or computed as a vertex cover of the edges in the coresets), since conditioned on $e^\star$ not being present in any coreset, the end point of this edge are chosen uniformly at random from all vertices in $L \setminus A$ and $R$ and hence a solution of size $o(n)$ would not contain either of them w.p. $1 - o(1)$. ∎

## 5    Communication Complexity Lower Bounds

We prove Result 3 in this section, showing that our randomized composable coresets in fact obtain the optimal communication complexity (among all possible protocols, not necessarily a coreset) in the simultaneous communication model. This result is a vast generalization of Result 2 proved in Section 4.

### 5.1    An $\Omega(nk/\alpha^2)$ Lower Bound on Communication Complexity of Matching

We prove a lower bound on the simultaneous communication complexity of the matching problem in the random partition model, formalizing Result 3 for matching.

**Theorem 5.** *For any $\alpha$ between $\Omega(\log n)$ and $o(\min\left\{\sqrt{n/k}, k\right\})$, the simultaneous communication complexity of $\alpha$-approximating the matching problem in the random partition model is $\Omega(nk/\alpha^2)$.*

By Yao's minimax principle [65], it suffices to analyze the communication complexity of *deterministic* protocols over a fixed (hard) distribution. We again use the distribution $\mathcal{D}_{\mathsf{Matching}}$ in Section 4.1. For the convenience of the reader, we repeat the description of this distribution here.

---

**Distribution $\mathcal{D}_{\mathsf{Matching}}$.** A hard input distribution for the matching problem.

- Let $G(L, R, E)$ (with $|L| = |R| = n$) be constructed as follows:

    1. Pick $A \subseteq L$ and $B \subseteq R$, each of size $n/\alpha$, uniformly at random.
    2. Define $E_{AB}$ as a set of edges between $A$ and $B$, chosen by picking each edge in $A \times B$ w.p. $k \cdot \alpha/n$.
    3. Define $E_{\overline{AB}}$ as a *random* perfect matching between $\overline{A}$ and $\overline{B}$.
    4. Let $E := E_{AB} \cup E_{\overline{AB}}$.

- Let $E^{(1)}, \ldots, E^{(k)}$ be a *random k-partitioning* of $E$ and let the input to player $P^{(i)}$ be the graph $G^{(i)}(L, R, E^{(i)})$.

---

Let $G$ be a graph sampled from the distribution $\mathcal{D}_{\mathsf{Matching}}$. Notice first that the graph $G$ always has a matching of size at least $n - n/\alpha \geq n/2$, i.e., the matching $E_{\overline{AB}}$. Additionally, it is easy to see that any matching of size more than $2n/\alpha$ in $G$ uses at least $n/\alpha$ edges from $E_{\overline{AB}}$: the edges in $E_{AB}$ can only form a matching of size $n/\alpha$ by construction. This implies that any $(\alpha/4)$-approximate

solution requires recovering at least $n/\alpha$ edges from $E_{\overline{AB}}$[6]. It is this task that we show requires $\Omega(nk/\alpha^2)$ communication.

The following definitions are identical to those in Section 4.1. For any $i \in [k]$, define the *induced matching* $M^{(i)}$ as the unique matching in $G^{(i)}$ that is incident on *vertices of degree exactly one*, i.e., both end-points of each edge in $M^{(i)}$ have degree one in $G^{(i)}$. Recall that by Lemma 4.1, $\left|M^{(i)}\right| = \Theta(n/\alpha)$ w.h.p.

Let $M^{\star(i)}$ be the subset of the matching $E_{\overline{AB}}$ assigned to $P^{(i)}$. It is clear that $M^{\star(i)} \subseteq M^{(i)}$ by the definition of $M^{(i)}$. By our previous discussion, it is these edges of $M^{\star(i)}$ that the player $P^{(i)}$ needs to communicate to the coordinator. Moreover, notice that the players can simply ignore all edges in $G^{(i)}$ that do not belong to $M^{(i)}$ as they clearly cannot be in $M^{\star(i)}$. However, a crucial property of the distribution $\mathcal{D}_{\mathsf{Matching}}$ is that the edges in $M^{\star(i)}$ and the remaining edges in $M^{(i)}$ are *indistinguishable* in $G^{(i)}$. In other words, conditioned on a specific assignment for $M^{(i)}$, *any* edge $e \in M^{(i)}$ belongs to the matching $M^{\star(i)}$ w.p. $\alpha/k$ independent of the other edges. Moreover, it is intuitive to think that only player $P^{(i)}$ is able to communicate the edges in $M^{\star(i)}$ as the input of other players, while dependent on the set of vertices in the matching $M^{(i)}$, are essentially independent of the edges in $M^{(i)}$. This discussion suggests the following intermediate problem in the one-way two-player communication model.

**Problem 1** (MatchingRecovery problem). *Let $H$ be a bipartite graph with $t$ vertices on each side. Alice is given a perfect matching $M_{\mathsf{Alice}}$ in $H$ and Bob is given the following input:*

- *Two sets $P$, $Q$ (each in one side of the bipartition of $H$) with $|P| = |Q| = p$ and the* promise *that in matching $M_{\mathsf{Alice}}$ vertices in $P$ are matched to vertices $Q$.*

- *A set $E_{\mathsf{Bob}}$ of edges in $H$ with the* promise *that the matching $M_{\mathsf{Alice}}$ does not contain any edge from $E_{\mathsf{Bob}}$.*

*The goal is for Alice to send a message to Bob and Bob needs to output the edges in the matching $M_{\mathsf{Alice}}$ that are between $P$ and $Q$.*

Consider the following distribution $\mathcal{D}_{\mathsf{MR}}$ for MatchingRecovery based on the distribution $\mathcal{D}_{\mathsf{Matching}}$: Fix any arbitrary $i \in [k]$; we sample an input instance $G^{(1)}, \ldots, G^{(k)}$ from $\mathcal{D}_{\mathsf{Matching}}$. Then, we let $H$ be the bipartite graph on the set of vertices in $M^{(i)}$ and let $M_{\mathsf{Alice}} = M^{(i)}$. We let the input sets $P$ and $Q$ to Bob be the set of vertices incident on $M^{\star(i)}$. Finally, we let $E_{\mathsf{Bob}}$ be the set of edges assigned to all graphs $G^{(j)}$ for $j \neq i$ that are between the vertices matched by $M^{(i)}$, i.e., are inside the graph $H$. This completes the description of the distribution $\mathcal{D}_{\mathsf{MR}}$.

In the following, we condition on the inputs chosen from $\mathcal{D}_{\mathsf{MR}}$ to have the following additional property: $\left|M^{(i)}\right| = \Theta(n/\alpha)$ and $\left|M^{\star(i)}\right| = \Theta(n/k)$. Notice that by Lemma 4.1, w.p. $1 - O(1/n)$, for any player $P^{(i)}$, $\left|M^{(i)}\right| = \Theta(n/\alpha)$. A simple application of Chernoff bound also ensures that the number of edges from $E_{\overline{AB}}$ assigned to each player, i.e., the edges in the matching $M^{\star(i)}$ is $\Theta(n/k)$ w.p. $1 - O(1/n)$. Consequently, conditioning on this event is essentially not changing the distribution and hence for simplicity from now on, we always assume the inputs chosen from $\mathcal{D}_{\mathsf{MR}}$ satisfy the mentioned properties. We establish the following lower bound for MatchingRecovery.

**Lemma 5.1** (Communication Complexity of MatchingRecovery). *Suppose $s = \Omega(k)$ denotes the communication cost of a protocol for MatchingRecovery and $X$ denotes the number of edges output by Bob in this protocol; for the inputs chosen from the distribution $\mathcal{D}_{\mathsf{MR}}$, we have $\mathbb{E}[X] \leq (\alpha/k) \cdot O(s)$.*

---

[6]Similar to Section 4.1, we prove the lower bound for $(\alpha/4)$-approximation protocols for simplicity of representation.

17

Notice that the distribution $\mathcal{D}_{\mathsf{MR}}$ for MatchingRecovery imposes a non-trivial correlation between the inputs of the two players which complicates the proof of this lower bound. We address this issue by expressing this distribution as a convex combination of a relatively small number of simpler (yet non-trivial) distributions and prove the lower bound for each distribution separately. These distributions are still not product distributions but we can show that the mild correlation in the input of the players in this case can be managed directly using a careful combinatorial argument. The proof is deferred to Section 5.2. Before that, we prove a formal reduction from the matching problem to MatchingRecovery and use Lemma 5.1 to finalize the proof of Theorem 5.

*Proof of Theorem 5.* Fix a protocol $\Pi_{\mathsf{Matching}}$ for the matching problem on $\mathcal{D}_{\mathsf{Matching}}$ with communication cost $o(nk/\alpha^2)$ and suppose each player $P^{(i)}$ communicates at most $s_i$ bits in this protocol. We assume that $s_i = \Omega(k)$ as otherwise we can simply augment it with $\Omega(k)$ bits to satisfy this bound while increasing the total communication cost of the protocol by $O(k^2) = o(nk/\alpha^2)$ bits (since $\alpha = o(\sqrt{n/k})$). Our goal is to show that in this protocol, at most $o(n/\alpha)$ edges from $E_{\overline{AB}}$ can be matched in expectation. The result then follows from the fact that obtaining better than $(\alpha/4)$-approximation requires outputting $\Omega(n/\alpha)$ edges from $E_{\overline{AB}}$.

We use $\Pi_{\mathsf{Matching}}$ to create $k$ protocols for MatchingRecovery, whereby in the $i$-th protocol $\Pi_i$, Alice plays the role of player $P^{(i)}$ and Bob plays the role of all other players plus the coordinator. Fix an $i \in [k]$; the protocol $\Pi_i$ works as follows.

Given their input in MatchingRecovery, Alice and Bob sample two random sets $X_M \subseteq L$ and $Y_M \subseteq R$, each of size $t$ using public randomness. They also sample two random sets $X_{\overline{M}} \subseteq L \setminus X_M$ and $Y_{\overline{M}} \subseteq R \setminus Y_M$, each of size $(n/\alpha - t + p)$.

Alice creates the graph $G^{(i)}$ by letting $M^{(i)}$ be the matching $M_{\mathsf{Alice}}$ and choose the remainder of the graph $G^{(i)}$ by sampling the edges between $X_{\overline{M}}$ and $Y_{\overline{M}}$ using the same distribution as the distribution of $G^{(i)}$ conditioned on $M^{(i)} = M_{\mathsf{Alice}}$ and the set of non-zero degree vertices in $G^{(i)}$ being subset of $X_{\overline{M}}$ and $Y_{\overline{M}}$.

Bob creates the input of the other players as follows. Bob picks a random mapping $\sigma : E_{\mathsf{Bob}} \to [k] \setminus \{i\}$ and in the graph $G^{(j)}$ for $j \neq i$, he assigns $\sigma^{-1}(j) \subseteq E_{\mathsf{Bob}}$ to be the edges between $X_M$ and $Y_M$. Finally, Bob samples the remainder of the graphs from the joint distribution of $\mathcal{D}_{\mathsf{Matching}}$ conditioned on the set $A = X_{\overline{M}} \cup X_M \setminus P$, $B = X_{\overline{M}} \cup X_M \setminus Q$, and the edges between $X_M$ and $Y_M$ be the ones already sampled (via $\sigma$). Note that since Bob has the knowledge of the sets $P$ and $Q$, he can sample the reminder of the matching $E_{\overline{AB}}$ for the $k-1$ remaining players as well.

One can verify that the distribution of the instances created by this reduction matches the distribution $\mathcal{D}_{\mathsf{Matching}}$. To finish the reduction, Alice and Bob simulate the protocol $\Pi_{\mathsf{Matching}}$ by Alice sending the message of $P^{(i)}$ to Bob (or equivalently the coordinator) and Bob creating the message of all other players locally and completing the protocol. Bob then outputs the part of the matching computed by $\Pi_{\mathsf{Matching}}$ which lies between $P$ and $Q$. This results in a protocol $\Pi_i$ for MatchingRecovery with communication cost of $s_i$ bits.

We can now invoke Lemma 5.1 to argue that the expected number of edges matched by $\Pi_i$ and hence by $\Pi_{\mathsf{Matching}}$ for the player $P^{(i)}$ in the distribution $\mathcal{D}_{\mathsf{Matching}}$ is at most $(\alpha/k) \cdot O(s_i)$. Summing over all players, we have that the total number of matched edges in $E_{\overline{AB}} = \bigcup_{i=1}^{k} M^{\star(i)}$ is $\sum_{i=1}^{k} (\alpha/k) \cdot O(s_i) = (\alpha/k) \cdot o(nk/\alpha^2) = o(n/\alpha)$, where the first equality is by the bound on the communication cost of the protocol. This completes the proof. ∎

We remark that the bound in Theorem 5 is tight (up to an $O(\log n)$ factor) for all ranges of $\alpha$.

**Remark 5.2.** *The protocol in which each player computes a maximum matching of the input graph, subsamples the edges of this matching w.p. $1/\alpha$, and sends it to the coordinator who outputs a*

*maximum matching of the received matchings is an $\alpha$-approximation protocol for the maximum matching problem with $\widetilde{O}(nk/\alpha^2)$ total communication.*

We briefly sketch the proof of correctness for protocol in Remark 5.2. Assume first that the maximum matching size in each player input is of size $\widetilde{O}(n/\alpha)$. The bound on the total communication cost follows immediately from this assumption. To see the correctness, recall that in the proof of Theorem 1, we showed each coreset (here, each player) can increase the size of the output matching by $\Omega(\mathsf{MM}(G)/k)$; since we are subsampling the maximum matching by a factor of $\alpha$, this increment would be $\Omega(\mathsf{MM}(G)/\alpha k)$ and hence over all $k$ players, we obtain a matching of size $\Omega(\mathsf{MM}(G)/\alpha)$. The assumption on the size of the maximum matching in each player input is essentially without loss of generality since otherwise one player can send an $\alpha$-approximate matching to the coordinator alone, resulting in a protocol with $\widetilde{O}(n/\alpha)$ communication. We point out that a simple concentration result proves that the maximum matching size between players is concentrated within an $O(\log n)$ factor. This easily implies that in this case we can ensure that only one player is sending his maximum matching and not all players.

## 5.2 Communication Complexity of MatchingRecovery

We start by reformulating the distribution $\mathcal{D}_{\mathsf{MR}}$ to make it more suitable for proving the lower bound in Lemma 5.1. Indeed, the distribution $\mathcal{D}_{\mathsf{MR}}$ is *not* a product distribution: the promise that Alice's matching $M_{\mathsf{Alice}}$ needs to always match the set $P$ to $Q$ correlates Alice's and Bob's input in a non-trivial way, complicating the analysis. To address this, we show that the distribution $\mathcal{D}_{\mathsf{MR}}$ can be expressed as a convex combination of a relatively small set of (essentially) product distributions; this significantly simplifies the proof of the lower bound.

Let us first define the distribution $\mathcal{D}_{\mathsf{MR}}$ directly, i.e., without depending on the distribution $\mathcal{D}_{\mathsf{Matching}}$. In $\mathcal{D}_{\mathsf{MR}}$ conditioned on $|M_{\mathsf{Alice}}| = t$ and $|P| = |Q| = p$, the matching $M_{\mathsf{Alice}}$ is chosen uniformly at random from the set of all matchings of size $t$ in $H$ and then $P$ and $Q$ are chosen uniformly at random from all pairs of sets of size $p$ that are matched together in $M_{\mathsf{Alice}}$. Finally, the edge-set $E_{\mathsf{Bob}}$ is chosen by picking each edge in $H$ that is not incident on $P$ and $Q$ and not in $M_{\mathsf{Alice}}$ w.p. $(k-1) \cdot \alpha/n$. One can check that this results in an equivalent definition of $\mathcal{D}_{\mathsf{MR}}$.

We now reformulate the distribution $\mathcal{D}_{\mathsf{MR}}$ as follows: we first randomly partition the vertices in $L_H$ and $R_H$ (i.e., the bipartition of the graph $H$) into $c := \lfloor t/p \rfloor$ *blocks* denoted by $\mathbf{B} := (P_1, Q_1), \ldots, (P_c, Q_c)$ such that $P_i \subseteq L_H$, $Q_i \subseteq R_H$ and $|P_i| = |Q_i| = p$ for all $i \in [c]$. We then create $M_{\mathsf{Alice}}$ by picking a random matching that matches each $P_i$ to $Q_i$[7]. Finally, the input to Bob is a pair $P$ and $Q$ chosen uniformly at random from these $c$ blocks. We pick the edge-set $E_{\mathsf{Bob}}$ of Bob as before. It is again easy to verify that this is indeed an equivalent formulation of the distribution $\mathcal{D}_{\mathsf{MR}}$.

Suppose we provide the identity of the blocks $\mathbf{B}$ to both Alice and Bob; this essentially breaks the dependence between Alice's and Bob's inputs (except for the mild correlation enforced by $E_B$ that we deal with directly). Note that revealing this extra information can only make our lower bound result stronger. In the following, we argue that even with $\mathbf{B}$ being public information, to solve the problem, Alice needs to communicate a large message.

For any fixed $\mathbf{B}$, let $\sigma_{\mathbf{B}}(M_{\mathsf{Alice}})$ be the (deterministic) mapping used by Alice to create the message sent to Bob. The mapping $\sigma_{\mathbf{B}}$ then partitions all matchings $M_{\mathsf{Alice}}$ that are valid with respect to $\mathbf{B}$ into $2^s$ classes $\Sigma_1, \ldots, \Sigma_{2^s}$ (one per each message). Moreover, for $\mathbf{B}$ and a fixed set of edges $E_{\mathsf{Bob}}$, we define $\mathcal{F}(E_{\mathsf{Bob}}, \mathbf{B})$ as the set of matchings $M_{\mathsf{Alice}}$ that conform to the restrictions

---

[7]Note that this way, it is possible that up to $p$ vertices in $L_H$ and $R_H$ become "left overs" i.e., do not belong to any block. We pick a random matching between these vertices also to complete the description of $M_{\mathsf{Alice}}$.

imposed by both $E_{\mathsf{Bob}}$ and $\mathbf{B}$ (we use $\mathcal{F}$ if $E_{\mathsf{Bob}}$ and $\mathbf{B}$ are clear from the context). Note that $\sigma_{\mathbf{B}}$ similarly partitions $\mathcal{F}$ into $2^s$ different classes as well.

An important observation is that given a message corresponding to some class $\Sigma_i$ and the edges $E_{\mathsf{Bob}}$, the matching $M_{\mathsf{Alice}}$ is chosen uniformly at random from all matchings in $\Sigma_i \cap \mathcal{F}$; consequently, Bob can only output an edge in the final answer if it belongs to *all* matchings in $\mathcal{F}$ that are mapped to $\Sigma_i$. For any set $F \subseteq \mathcal{F}$, we define $M_F$ as the intersection of all matchings $M_{\mathsf{Alice}}$ in $F$. Intuitively, whenever $M_F$ is large, the set $F$ itself should be small as many edges of the matchings $M_{\mathsf{Alice}} \in F$ are forced to be the same. We formalize this intuition in the following lemma.

**Lemma 5.3.** *For any set $F \subseteq \mathcal{F}(E_{\mathsf{Bob}}, \mathbf{B})$, if $M_F$ contains $\ell$ edges, then $|F| \leq 2^{-(\ell - \Theta(k))} \cdot |\mathcal{F}|$ w.h.p. (the probability is taken over the choice of $E_B$ after fixing $\mathbf{B}$).*

*Proof.* Note that we can switch the order in which we pick $M_{\mathsf{Alice}}$ and $E_{\mathsf{Bob}}$ in the distribution $\mathcal{D}_{\mathsf{MR}}$. Fix any block $(P_i, Q_i) \in \mathbf{B}$ and any vertex $u \in P_i$; each edge $(u, v)$ for $v \in Q_i$ is chosen w.p. $(k-1) \cdot \alpha/n$ in $E_{\mathsf{Bob}}$. Hence, w.h.p, at most $\beta := 2 \cdot |Q_i| \cdot (k-1) \cdot \alpha/n = \Theta(\alpha)$ neighbors are chosen for $u$ in $Q_i$ in $E_{\mathsf{Bob}}$ (here, we used a standard application of Chernoff bound and the assumption that $\alpha = \Omega(\log n)$). This means that for each vertex in $P_i$ there are $p - \beta$ possible choices for its neighbor in $M_{\mathsf{Alice}}$; hence, there are at least $(p - \beta)!$ choices for $M_{\mathsf{Alice}}$ to match $P_i$ to $Q_i$. Since there are $c$ different blocks, we have $|\mathcal{F}| \geq ((p - \beta)!)^c$.

Now suppose we fix $\ell$ edges for the matching $M_{\mathsf{Alice}}$ (as happens in the set $F$), and let $\ell_i$ be the number of fixed edges between $(P_i, Q_i) \in \mathbf{B}$. There can be at most $(p - \ell_i)!$ choices for the matching between $P_i$ and $Q_i$ (we ignore the restriction implied by $E_{\mathsf{Bob}}$ for the purpose of obtaining an upper bound). Hence,

$$\frac{|\mathcal{F}|}{|F|} \geq \prod_{i=1}^{c} \frac{(p - \beta)!}{(p - \ell_i)!} \geq \prod_{i=1}^{c} (p - \beta) \dots (p - \beta - \ell_i + 1)$$
$$\geq \prod_{i=1}^{c} 2^{\ell_i - \beta} = 2^{\ell - c \cdot \beta} = 2^{\ell - \Theta(k)}$$

where the last equality is by the fact that $\beta = \Theta(\alpha)$ and $c = \Theta(k/\alpha)$. ∎

We are now ready to finalize the proof of Lemma 5.1.

*Proof of Lemma 5.1.* Fix a set of blocks $\mathbf{B}$ and edges $E_{\mathsf{Bob}}$ and assume that the event in Lemma 5.3 happens. Notice that the mapping $\sigma_{\mathbf{B}}$ maps the set $\mathcal{F}(E_{\mathsf{Bob}}, \mathbf{B})$ to $2^s$ different choices $\Sigma_1, \dots, \Sigma_{2^s}$. We define $\mathcal{E}$ as the event that $\sigma_{\mathbf{B}}(M_{\mathsf{Alice}})$ maps to a class $\Sigma$ with $|\Sigma \cap \mathcal{F}| \geq |\mathcal{F}|/2^{2s}$. The following claim on the probability of $\mathcal{E}$ can be proven using a simple counting argument; the proof is deferred to after the current proof.

**Claim 5.4.** $\Pr(\mathcal{E}) = 1 - O(1/n)$.

Now fix a set $\Sigma$ that corresponds to the message Alice sent to Bob and suppose $\mathcal{E}$ happens. As argued earlier, given a message Bob can only output an edge in $M_{\mathsf{Alice}}$ if it belongs to all matchings that are mapped to $\Sigma$, i.e., to $M_\Sigma$. By Lemma 5.3, the matching $M_\Sigma$ contains at most $\ell = 2s + \Theta(k) = \Theta(s)$ edges (since $s = \Omega(k)$).

Now recall that the block $(P, Q)$ of Bob is chosen *uniformly at random* from the blocks in $\mathbf{B}$, *even* conditioned on a specific input matching $M_{\mathsf{Alice}}$ to Alice; this implies that in expectation $\ell/c = O(\alpha/k) \cdot \ell = (\alpha/k) \cdot O(s)$ edges of $M_\Sigma$ are between $(P, Q)$. Consequently, Bob can only output $(\alpha/k) \cdot O(s)$ edges between $P$ and $Q$ in expectation. To finalize the proof, note that $\mathcal{E}$ and the event in Lemma 5.3 happens w.h.p., and hence conditioning on these events can only change the expectation by an $O(1)$ additive factor. ∎

For completeness, we provide a proof of Claim 5.4 here.

*Proof of Claim 5.4.* We say that $\Sigma_i$ in the partition $\Sigma_1, \ldots, \Sigma_{2^s}$ is *light* iff $|\Sigma_i \cap \mathcal{F}| < |\mathcal{F}| / 2^{2s}$. Since the matching $M_A$ is chosen uniformly at random, the probability that $\sigma_{\mathbf{B}}(M_A)$ maps to some $\Sigma_i$ is exactly $|\Sigma_i \cap \mathcal{F}| / |\mathcal{F}|$. Hence, the probability that $M_A$ maps to some light set is at most

$$2^s \cdot \frac{|\Sigma_i \cap \mathcal{F}|}{|\mathcal{F}|} \leq 2^s \cdot \frac{|\mathcal{F}|}{(2^{2s} \cdot |\mathcal{F}|)} = \frac{1}{2^s} = O(1/n)$$

where we used the fact that $s = \Omega(k) = \Omega(\alpha) = \Omega(\log n)$. ∎

## 5.3 An $\Omega(nk/\alpha)$ Lower Bound on Communication Complexity of Vertex Cover

We prove the following theorem on the simultaneous communication complexity of vertex cover, formalizing Result 3 for vertex cover.

**Theorem 6.** *For any $\alpha$ between $\Omega(\log n)$ and $o(\min\{n/k, k\})$, the simultaneous communication complexity of $\alpha$-approximating the vertex cover problem in the random partition model with success probability at least $0.9$ is $\Omega(nk/\alpha)$.*

For simplicity of exposition, we prove the lower bound for protocols that can obtain a $c \cdot \alpha$ approximation for some small constant $c > 0$ to be determined later. By re-parametrizing $\alpha$ by a constant factor in the following, one can obtain the result for $\alpha$-approximation protocols as well. We again use the distribution $\mathcal{D}_{\mathsf{VC}}$ in Section 4.2. For the convenience of the reader, we repeat the description of this distribution here.

---

**Distribution $\mathcal{D}_{\mathsf{VC}}$.** A hard input distribution for the vertex cover problem.

- Construct $G(L, R, E)$ (with $|L| = |R| = n$) as follows:

   1. Pick $A \subseteq L$ of size $n/\alpha$ uniformly at random.
   2. Let $E_A$ be a set of edges chosen by picking each edge in $A \times R$ w.p. $k/2n$.
   3. Pick a single vertex $v^\star$ uniformly at random from $\overline{A}$ and let $e^\star$ be an edge incident on $v^\star$ chosen uniformly at random.
   4. Let $E := E_A \cup \{e^\star\}$.

- Let $E^{(1)}, \ldots, E^{(k)}$ be a *random $k$-partitioning* of $E$ and let the input to player $P^{(i)}$ be the graph $G^{(i)}(L, R, E^{(i)})$.

---

The intuition behind the proof is as follows. The distribution ensures that w.h.p., the input to each player $P^{(i)}$ contains $\Theta(n/\alpha)$ vertices in $L$ with degree exactly one. Let us denote this set with $D_i$. Now consider the input of the player $P^{(i^\star)}$ which is given the edge $e^\star$ also. It is easy to see that in this case, player $P^{(i^\star)}$ is oblivious to which vertex in $D_i$ is $u^\star$; more formally, conditioned on the input $D_{i^\star}$, the vertex $u^\star$ is chosen uniformly at random from $D_{i^\star}$. This means the if $P^{(i^\star)}$ communicates $o(|D_i|) = o(n/\alpha)$ bits, he is essentially not "revealing any information" about $v^\star$ (or the other end point of $e^\star$). On the other hand, as only $P^{(i^\star)}$ has a knowledge about $v^\star$, this intuitively means that coordinator is not provided with enough information about $v^\star$ as well. This forces the coordinator to cover $\Omega(n)$ vertices to ensure that $e^\star$ is also being covered.

Making this intuition formal is complicated by the fact that the message by other players is still revealing "some information" about the input of player $P^{(i^\star)}$, for instance, the identity of the set $A$

or the set of edges that may possibly be in the input of $P^{(i^\star)}$. To overcome this, we show that by proper conditioning on some part of the input, one can embed an instance of the well-known set disjointness problem in this distribution. On a high level, solving the disjointness on this embedded instance amounts to finding the vertex $u^\star$. This further allows us to design a reduction from our problem to the disjointness problem and prove the lower bound. Interestingly, while we consider the vertex cover only in the simultaneous model, our reduction requires a two-way communication between the players (however note that disjointness is still hard even in the two-way communication model). We now continue with the formal proof.

We can interpret the last line of distribution $\mathcal{D}_{\mathsf{VC}}$ as follows: Pick a random $k$-partitioning $\widehat{E}^{(1)}, \ldots, \widehat{E}^{(k)}$ of *all possible edges between $L$ and $R$*, and let $E^{(i)} = \widehat{E}^{(i)} \cap E$ for all $i \in [k]$. In the following, we assume that this *initial partitioning* $\widehat{E}^{(1)}, \ldots, \widehat{E}^{(k)}$ is public knowledge, as this allows us to reduce the dependence between the inputs of players which is crucial in our lower bound proof. Clearly, this assumption can only strengthen our results. Throughout the proof, we fix an arbitrary small constant $\varepsilon > 0$. We say that the initial partitioning $\widehat{E}^{(1)}, \ldots, \widehat{E}^{(k)}$ is $\varepsilon$-*balanced* if the degree of each vertex $v \in L$ in each graph $G(L, R, \widehat{E}^{(i)})$ for all $i \in [k]$ is in $(1 \pm \varepsilon) \cdot n/k$ (note that this graph is *not* the input graph to player $P^{(i)}$). As $n/k \geq \alpha = \Omega(\log n)$, by Chernoff bound, any initial partitioning is $\varepsilon$-balanced w.p. $1 - O(1/n)$. Consequently, conditioning on this event is essentially not changing the distribution and hence for simplicity from now on, we always assume the inputs chosen from $\mathcal{D}_{\mathsf{VC}}$ satisfy the $\varepsilon$-balanced property.

We say that the player $P^{(i)}$ (for $i \in [k]$) is the *critical* player iff the edge $e^\star$ is assigned to $E^{(i)}$, i.e., it appears in the input to $P^{(i)}$. We use $i^\star$ to denote the index of the critical player. In the following, we show that the critical player and the coordinator need to (implicitly) solve a "hard" communication task (named the *Hidden Vertex Problem*, $\mathsf{HVP}$ for short) which requires a large communication from $P^{(i)}$.

Fix $\Pi_{\mathsf{VC}}$ as a $\delta$-error $(c \cdot \alpha)$-approximation protocol for vertex cover over distribution $\mathcal{D}_{\mathsf{VC}}$ (for some sufficiently small constant $c > 0$ to be determined later). For any $i \in [k]$, let $\delta_i$ be the probability that $\Pi_{\mathsf{VC}}$ errs conditioned on $i^\star = i$. A simple averaging argument ensures that for any $i \in [k]$ there exists a set $\widehat{E}^{(i)}$ such that $\Pr\left( \Pi_{\mathsf{VC}} \text{ errs} \mid \widehat{E}^{(i)}, i^\star = i \right) \leq \delta_i$. We refer to such $\widehat{E}^{(i)}$ as a *good* initial partition for $P^{(i)}$.

Fix a player $i \in [k]$ and a good initial partition $\widehat{E}^{(i)}$ for $i$. Let $D_0^{(i)}, D_1^{(i)}$ and $D_{\geq 2}^{(i)}$ be the set of vertices in $A$ that have degree, respectively, zero, one, and *at least* two in the graph $G^{(i)}$. We further define $D_{\leq 1}^{(i)} := D_0^{(i)} \cup D_1^{(i)} = A \setminus D_{\geq 2}^{(i)}$.

**Claim 5.5.** *For any $i \in [k]$ and any good initial partition $\widehat{E}^{(i)}$ for $i$, there exists a set $D_{\geq 2}^{(i)}$ with $\left| A \setminus D_{\geq 2}^{(i)} \right| = \Omega(n/\alpha)$ such that, $\Pr\left( \Pi_{\mathsf{VC}} \text{ errs} \mid \widehat{E}^{(i)}, D_{\geq 2}^{(i)}, i^\star = i \right) \leq \delta_i + o(1)$.*

*Proof.* Each vertex in $v \in A$ has degree more than 1 in $G^{(i)}$ independently and with some constant probability $p$ bounded away from 1 (see the exact calculation in Claim 5.6). Hence, in expectation $p \cdot n/\alpha$ vertices in $A$ have degree more than 1. As $n/\alpha = \Omega(\log n)$, by Chernoff bound plus a union bound, w.p. $1 - o(1)$ at most $p \cdot n/\alpha + o(n/\alpha)$ vertices in $A$ have degree more than 1 in $G^{(i)}$. The claim now follows immediately from this. ∎

In the following, we further condition on a set $D_{\geq 2}^{(i)}$ as in Claim 5.5. This implies that $A \setminus D_{\geq 2}^{(i)}$, i.e., the set $D_{\leq 1}^{(i)}$ is a set of size $\Omega(n/\alpha)$ chosen uniformly at random from $L \setminus D_{\geq 2}^{(i)}$. We are now ready to define the hidden vertex problem in the one-way two-player communication model.

**Problem 2** (The Hidden Vertex Problem ($\mathsf{HVP}$))**.** *There are two disjoint sets $U$ and $V$ and a mapping $\sigma : U \to V$ known to both Alice and Bob. Bob is given a set $T \subseteq U$. Alice is given a*

*set $S \subseteq T$ and a single vertex $u^\star$ chosen from $U \setminus T$ (identity of $u^\star$ is unknown to players). Alice sends a single message to Bob and Bob needs to output two sets $X \subseteq U$ and $Y \subseteq V$ such that either $u^\star \in X$ or $\sigma(u^\star) \in Y$. The goal of the players is to* minimize the size *of $X \cup Y$.*

Consider the following distribution $\mathcal{D}_{\mathsf{HVP}}$ for $\mathsf{HVP}$: we sample an instance of vertex cover from $\mathcal{D}_{\mathsf{VC}} \mid \widehat{E}^{(i)}, D_{\geq 2}^{(i)}, i^\star = i$, where $\widehat{E}^{(i)}$ is a good initial partition and $D_{\geq 2}^{(i)}$ is a fixed set defined in Claim 5.5. We set $U = L \setminus D_{\geq 2}^{(i)}$, $V = R$, and choose $\sigma : U \to V$ by mapping each $u \in U$ to one of the neighbors of $u$ (in $L$) in $\widehat{E}^{(i)}$ uniformly at random. Next, we set $T = D_{\leq 1}^{(i)}$ and $S = D_1^{(i)} \cup \{v^\star\}$; this way the vertex $u^\star$ in $\mathsf{HVP}$ is the special vertex $v^\star$ in distribution $\mathcal{D}_{\mathsf{VC}}$. We make the following simple observation about distribution $\mathcal{D}_{\mathsf{HVP}}$.

**Claim 5.6.** *Each vertex in $T$ independently belongs to $S$ w.p. $(1 \pm O(\varepsilon)) \cdot 1/3$.*

*Proof.* Note that $T = D_{\leq 1}^{(i)}$, i.e., the vertices that have degree 0 or 1 in $G^{(i)}$, and vertices in $S \cap T$ are vertices that have degree exactly 1 in $G^{(i)}$. Fix a vertex $v \in A$ and consider only conditioning on $\widehat{E}^{(i)}$. We know that since $\widehat{E}^{(i)}$ is a good initial partition, $v$ is incident on $(1 \pm \varepsilon) \cdot n/k$ edges in $\widehat{E}^{(i)}$, and each of these edges appear independently in $G^{(i)}$ w.p. $k/2n$ (by definition of $\mathcal{D}_{\mathsf{VC}}$). Let $d(v)$ denote the degree of $v$ in $G^{(i)}$.

$$\Pr\left(d(v) = 0 \mid \widehat{E}^{(i)}\right) = (1 - \frac{k}{2n})^{(1\pm\varepsilon)\cdot\frac{n}{k}} = e^{-\frac{1}{2}} \cdot (1 \pm O(\varepsilon))$$

$$\Pr\left(d(v) = 1 \mid \widehat{E}^{(i)}\right) = (1 \pm \varepsilon)\frac{n}{k} \cdot \left(\frac{k}{2n}\right) \cdot (1 - \frac{k}{2n})^{(1\pm\varepsilon)\cdot\frac{n}{k}-1}$$

$$= \frac{1}{2} \cdot e^{-\frac{1}{2}} \cdot (1 \pm O(\varepsilon))$$

We can now conclude that $\Pr\left(d(v) = 0 \mid \widehat{E}^{(i)}, D_{\geq 2}^{(i)}\right) = (1 \pm O(\varepsilon)) \cdot 2 \cdot \left(d(v) = 1 \mid \widehat{E}^{(i)}, D_{\geq 2}^{(i)}\right)$, as conditioning on $D_{\geq 2}^{(i)}$ imply that for each vertex $v \in T$, $d(v) \in \{0, 1\}$. The assertion of the claim now immediately follows. ∎

We establish the following lower bound on the communication complexity of $\mathsf{HVP}$ on $\mathcal{D}_{\mathsf{HVP}}$.

**Lemma 5.7** (Communication Complexity of $\mathsf{HVP}$). *There exists a* universal constant $C_{\mathsf{HVP}} > 0$ *such that any protocol for $\mathsf{HVP}$ on the distribution $\mathcal{D}_{\mathsf{HVP}}$ that computes an answer with $|X \cup Y| \leq (C_{\mathsf{HVP}} \cdot n)$ w.p. at least $2/3$ needs $\Omega(n/\alpha)$ communication.*

We prove Lemma 5.7 in Section 5.3.1. Before that, we establish Theorem 6 using this lemma.

*Proof of Theorem 6.* Recall that protocol $\Pi_{\mathsf{VC}}$ is a $\delta$-error $(c \cdot \alpha)$-approximation protocol for vertex cover on $\mathcal{D}_{\mathsf{VC}}$. Note that w.p. $1 - o(1)$, $\widehat{E}^{(1)}, \ldots, \widehat{E}^{(k)}$ is an $\varepsilon$-balanced initial partitioning. Let $\mathcal{E}$ denote this event. Conditioned on $\mathcal{E}$, each edge in $G$ and in particular the edge $e^\star$ belong to each player $i \in [k]$ w.p. $(1 \pm \varepsilon)1/k$. Hence, each player is the critical player w.p. $(1 \pm \varepsilon)1/k$. Let $\mathcal{C}$ be the set of players such that $\delta_i \leq 2\delta$. We have $|\mathcal{C}| \geq k/3$ as otherwise,

$$\Pr\left(\Pi_{\mathsf{VC}} \text{ errs} \mid \mathcal{E}\right) \geq \Pr\left(\Pi_{\mathsf{VC}} \text{ errs} \mid \mathcal{E}, i^\star \notin \mathcal{C}\right) \cdot \Pr\left(i^\star \notin \mathcal{C}\right)$$
$$> 2\delta \cdot |\overline{\mathcal{C}}| \cdot (1 - \varepsilon) \cdot 1/k > 4/3\delta(1 - \varepsilon) > \delta$$

for small enough $\varepsilon > 0$, which contradicts the fact that $\Pr\left(\Pi_{\mathsf{VC}} \text{ errs} \mid \mathcal{E}\right) \leq \delta + o(1)$.

Fix an $i \in \mathcal{C}$ and let $\Pi_i$ be the message sent by $P^{(i)}$ in protocol $\Pi_{\mathsf{VC}}$. We use $\Pi_i$ to design a protocol $\Pi'$ for $\mathsf{HVP}$ on $\mathcal{D}_{\mathsf{HVP}}$ (recall that $\mathcal{D}_{\mathsf{HVP}}$ is a function of $\Pi_{\mathsf{VC}}$ and also index $i$). Given an instance of $\mathsf{HVP}$ from $\mathcal{D}_{\mathsf{HVP}}$, Alice and Bob create an instance of vertex cover sampled from $\mathcal{D}_{\mathsf{VC}} \mid \widehat{E}^{(i)}, D_{\geq 2}^{(i)}, i^\star = i$ as follows:

1. Alice plays the role of $P^{(i)}$ and Bob plays the role of all other players plus the coordinator.

2. Alice constructs the input of $P^{(i)}$ (i.e., the graph $G^{(i)}$) as follows: $(i)$ for each $u \in D_{\geq 2}^{(i)}$, Alice samples the neighbors of $u$ from $\widehat{E}^{(i)}$ according to distribution $\mathcal{D}_{\mathsf{VC}}$, and $(ii)$ for each $u \in D_1^{(i)} \cup \{u^\star\}$, she adds the edge $(u, \sigma(u))$ to $G^{(i)}$.

3. Bob constructs the inputs of all other players by letting the set $A = D_{\geq 2}^{(i)} \cup T$ and sampling their inputs according to distribution $\mathcal{D}_{\mathsf{VC}}$. This is indeed possible since the input to players in $\mathcal{D}_{\mathsf{VC}}$ are *independent* conditioned on $A, \widehat{E}^{(1)}, \dots, \widehat{E}^{(k)}, i^\star$.

Next, Alice sends the message of $P^{(i)}$ to Bob and Bob simulates the messages of all other players (without any communication) and outputs the vertex cover computed by $\Pi_{\mathsf{VC}}$ as the answer to the HVP instance. Using the definition of the distribution $\mathcal{D}_{\mathsf{HVP}}$ one can verify that the distribution of the instances sampled in this reduction matches distribution $\mathcal{D}_{\mathsf{VC}} \mid \widehat{E}^{(i)}, D_{\geq 2}^{(i)}, i^\star = i$. Hence, since the minimum vertex cover size in $G$ is at most $n/\alpha + 1$ (by picking $A \cup \{u^\star\}$), the output of $\Pi$ (i.e., the sets $X \cup Y$) is of size at most $c \cdot n$ w.p. $1 - 2\delta$ (as $i \in \mathcal{C}$). Moreover. since the edge $e^\star$ in the vertex cover instance corresponds to the pair $(u^\star, \sigma(u^\star))$ in the HVP instance, the returned solution is feasible. As $\delta \leq 0.1$, by picking the constant $c$ (in the $(c \cdot \alpha)$-approximation factor) to be smaller than $C_{\mathsf{HVP}}$ (in Lemma 5.7), we obtain that the size of $\Pi_i$ must be $\Omega(n/\alpha)$. Finally, since $|\mathcal{C}| \geq k/3$ (i.e., there are at least $k/3$ choices for player $P^{(i)}$), we obtain that the communication cost of $\Pi_{\mathsf{VC}}$ is $\Omega(nk/\alpha)$, proving the theorem. ∎

We finish this section by noting that the bound stated in Theorem 6 is in fact tight (up to an $O(\log n)$ factor) for any approximation ratio $\alpha$.

**Remark 5.8.** *The protocol in which the players group the vertices in the original graph into groups of size $\Theta(\alpha/\log n)$ (deterministically but consistently across players) and then run the algorithm in Theorem 2 on the resulting graph is an $\alpha$-approximation protocol with $\widetilde{O}(nk/\alpha)$ communication for the minimum vertex cover problem.*

Note that in Remark 5.8, we used the fact that Theorem 2 works even when the input graph has parallel edges, i.e., is a *multi-graph*.

### 5.3.1 Communication Complexity of HVP

In this section, we prove Lemma 5.7 by a reduction from the well-known set disjointness in the two-player communication model. In this problem, Alice is given a set $A \subseteq [N]$ and Bob is given a set $B \subseteq [N]$ with the promise that $|A \cap B| \in \{0, 1\}$ and their goal is to distinguish between these two cases *via two-way communication*. Let $\mathcal{D}_{\mathsf{Disj}}$ be the following distribution: start with $A = B = [N]$ and for each element $e \in A$, w.p. $1/2$, drop $e$ from both $A$ and $B$, w.p. $1/4$ drop $e$ from $A$, and with the remaining $1/4$ probability, drop $e$ from $B$. Next, pick an element $e^\star \in [N]$ uniformly at random and w.p. $1/2$ add $e^\star$ to both $A$ and $B$. It is known that solving disjointness under $\mathcal{D}_{\mathsf{Disj}}$ requires $\Omega(N)$ communication (see, e.g., [13, 63]). It also immediately follows from [13] that if instead of dropping each element w.p. exactly $1/4$, we drop them w.p. $(1 \pm \varepsilon) \cdot 1/4$ (for sufficiently small constant $\varepsilon > 0$), the distribution still remains hard.

Now let $\Pi_{\mathsf{HVP}}$ be a $\delta$-error protocol for HVP on distribution $\mathcal{D}_{\mathsf{HVP}}$. We use $\Pi_{\mathsf{HVP}}$ to create a protocol $\Pi'$ for disjointness on distribution $\mathcal{D}_{\mathsf{Disj}}$. Note that while $\Pi_{\mathsf{HVP}}$ is a one-way protocol, the protocol $\Pi'$ is allowed to use two-way communication. Given a pair of sets $(A, B)$ in $\mathcal{D}_{\mathsf{Disj}}$, we create an instance of HVP as follows:

1. Bob first communicates the *size* of $B$ to Alice.

2. The players choose a set $Z$ of size $N = |T| + |B|$ vertices from $U$ uniformly at random and consider a fixed mapping between $[N]$ and $Z$; note that $|T|$ is fixed in distribution $\mathcal{D}_{\mathsf{HVP}}$ and $|B|$ is known at this point by both players.

3. Alice lets $S = A$ and Bob lets $T = Z \setminus B$ and they pick $\sigma$ uniformly at random from $\mathcal{D}_{\mathsf{HVP}}$.

4. The players run $\Pi_{\mathsf{HVP}}$; Bob computes the sets $X$ and $Y$ and let $B' = \left( X \cup \sigma^{-1}(Y) \right) \cap B$.

5. If $|B'| > 3C_{\mathsf{HVP}} \cdot N$, Bob terminates the protocol. Otherwise the players run a lopsided set disjointness protocol (see, e.g., [23,60]) for solving the disjointness instance $(A, B')$ (with error guarantee, say, $1/10$) and output the same answer as this protocol.

Whenever $(A, B)$ is a no instance of $\mathcal{D}_{\mathsf{Disj}}$, i.e., $|A \cap B| = 1$, the distribution of the instances constructed by $\Pi'$ is $\mathcal{D}_{\mathsf{HVP}}$ (with $u^\star = A \cap B$). To see this, notice that for a fixed set $T$ in $\mathcal{D}_{\mathsf{HVP}}$, each element in $T$ is in $S$ w.p. $1/3 \cdot (1 \pm O(\varepsilon))$ and is outside $S$ with remaining probability (by Claim 5.6). This is exactly the distribution of the set $[N] \setminus B$ in $\mathcal{D}_{\mathsf{Disj}}$ conditioned on $B$. The rest follows since we are choosing the set $T$ (by the random choice of $Z$) and $\sigma$ according to distribution $\mathcal{D}_{\mathsf{HVP}}$.

On the other hand, when $|A \cap B| = 0$, the distribution of instances do *not* correspond to $\mathcal{D}_{\mathsf{HVP}}$. In fact, this is not even a valid instance of $\mathsf{HVP}$ as there is no element $u^\star$ in this instance. This means that in this case, $\Pi_{\mathsf{HVP}}$ may terminate, output a non-valid answer, or still output two sets $X \subseteq U$ and $Y \subseteq V$ with $|X \cup Y| \le C_{\mathsf{HVP}} \cdot n$. Unless the later happens, Bob is always able to distinguish this case as a Yes case of disjointness and solve the problem correctly (w.p. $1 - \delta$). Hence, in the following, we assume the worst case that $\Pi_{\mathsf{HVP}}$ outputs two sets $X$ and $Y$ even if the instance created is not a legal input of $\mathsf{VertexCollection}$. We can now argue the following key lemma.

**Lemma 5.9.** *In any instance $(S, T)$ of $\mathsf{VertexCollection}$ created by $\Pi'$, $|B'| \le 3C_{\mathsf{HVP}} \cdot |B|$ w.h.p.*

*Proof.* Consider the set $B^- := B \setminus \{u^\star\}$: this set is chosen from $U \setminus S \cup T$ uniformly at random. On the other hand, conditioned on $S, T$ (and $\sigma$), i.e., all the inputs in distribution $\mathcal{D}_{\mathsf{HVP}}$, the output of $\Pi_{\mathsf{HVP}}$ are two fixed sets $X$ and $Y$ chosen independent of $B^-$. This means that each vertex in $B^-$ belongs to $X$ w.p. $|X| / |U \setminus S \cup T|$. Similarly, each vertex in $\sigma(B^-)$ also belongs to $Y$ w.p. $|Y| / |U \setminus S \cup T|$ (as $\sigma$ is a random mapping). This ensures that $\left| B^- \cap \left( X \cup \sigma^{-1}(Y) \right) \right| \le (|X| + |Y|) / (n/2) \le 2C_{\mathsf{HVP}}$ in expectation. A simple application of Chernoff bound finalizes the proof. $\blacksquare$

*Proof of Lemma 5.7.* We first argue the correctness of the protocol $\Pi'$ and then bound its communication cost. Clearly we have $B' \subseteq B$ and moreover in the no instances of disjointness, the reduction ensures that $A \cap B \subseteq B'$; the reason is that $u^\star = A \cap B$ and since $\Pi_{\mathsf{HVP}}$ is computing two sets $X$ and $Y$ which contain either $u^\star$ or $\sigma(u^\star)$, we obtain that $u^\star \in B'$. Consequently, the probability that $\Pi'$ errs is at most $1/3$ (if $\Pi_{\mathsf{HVP}}$ errs), plus $o(1)$ (if Bob terminates the protocol (by Lemma 5.9)), plus $1/10$ (by error guarantee of the lopsided disjointness instance). This means that $\Pi'$ is a $\delta'$-error protocol for disjointness with $\delta' < 1/2$ (bounded away from half).

We now bound the communication cost of $\Pi'$. In the following, let $c$ be a constant such that communication complexity of disjointness on $\mathcal{D}_{\mathsf{Disj}}$ is at least $c \cdot N$. Since if the protocol is not terminated, $|B'| \le 3C_{\mathsf{HVP}} \cdot |B|$, the lopsided disjointness problem $(A, B')$ can be solved with $3C_{\mathsf{HVP}} \cdot O(|B|) = C_{\mathsf{HVP}} \cdot O(N)$ communication (using, e.g., the protocol of Håstad and Wigderson [34]). Now assume by contradiction that cost $\Pi_{\mathsf{HVP}}$ is $o(n/\alpha) = o(N)$. This means that the total communication cost of $\Pi'$ is $O(\log N) + o(N) + C_{\mathsf{HVP}} \cdot O(N)$. By taking $C_{\mathsf{HVP}} \ll c$ but still a constant to suppress the constant in the $O(N)$ term above, the total cost of $\Pi'$ can be made smaller than $c \cdot N$. This

contradicts the fact that communication complexity of disjointness on $\mathcal{D}_{\mathsf{Disj}}$ is at least $c \cdot N$, finalizing the proof. $\blacksquare$

## Acknowledgements

## References

[1] S. Abbar, S. Amer-Yahia, P. Indyk, S. Mahabadi, and K. R. Varadarajan. Diverse near neighbor problem. In *Symposuim on Computational Geometry 2013, SoCG '13, Rio de Janeiro, Brazil, June 17-20, 2013*, pages 207–214, 2013.

[2] P. K. Agarwal, S. Har-Peled, and K. R. Varadarajan. Approximating extent measures of points. *J. ACM*, 51(4):606–635, 2004.

[3] K. J. Ahn and S. Guha. Linear programming in the semi-streaming model with application to the maximum matching problem. *Inf. Comput.*, 222:59–79, 2013.

[4] K. J. Ahn and S. Guha. Access to data and number of iterations: Dual primal algorithms for maximum matching under resource constraints. In *Proceedings of the 27th ACM on Symposium on Parallelism in Algorithms and Architectures, SPAA 2015, Portland, OR, USA, June 13-15, 2015*, pages 202–211, 2015.

[5] K. J. Ahn, S. Guha, and A. McGregor. Analyzing graph structure via linear measurements. In *Proceedings of the Twenty-third Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA '12, pages 459–467. SIAM, 2012.

[6] K. J. Ahn, S. Guha, and A. McGregor. Graph sketches: sparsification, spanners, and subgraphs. In *Proceedings of the 31st ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems, PODS 2012, Scottsdale, AZ, USA, May 20-24, 2012*, pages 5–14, 2012.

[7] Y. Ai, W. Hu, Y. Li, and D. P. Woodruff. New characterizations in turnstile streams with applications. In *31st Conference on Computational Complexity, CCC 2016, May 29 to June 1, 2016, Tokyo, Japan*, pages 20:1–20:22, 2016.

[8] N. Alon, N. Nisan, R. Raz, and O. Weinstein. Welfare maximization with limited interaction. In *IEEE 56th Annual Symposium on Foundations of Computer Science, FOCS 2015, Berkeley, CA, USA, 17-20 October, 2015*, pages 1499–1512, 2015.

[9] S. Assadi, S. Khanna, and Y. Li. On estimating maximum matching size in graph streams. In *Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2017, Barcelona, Spain, Hotel Porta Fira, January 16-19*, pages 1723–1742, 2017.

[10] S. Assadi, S. Khanna, Y. Li, and G. Yaroslavtsev. Maximum matchings in dynamic graph streams and the simultaneous communication model. In *Proceedings of the Twenty-Seventh Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2016, Arlington, VA, USA, January 10-12, 2016*, pages 1345–1364, 2016.

[11] A. Badanidiyuru, B. Mirzasoleiman, A. Karbasi, and A. Krause. Streaming submodular maximization: massive data summarization on the fly. In *The 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '14, New York, NY, USA - August 24 - 27, 2014*, pages 671–680, 2014.

[12] M. Balcan, S. Ehrlich, and Y. Liang. Distributed k-means and k-median clustering on general communication topologies. In *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States.*, pages 1995–2003, 2013.

[13] Z. Bar-Yossef, T. S. Jayram, R. Kumar, and D. Sivakumar. An information statistics approach to data stream and communication complexity. In *43rd Symposium on Foundations of Computer Science (FOCS 2002), 16-19 November 2002, Vancouver, BC, Canada, Proceedings*, pages 209–218, 2002.

[14] S. Baswana, M. Gupta, and S. Sen. Fully dynamic maximal matching in o(log n) update time. *SIAM J. Comput.*, 44(1):88–113, 2015.

[15] M. Bateni, A. Bhaskara, S. Lattanzi, and V. S. Mirrokni. Distributed balanced clustering via mapping coresets. In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, pages 2591–2599, 2014.

[16] S. Bhattacharya, M. Henzinger, and G. F. Italiano. Deterministic fully dynamic data structures for vertex cover and matching. In *Proceedings of the Twenty-Sixth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2015, San Diego, CA, USA, January 4-6, 2015*, pages 785–804, 2015.

[17] S. Bhattacharya, M. Henzinger, D. Nanongkai, and C. E. Tsourakakis. Space- and time-efficient algorithm for maintaining dense subgraphs on one-pass dynamic streams. In *Proceedings of the Forty-Seventh Annual ACM on Symposium on Theory of Computing, STOC 2015, Portland, OR, USA, June 14-17, 2015*, pages 173–182, 2015.

[18] L. Bulteau, V. Froese, K. Kutzkov, and R. Pagh. Triangle counting in dynamic graph streams. *Algorithmica*, 76(1):259–278, 2016.

[19] A. Chakrabarti, G. Cormode, and A. McGregor. Robust lower bounds for communication and stream computation. In *Proceedings of the 40th Annual ACM Symposium on Theory of Computing, Victoria, British Columbia, Canada, May 17-20, 2008*, pages 641–650, 2008.

[20] R. Chitnis, G. Cormode, H. Esfandiari, M. Hajiaghayi, A. McGregor, M. Monemizadeh, and S. Vorotnikova. Kernelization via sampling with applications to finding matchings and related problems in dynamic graph streams. In *Proceedings of the Twenty-Seventh Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2016, Arlington, VA, USA, January 10-12, 2016*, pages 1326–1344, 2016.

[21] R. H. Chitnis, G. Cormode, M. T. Hajiaghayi, and M. Monemizadeh. Parameterized streaming: Maximal matching and vertex cover. In *Proceedings of the Twenty-Sixth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2015, San Diego, CA, USA, January 4-6, 2015*, pages 1234–1251, 2015.

[22] M. Crouch and D. S. Stubbs. Improved streaming algorithms for weighted matching, via unweighted matching. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques, APPROX/RANDOM 2014, September 4-6, 2014, Barcelona, Spain*, pages 96–104, 2014.

[23] A. Dasgupta, R. Kumar, and D. Sivakumar. Sparse and lopsided set disjointness via information theory. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques - 15th International Workshop, APPROX 2012, and 16th International Workshop, RANDOM 2012, Cambridge, MA, USA, August 15-17, 2012. Proceedings*, pages 517–528, 2012.

[24] S. Dobzinski, N. Nisan, and S. Oren. Economic efficiency requires interaction. In *Symposium on Theory of Computing, STOC 2014, New York, NY, USA, May 31 - June 03, 2014*, pages 233–242, 2014.

[25] D. P. Dubhashi and A. Panconesi. *Concentration of Measure for the Analysis of Randomized Algorithms.* Cambridge University Press, 2009.

[26] S. Eggert, L. Kliemann, and A. Srivastav. Bipartite graph matchings in the semi-streaming model. In *Algorithms - ESA 2009, 17th Annual European Symposium, Copenhagen, Denmark, September 7-9, 2009. Proceedings*, pages 492–503, 2009.

[27] L. Epstein, A. Levin, J. Mestre, and D. Segev. Improved approximation guarantees for weighted matching in the semi-streaming model. *SIAM J. Discrete Math.*, 25(3):1251–1265, 2011.

[28] H. Esfandiari, M. Hajiaghayi, and M. Monemizadeh. Finding large matchings in semi-streaming. In *IEEE International Conference on Data Mining Workshops, ICDM Workshops 2016, December 12-15, 2016, Barcelona, Spain.*, pages 608–614, 2016.

[29] H. Esfandiari, M. T. Hajiaghayi, V. Liaghat, M. Monemizadeh, and K. Onak. Streaming algorithms for estimating the matching size in planar graphs and beyond. In *Proceedings of the Twenty-Sixth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2015, San Diego, CA, USA, January 4-6, 2015*, pages 1217–1233, 2015.

[30] J. Feigenbaum, S. Kannan, A. McGregor, S. Suri, and J. Zhang. On graph problems in a semi-streaming model. *Theor. Comput. Sci.*, 348(2-3):207–216, 2005.

[31] A. Goel, M. Kapralov, and S. Khanna. On the communication and streaming complexity of maximum bipartite matching. In *Proceedings of the Twenty-third Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA '12, pages 468–485. SIAM, 2012.

[32] V. Guruswami and K. Onak. Superlinear lower bounds for multipass graph processing. In *Proceedings of the 28th Conference on Computational Complexity, CCC 2013, K.lo Alto, California, USA, 5-7 June, 2013*, pages 287–298, 2013.

[33] A. Hassidim, J. A. Kelner, H. N. Nguyen, and K. Onak. Local graph partitions for approximation and testing. In *50th Annual IEEE Symposium on Foundations of Computer Science, FOCS 2009, October 25-27, 2009, Atlanta, Georgia, USA*, pages 22–31, 2009.

[34] J. Håstad and A. Wigderson. The randomized communication complexity of set disjointness. *Theory of Computing*, 3(1):211–219, 2007.

[35] Z. Huang, B. Radunovic, M. Vojnovic, and Q. Zhang. Communication complexity of approximate matching in distributed graphs. In *32nd International Symposium on Theoretical Aspects of Computer Science, STACS 2015, March 4-7, 2015, Garching, Germany*, pages 460–473, 2015.

[36] P. Indyk, S. Mahabadi, M. Mahdian, and V. S. Mirrokni. Composable core-sets for diversity and coverage maximization. In *Proceedings of the 33rd ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems, PODS'14, Snowbird, UT, USA, June 22-27, 2014*, pages 100–108, 2014.

[37] M. Kapralov. Better bounds for matchings in the streaming model. In *Proceedings of the Twenty-Fourth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2013, New Orleans, Louisiana, USA, January 6-8, 2013*, pages 1679–1697, 2013.

[38] M. Kapralov, S. Khanna, and M. Sudan. Approximating matching size from random streams. In *Proceedings of the Twenty-Fifth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2014, Portland, Oregon, USA, January 5-7, 2014*, pages 734–751, 2014.

[39] M. Kapralov, S. Khanna, and M. Sudan. Streaming lower bounds for approximating MAX-CUT. In *SODA*, 2015.

[40] M. Kapralov, Y. T. Lee, C. Musco, C. Musco, and A. Sidford. Single pass spectral sparsification in dynamic streams. In *55th IEEE Annual Symposium on Foundations of Computer Science, FOCS 2014, Philadelphia, PA, USA, October 18-21, 2014*, pages 561–570, 2014.

[41] M. Kapralov and D. Woodruff. Spanners and sparsifiers in dynamic streams. *PODC*, 2014.

[42] H. J. Karloff, S. Suri, and S. Vassilvitskii. A model of computation for mapreduce. In *Proceedings of the Twenty-First Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2010, Austin, Texas, USA, January 17-19, 2010*, pages 938–948, 2010.

[43] C. Konrad. Maximum matching in turnstile streams. In *Algorithms - ESA 2015 - 23rd Annual European Symposium, Patras, Greece, September 14-16, 2015, Proceedings*, pages 840–852, 2015.

[44] C. Konrad, F. Magniez, and C. Mathieu. Maximum matching in semi-streaming with few passes. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques - 15th International Workshop, APPROX 2012, and 16th International Workshop, RANDOM 2012, Cambridge, MA, USA, August 15-17, 2012. Proceedings*, pages 231–242, 2012.

[45] E. Kushilevitz and N. Nisan. *Communication complexity*. Cambridge University Press, 1997.

[46] S. Lattanzi, B. Moseley, S. Suri, and S. Vassilvitskii. Filtering: a method for solving graph problems in mapreduce. In *SPAA 2011: Proceedings of the 23rd Annual ACM Symposium on Parallelism in Algorithms and Architectures, San Jose, CA, USA, June 4-6, 2011 (Co-located with FCRC 2011)*, pages 85–94, 2011.

[47] Y. Li, H. L. Nguyen, and D. P. Woodruff. Turnstile streaming algorithms might as well be linear sketches. In *Symposium on Theory of Computing, STOC 2014, New York, NY, USA, May 31 - June 03, 2014*, pages 174–183, 2014.

[48] A. McGregor. Finding graph matchings in data streams. In *Approximation, Randomization and Combinatorial Optimization, Algorithms and Techniques, 8th International Workshop on Approximation Algorithms for Combinatorial Optimization Problems, APPROX 2005 and 9th InternationalWorkshop on Randomization and Computation, RANDOM 2005, Berkeley, CA, USA, August 22-24, 2005, Proceedings*, pages 170–181, 2005.

[49] A. McGregor. Graph stream algorithms: a survey. *SIGMOD Record*, 43(1):9–20, 2014.

[50] A. McGregor, D. Tench, S. Vorotnikova, and H. T. Vu. Densest subgraph in dynamic graph streams. In *Mathematical Foundations of Computer Science 2015 - 40th International Symposium, MFCS 2015, Milan, Italy, August 24-28, 2015, Proceedings, Part II*, pages 472–482, 2015.

[51] A. McGregor and S. Vorotnikova. Planar matching in streams revisited. *To appear in Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques - 19th International Workshop, APPROX 2016, and 20th International Workshop, RANDOM 2016*, 2016.

[52] V. S. Mirrokni and M. Zadimoghaddam. Randomized composable core-sets for distributed submodular maximization. In *Proceedings of the Forty-Seventh Annual ACM on Symposium on Theory of Computing, STOC 2015, Portland, OR, USA, June 14-17, 2015*, pages 153–162, 2015.

[53] B. Mirzasoleiman, A. Karbasi, R. Sarkar, and A. Krause. Distributed submodular maximization: Identifying representative elements in massive data. In *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States.*, pages 2049–2057, 2013.

[54] S. Muthukrishnan. *Data streams: Algorithms and applications*. Now Publishers Inc, 2005.

[55] O. Neiman and S. Solomon. Simple deterministic algorithms for fully dynamic maximal matching. In *Symposium on Theory of Computing Conference, STOC'13, Palo Alto, CA, USA, June 1-4, 2013*, pages 745–754, 2013.

[56] H. N. Nguyen and K. Onak. Constant-time approximation algorithms via local improvements. In *49th Annual IEEE Symposium on Foundations of Computer Science, FOCS 2008, October 25-28, 2008, Philadelphia, PA, USA*, pages 327–336, 2008.

[57] K. Onak, D. Ron, M. Rosen, and R. Rubinfeld. A near-optimal sublinear-time algorithm for approximating the minimum vertex cover size. In *Proceedings of the Twenty-Third Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2012, Kyoto, Japan, January 17-19, 2012*, pages 1123–1131, 2012.

[58] K. Onak and R. Rubinfeld. Maintaining a large matching and a small vertex cover. In *Proceedings of the 42nd ACM Symposium on Theory of Computing, STOC 2010, Cambridge, Massachusetts, USA, 5-8 June 2010*, pages 457–464, 2010.

[59] M. Parnas and D. Ron. Approximating the minimum vertex cover in sublinear time and a connection to distributed algorithms. *Theor. Comput. Sci.*, 381(1-3):183–196, 2007.

[60] M. Patrascu. Unifying the landscape of cell-probe lower bounds. *SIAM J. Comput.*, 40(3):827–847, 2011.

[61] A. Paz and G. Schwartzman. A $(2 + \varepsilon)$-approximation for maximum weight matching in the semi-streaming model. In *Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2017, Barcelona, Spain, Hotel Porta Fira, January 16-19*, pages 2153–2161, 2017.

[62] J. M. Phillips, E. Verbin, and Q. Zhang. Lower bounds for number-in-hand multiparty communication complexity, made easy. In *Proceedings of the Twenty-Third Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2012, Kyoto, Japan, January 17-19, 2012*, pages 486–501, 2012.

[63] A. A. Razborov. On the distributional complexity of disjointness. *Theor. Comput. Sci.*, 106(2):385–390, 1992.

[64] S. Solomon. Fully dynamic maximal matching in constant update time. *CoRR*, abs/1604.08491. To appear in FOCS, 2016.

[65] A. C. Yao. Lower bounds to randomized algorithms for graph properties (extended abstract). In *28th Annual Symposium on Foundations of Computer Science, Los Angeles, California, USA, 27-29 October 1987*, pages 393–400, 1987.

[66] Y. Yoshida, M. Yamamoto, and H. Ito. Improved constant-time approximation algorithms for maximum matchings and other optimization problems. *SIAM J. Comput.*, 41(4):1074–1093, 2012.

# A   Proof of Lemma 4.1

Consider the edges of $E_{AB}$ assigned to the graph $G^{(i)}$. It is easy to see that by the choice of $E_{AB}$ and further partitioning the edges between the $k$ players, the graph $G^{(i)}$ on the set of edges $E_{AB}^{(i)}$ forms a random bipartite graph. Hence, proving Lemma 4.1 reduces to proving the following property of random bipartite graphs.

Let $\mathcal{G}(n, n, 1/n)$ denote the family of random bipartite graphs where each side of the bipartition contains $n$ vertices, and each edge is present w.p. $1/n$. We will show that if we sample a random graph $G \in \mathcal{G}(n, n, 1/n)$, then w.p. at least $1 - 1/n^2$, it contains an *induced matching* of size $\Omega(n)$. We emphasize here that the notion of induced matching is with respect to the entire graph and not only with respect to the vertices included in the induced matching.

Our proof will use the following pair of elementary propositions.

**Proposition A.1.** *Suppose we assign $N$ balls uniformly at random to $M > N$ bins. Let $B$ be an arbitrary fixed subset of bins. Then with probability at least $1 - \frac{1}{N^3}$, there are at least $\left(\frac{|B|}{M}\right) \cdot \frac{N}{e} - o(N)$ bins in $B$ that contain exactly one ball, assuming $N$ is sufficiently large.*

*Proof.* We arbitrarily number the balls 1 through $N$ and the bins 1 through $M$. W.l.o.g. assume that the bins in $B$ are numbered 1 through $|B|$. For $1 \leq i \leq |B|$, let $Z_i$ be the 0/1 random variable that indicates whether or not bin $i \in B$ receives exactly one ball, and furthermore, let $Z = \sum_{i=1}^{|B|} Z_i$. Then

$$\Pr[Z_i = 1] = \binom{N}{1} \cdot \left(\frac{1}{M}\right) \cdot \left(1 - \frac{1}{M}\right)^{N-1} \geq \frac{N}{M} \cdot \left(\frac{1}{e} - o(1)\right).$$

Hence $E[Z] \geq \left(\frac{|B|}{M}\right) \cdot \frac{N}{e} - o(N)$. We now wish to argue that the value of $Z$ is concentrated around its expectation. However, we can not directly invoke the standard Chernoff bound since the variables $Z_i$'s are not independent. We will instead utilize the more general version stated in Proposition 2.2.

Let $X_j \in [1..M]$ denote the index of the bin in which the $j_{th}$ ball lands. Given the variables $X_1, X_2, ..., X_N$, we can define the function $f(X_1, X_2, ...., X_N)$ to be the number of bins in $B$ that receive exactly one ball. Note that $f$ is completely determined by the variables $X_1, X_2, ..., X_N$ and that $E[f] = E[Z] \geq \left(\frac{|B|}{M}\right) \cdot \frac{N}{e} - o(N)$. It is easy to see that the function $f$ satisfies the Lipschitz property with $d = 2$ since changing the assignment of any single ball, can reduce or increase the number of bins in $B$ with exactly one ball by at most 2. We can thus invoke Proposition 2.2 with $t = 4\sqrt{N \ln N}$, completing the proof. ∎

**Proposition A.2.** *For sufficiently large $n$, with probability at least $1 - 1/n^3$, a graph $G(L \cup R, E)$ drawn from $\mathcal{G}(n, n, 1/n)$ satisfies the following properties:*

(a) *The set $S \subseteq L$ of all vertices in $L$ with degree exactly one in $G$ has size $n/e \pm o(n)$.*

(b) *The set $T \subseteq R$ of vertices defined as all vertices in $R$ with no edges to $L \setminus S$ has size at least $n/e - o(n)$.*

*Proof.* To see property (a), let us define 0/1 random variables $X_1, X_2, ..., X_n$ where $X_i = 1$ iff vertex $i \in L$ has degree exactly one in $G$. Then $\Pr[X_i = 1] = (1 - \frac{1}{n})^{n-1} = 1/e - o(1)$ for sufficiently large $n$. Thus $E[\sum_{i=1}^{n} X_i] = n/e - o(n)$, and using Chernoff bound (Proposition 2.1 with $t = 4\sqrt{n \ln n}$) implies that with probability at least $1 - 2/n^4$, there is a set $S \subseteq L$ of size $n/e \pm o(n)$ whose vertices have degree exactly one in $G$.

To see property (b), fix a set $S$ of degree 1 vertices in $L$. Let us define 0/1 random variables $Y_1, Y_2, ..., Y_n$ where $Y_i = 1$ iff vertex $i \in R$ receives no edges from vertices in $L \setminus S$. Then $\Pr[Y_i = 1] = (1 - \frac{1}{n})^{|L \setminus S|} \geq 1/e - o(1)$ for sufficiently large $n$. Thus $E[\sum_{i=1}^{n} Y_i] \geq n/e - o(n)$, and using Chernoff bound (Proposition 2.1 with $t = 4\sqrt{n \ln n}$) implies that with probability at least $1 - 2/n^4$, there is a set $T \subseteq R$ of size at least $n/e - o(n)$ whose vertices do not have any edges to $L \setminus S$.

Thus both properties (a) and (b) hold with probability at least $1 - 1/n^3$, as desired. ∎

**Lemma A.3.** *Let $G(L \cup R, E)$ be drawn from $\mathcal{G}(n, n, 1/n)$. Then for sufficiently large $n$, with probability at least $1 - 1/n^2$, $G$ contains an* induced matching *of size $n/e^3 - o(n)$.*

*Proof.* By Proposition A.2, we know with probability at least $1 - 1/n^3$, the graph $G(L \cup R, E)$ satisfies properties (a) and (b). We will assume from here on that this event, denoted by $\mathcal{E}$, has occurred. We first observe that conditioned on the event $\mathcal{E}$, and for any choice of sets $S$ and $T$ as defined in Proposition A.2 as well as edges from the set $L \setminus S$ to $R \setminus T$, sampling a graph $G$ from $\mathcal{G}(n, n, 1/n)$ is equivalent to assigning each vertex in $S$ a uniformly at random neighbor in $R$.

Now invoking Proposition A.1, with $N = |S|$, $B = T$, we know that with probability at least $1 - O(1/n^3)$, there is a set $T' \subseteq T$ of size at least

$$\frac{|T|}{n} \cdot \frac{|S|}{e} - o(|S|) = \left(\frac{1}{e} - o(1)\right) \cdot \left(\frac{n/e - o(n)}{e}\right) - o(n) \geq \frac{n}{e^3} - o(n)$$

such that each vertex in $T'$ receives exactly one ball from $S$ (i.e. receives exactly one edge from the vertices in $S$). Let $S' \subseteq S$ be the set of vertices that "supply a ball" (i.e. an edge) to vertices in $T'$. Since by definition the vertices in $T$ receive edges only from $S$, and since all vertices in $S$ have degree exactly one, the set $S' \cup T'$ of vertices induces a matching of size at least $n/e^3 - o(n)$ in $G$, as asserted in the lemma. ∎

The lower bound in Lemma 4.1 now follows from Lemma A.3 for the family of bipartite graphs with $n/\alpha$ vertices on each side. The upper bound is a simple application of Chernoff bound on the number of edges from $E_{\overline{AB}}$ that are assigned to $G^{(i)}$.