

## Weeks 4-5: Streaming and Parallel Algorithms

*Aleksandar Nikolov*

## 1 Motivation and Model

In the next few lectures we will study algorithms which are *extremely* efficient, in terms of space, and usually in terms of time complexity as well. This is motivated by real-world scenarios in which we want to process enormous amounts of data quickly. For example:

- A network service provider wants to keep track of the traffic going through the network, in order to detect unusual patterns, which may be a sign of malicious activity. For example, an increased traffic to a particular IP address may suggest a Denial of Service attack is under way. Such unusual traffic patterns should be detected almost in real time, however the throughput of the network is so large that it's not feasible to store all the traffic data and analyze it. So, you need a solution which works online (i.e. processes the traffic as it passes through the routers), and accurately summarizes the data in a very concise “sketch.”
- Processing data in a very large database can be very expensive in terms of time. Often it's a good idea to quickly compute some rough statistics about the data before embarking on a more expensive analysis. Because of its size, the database usually cannot be entirely loaded into the computer memory, and most of it will be on a relatively slow hard disk. The fastest way to process such a database is to do so sequentially, making only a few passes over the data, as opposed jumping back and forth between different records. So, we are interested in statistics which can be computed approximately in a few sequential passes over the data.

Scenarios like these motivate the simple and elegant streaming computation model. In the streaming model, the algorithm  $\mathcal{A}$  works in time steps, and receives one update per time step. In the simplest version of the model, an update is just an integer in the set  $\{1, \dots, n\}$ , which is meant as an identifier of some object. For example, the integer can encode an IP address. The sequence of updates is called the *stream*. The total number of updates  $m$  is called the length of the stream, and may or may not be explicitly given to the algorithm. After receiving an update, and before receiving the next one,  $\mathcal{A}$  updates its memory. Moreover, at any point during its execution, the algorithm should be able to report an answer to a given question, e.g. report the number of distinct integers seen so far. The crucial constraint is that, at any time step, the algorithm is only allowed to store in its memory a number of bits which is bounded by a polynomial in  $\log n$  and  $\log m$ . I.e. the algorithm's memory may never exceed  $O(\log^c(nm))$  bits, where  $c$  is some fixed constant, say 1 or 2. This constraint models the fact that, in the applications we mentioned above, the size of the data (corresponding to  $m$ ), and the size of the universe it is coming from (corresponding to  $n$ ), are very large. While the model does not pose an explicit constraint on the running time, it's desirable that updates are as efficient as possible.

Often when we analyze algorithm in the streaming model, we count the number of words used by the algorithm, where a word is a block of  $O(\log(nm))$  bits, which can store, e.g. a single variable, or a single cell in an array. We will adopt this convention in these notes.

Many fundamental problems in the streaming model are conveniently summarized by the *frequency vector*  $f$ , which gives the number of times each element of the universe appears in the stream. I.e. if we have a stream  $\sigma = (i_1, \dots, i_t)$  consisting of updates in  $\{1, \dots, n\}$ , then  $f \in \mathbb{Z}^n$  is defined by  $f_i = |\{t : i_t = i\}|$ . Note that the streaming algorithm will not actually store the frequency vector  $f$ , because that would require space linear in  $n$ . However, it helps to refer to the vector when defining problems or analyzing algorithms.

In some versions of the model, updates can have a richer meaning. For example, in the *turnstile* model, an update is a pair  $(i, s)$ , where  $i \in \{1, \dots, n\}$ , and  $s \in \{-1, 1\}$ . As before,  $i$  is just an identifier;  $s$  indicates “entering,” when  $s = +1$ , or “leaving,” when  $s = -1$ . For example,  $i$  can identify a particular subway station, and  $s$  can indicate whether at the given time step a customer has entered, or left the station. In the turnstile model we define the frequency vector  $f$  for a stream  $\sigma = \{(i_1, s_1), \dots, (i_m, s_m)\}$  by  $f_i = \sum_{t:i_t=i} s_t$ .

**Exercise 1.** *Suppose the stream has length  $n - 1$ , and consists of  $n - 1$  of the integers  $\{1, \dots, n\}$ . All updates are distinct, and the integers can appear in an arbitrary order. Design an algorithm which makes a single pass over the data, keeps only a constant number of words of memory, and finds the missing integer from the stream. Adapt your solution to finding two missing integers in a stream of  $n - 2$  distinct updates from  $\{1, \dots, n\}$ .*

**Exercise 2.** *Give an algorithm which samples a random element in a stream  $\sigma = (i_1, \dots, i_m)$  of updates in  $\{1, \dots, n\}$  so that the probability that element  $i$  is sampled equals  $\frac{|\{t:i_t=i\}|}{m}$ . The algorithm should keep only a single variable, which at every time step is equal to the random sample from the portion of the stream seen so far. Do not assume the algorithm knows the length of the stream.*

## 2 Frequent Elements

It should be surprising to you that it is at all possible to do anything in the streaming model. Often approximation and randomization are both absolutely necessary to solve a problem. We will start with a rare example of a problem for which there are efficient deterministic algorithms.

As a warm-up, we consider the Majority problem. You are given a stream  $\sigma = (i_1, \dots, i_m)$  of updates in  $[n] = \{1, \dots, n\}$ . If there exists an  $i \in [n]$  such that more than half the updates in  $\sigma$  are equal to  $i$ , the algorithm should output  $i$ . If no such majority element exists, the algorithm can output any element. (Notice that only one such element can exist.) The algorithm below (due to Boyer and Moore) solves this problem with *only two words of memory*:

MAJORITY( $\sigma$ )

```

1  element =  $i_1$ 
2  count = 1
3  for  $t = 2$  to  $m$ 
4      if element ==  $i_t$ 
5          count = count + 1
6      elseif count > 0
7          count = count - 1
8      else element =  $i_t$ 
9          count = 1
10 return element
```

**Theorem 1.** *If there exists an element  $i$  such that more than half the updates in  $\sigma$  are equal to  $i$ , then the MAJORITY algorithm outputs  $i$ . Moreover, at any point during the execution of the algorithm,  $f_{element} \leq count + m/2$ .*

*Proof.* We group the updates in the stream in pairs as follows. At the start, we leave the first update  $i_1$  unpaired and move to  $t = 2$ . If  $i_t = element$ , or if  $i_t \neq element$  but  $count = 0$ , we leave the update  $i_t$  unpaired for now and move on to the next value of  $t$ . If  $i_t \neq element$  and  $count > 0$ , we pair  $i_t$  with one of the prior updates  $i_s$ ,  $s < t$ , for which  $i_s = element$ , and  $i_s$  is still unpaired. Then we move on to the next value of  $t$ . The idea behind this pairing procedure is that we view the event of  $count$  decreasing at time step  $t$  as  $i_t$  “taking away” one of the prior instances of  $element$ .

The following fact is easily shown by induction:

- At any time step  $t$ ,  $count$  is equal to the number of yet unpaired updates equal to  $element$ . For all  $j \neq element$ , all updates equal to  $j$  are paired.

Let  $i$  be as in the statement of the theorem. Observe that for every pair  $(i_s, i_t)$  we have  $i_s \neq i_t$ . So, since there are at most  $m/2$  pairs in the pairing, at most  $m/2$  of the updates equal to  $i$  are paired. This means that at least  $f_i - m/2 > 0$  updates equal to  $i$  are left unpaired at the end of the execution of the algorithm. By the claim above, this means the final value of  $element$  returned by the algorithm is  $i$ , and  $count \geq f_i - m/2$ . Notice that this also proves the claim after “moreover” for the majority element  $i$ . If there is no majority element, that claim holds trivially.  $\square$

If we want to also verify that the element  $i$  returned by the algorithm is a majority element, we can just make a second pass over the stream and just count the number of occurrences of  $i$ .

Next we study a generalization of MAJORITY which finds all elements that appear in more than  $1/k$  fraction of the updates. This algorithm is due to Misra and Gries.

FREQUENT( $\sigma, k$ )

```

1   $S = \emptyset$ 
2  for  $t = 1$  to  $m$ 
3      if  $\exists x \in S$  such that  $x.elem == i_t$ 
4           $x.count = x.count + 1$ 
5      elseif  $|S| < k - 1$ 
6          Create an element  $x$  with  $x.elem = i_t$  and  $x.count = 1$ 
7           $S = S \cup x$ 
8      else for  $x \in S$ 
9           $x.count = x.count - 1$ 
10         if  $x.count == 0$ 
11              $S = S \setminus \{x\}$ 
12  Return  $S$ 

```

**Theorem 2.** *The set  $S$  output by FREQUENT contains all  $i \in [n]$  such that  $f_i > m/k$ . Moreover, for any  $x \in S$ ,  $f_{x.elem} \leq x.count + m/k$ .*

*Proof.* We group the updates in the stream in groups of size  $k$ . Let us say that an element  $i \in [n]$  is represented in  $S$  if there exists an  $x \in S$  such that  $x.elem = i$ . Starting from  $t = 1$ , we look at the following cases:

1. If  $i_t$  is represented in  $S$ , we leave  $i_t$  ungrouped for now, and move to the next value of  $t$ .
2. Similarly, if  $i_t$  is not represented in  $S$ , but  $|S| < k - 1$ , we leave  $i_t$  ungrouped for now, and move to the next value of  $t$ .
3. If  $i_t$  is not represented in  $S$  and  $|S| = k - 1$ , then we group  $i_t$  with  $k - 1$  previously ungrouped distinct updates represented in  $S$ . I.e. we find  $k - 1$  updates, occurring at times  $s_1, \dots, s_{k-1}$ , all preceding  $i_t$  and ungrouped, with the property that  $i_{s_1}, \dots, i_{s_{k-1}}$  are all distinct, and represented in  $S$ . We group  $i_t$  with  $i_{s_1}, \dots, i_{s_{k-1}}$  and move to the next value of  $t$ .

It is easy to prove the following crucial claim by induction:

- At any time step  $t$ , for any  $x \in S$ ,  $x.count$  is equal to the number of yet ungrouped updates equal to  $x.elem$ . For all  $j \in [n]$  not represented in  $S$ , all updates equal to  $j$  are grouped.

Let  $i$  be such that  $f_i > m/k$ . Observe that the groups are defined so that they contain distinct elements of  $[n]$ . Because there can be at most  $m/k$  groups, at least  $f_i - m/k > 0$  updates equal to  $i$  are left ungrouped. By the claim above, this means that  $i$  is represented in  $S$ , and that  $x.count \geq f_i - m/k$ , where  $x \in S$  is such that  $x.elem = i$ . This also proves the claim after “moreover” for all  $i$  such that  $f_i > m/k$ . The claim is trivial for all other elements.  $\square$

Notice that  $\text{FREQUENT}(\sigma, k)$  will output all elements with  $f_i > m/k$ , but it may also output some infrequent elements. To be sure which elements have frequency greater than  $m/k$ , and which do not, you need to make another pass over the stream.

**Exercise 3.** Show how to implement  $\text{FREQUENT}$  so that each update is as efficient as possible. What data structure would you use to represent  $S$  in memory? What is the worst-case time complexity of an update? What is the amortized time complexity of an update (i.e. the total time taken by the algorithm, divided by  $m$ ).

**Exercise 4.** Design an algorithm in the streaming model that, given values  $\phi$  and  $\varepsilon$  such that  $0 < \varepsilon < \phi < 1$ , outputs a set  $S \subseteq [n]$  such that:

1. If  $i$  is such that  $f_i > \phi m$ , then  $i \in S$ ;
2. If  $i$  is such that  $f_i < (\phi - \varepsilon)m$ , then  $i \notin S$ .

Your algorithm should use  $O(\frac{1}{\varepsilon})$  words of memory.

### 3 Distinct Elements Count

In the distinct elements count problem we want to know how many distinct integers we have seen in the stream so far. In terms of the frequency vector, we want to approximate  $F_0 = |\{i : f_i > 0\}|$ . For this problem it turns out that we need to allow both approximation and randomization in order to satisfy the constraints of the streaming model.

### 3.1 Warm-Up

We first give an algorithm which makes the (unrealistic) assumption that we know a number  $\tilde{F}_0$  such that

$$F_0 \leq \tilde{F}_0 \leq 2F_0. \quad (1)$$

Our goal in this case is to refine this “rough” estimate  $\tilde{F}_0$  to a much more precise one by processing the stream. In the next subsection we will see how to remove this unrealistic assumption by building on the ideas from this subsection.

The algorithm to refine the rough estimate  $\tilde{F}_0$  works by a simple sampling strategy:

DISTINCT-SIMPLE( $\sigma, k, \tilde{F}_0$ )

```
1  $S = \emptyset$ 
2  $d = \lceil \log_2(\tilde{F}_0/k) \rceil$ 
3  $L = \lceil \log_2 n \rceil$ 
4 Pick a hash function  $h : [n] \rightarrow \{0, 1\}^L$ 
5 for  $t = 1$  to  $m$ 
6     if  $h(i_t) \in 0^d\{0, 1\}^{L-d}$  and  $i_t \notin S$ 
7          $S = S \cup \{i_t\}$ 
8 return  $\hat{F}_0 = 2^d \cdot |S|$ 
```

A few clarifications are in order. In the analysis we assume that  $h$  behaves like a random function. I.e., we assume that  $h(1), \dots, h(n)$  are  $n$  independent random variables, uniformly distributed in  $\{0, 1\}^L$ . This is essentially the simple uniform hashing assumption you may remember from your data structures course. The assumption can be removed using a construction similar to universal hashing, and we discuss this a little bit later. We also clarify the notation: by  $h(i_t) \in 0^d\{0, 1\}^{L-d}$  we simply mean that the leftmost  $d$  bits of  $h(i_t)$  are all 0.

So, in words, the algorithm works as follows: it uses the hash function to keep a set  $S$  which contains each element that appears in the stream with probability  $2^{-d}$ . In expectation,  $S$  then contains  $2^{-d}$  fraction of the elements that appear in the stream, so our estimate of the number of distinct elements is  $2^d|S|$ . The value  $d$  is chosen so that the expected size of  $S$  is at most  $k$ : this is the only place where we use  $\tilde{F}_0$ ; if we did not know such an estimate, we would not know how to set  $d$ . There is a natural trade-off here: the larger we pick  $k$  (and, therefore, the smaller we pick  $d$ ), the more space our algorithm uses, but the more accurate its estimate becomes.

To analyze the algorithm, we will recall a basic concept from probability theory: *variance*. The variance of a real-valued random variable  $X$  is defined as  $\text{Var}(X) = \mathbb{E}[(X - \mathbb{E}[X])^2]$ . It measures how much  $X$  deviates from its expectation on average. The following calculation is often very useful:

$$\begin{aligned} \text{Var}(X) &= \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{E}[X^2 - 2X \cdot \mathbb{E}[X] + \mathbb{E}[X]^2] \\ &= \mathbb{E}[X^2] - 2\mathbb{E}[X]^2 + \mathbb{E}[X]^2 \\ &= \mathbb{E}[X^2] - \mathbb{E}[X]^2. \end{aligned}$$

A basic fact about variance is that the variance of the sum of independent random variables is equal to the sum of their variances.

**Proposition 3.** Let  $X_1, \dots, X_n$  be independent random variables, and let  $X = \sum_{i=1}^n X_i$ . Then  $\text{Var}(X) = \sum_{i=1}^n \text{Var}(X_i)$ .

*Proof.* For simplicity, let us define new random variables  $Y_i = X_i - \mathbb{E}[X_i]$ , and  $Y = \sum_{i=1}^n Y_i$ . Notice  $Y_1, \dots, Y_n$  are also independent, that  $Y = X - \mathbb{E}[X]$ ,  $\mathbb{E}[Y_i] = 0$  for all  $i$ ,  $\mathbb{E}[Y] = 0$ , and each of  $Y, Y_1, \dots, Y_n$  has the same variance as  $X, X_1, \dots, X_n$ , respectively. Then, it is enough to prove the proposition for  $Y$  and  $Y_1, \dots, Y_n$ . We have:

$$\begin{aligned} \text{Var}(Y) &= \mathbb{E}[Y^2] = \mathbb{E}[(Y_1 + \dots + Y_n)^2] \\ &= \sum_{i=1}^n \mathbb{E}[Y_i^2] + \sum_{i \neq j} \mathbb{E}[Y_i Y_j] \\ &= \sum_{i=1}^n \mathbb{E}[Y_i^2] + \sum_{i \neq j} \mathbb{E}[Y_i] \mathbb{E}[Y_j] \\ &= \sum_{i=1}^n \mathbb{E}[Y_i^2] = \sum_{i=1}^n \text{Var}(Y_i). \end{aligned}$$

The third line above follows from the assumption that  $Y_i$  and  $Y_j$  are independent for each  $i \neq j$ .  $\square$

The main reason variance is useful is that it gives us some control on how far a random variable can be from its expectation. Let us first recall Markov's inequality, proved in the lecture notes on Locality Sensitive Hashing.

**Theorem 4** (Markov's Inequality). Let  $Z \geq 0$  be a random variable. Then, for any  $z > 0$ ,

$$\mathbb{P}(Z > z) < \frac{\mathbb{E}[Z]}{z}.$$

We will use Markov's inequality to prove Chebyshev's inequality.

**Theorem 5** (Chebyshev's Inequality). Let  $X$  be a random variable. For any  $t > 0$ ,

$$\mathbb{P}(|X - \mathbb{E}[X]| > t) < \frac{\text{Var}(X)}{t^2}.$$

*Proof.* Define the random variable  $Z = (X - \mathbb{E}[X])^2$ . By definition,  $Z \geq 0$ . Notice that  $\text{Var}(X) = \mathbb{E}[Z]$  and that

$$|X - \mathbb{E}[X]| > t \Leftrightarrow Z > t^2.$$

Then, by Markov's inequality,

$$\mathbb{P}(|X - \mathbb{E}[X]| > t) = \mathbb{P}(Z > t^2) < \frac{\text{Var}(X)}{t^2}.$$

$\square$

Notice that, unlike Markov's inequality, Chebyshev's inequality does not need to assume that the random variable is non-negative. Moreover, Chebyshev's inequality bounds the probability that the random variable is far from its expectation in either direction.

**Exercise 5.** Give, for each value of  $t > 0$ , a random variable  $X$  so that  $\mathbb{E}[X] = 0$ ,  $\text{Var}(X) = 1$ , and such that Chebyshev's inequality holds with equality, i.e.  $\mathbb{P}(|X - \mathbb{E}[X]| \geq t) = \frac{1}{t^2}$ .

**Exercise 6.** Let  $X_1, \dots, X_n$  be  $n$  random bits, where  $\mathbb{P}(X_i = 1) = p$ . Assume that  $X_1, \dots, X_n$  are mutually independent. What is the variance of  $-X_1 + X_2 - X_3 + \dots + (-1)^n X_n$ ?

**Exercise 7.** Suppose we roll a fair die  $n$  times, and let  $X$  be the sum of the numbers that appeared over the  $n$  rolls. Use Chebyshev's inequality to bound the probability that  $|X - 3.5n| > 0.5n$ .

**Exercise 8.** Let  $A[1..n]$  be an array of distinct integers. Recall that the rank of an integer  $x$  in  $A$  equals  $r$  if and only if there are exactly  $r - 1$  integers in  $A$  that are strictly smaller than  $x$ .

Consider  $k$  indexes  $i_1, \dots, i_k$ , sampled independently and uniformly at random with replacement from  $\{1, \dots, n\}$ . Let  $\ell$  be the number of integers in  $A[i_1], \dots, A[i_k]$ , counted with repetition, that have rank  $> (\frac{1}{2} + \varepsilon)n$ . What is the expected value of  $\ell$  (as a function of  $\varepsilon$ )? What is the variance of  $\ell$ ?

Use Chebyshev's inequality to give an upper bound on the probability that the median of  $A[i_1], \dots, A[i_k]$  has rank  $> (\frac{1}{2} + \varepsilon)n$  in  $A$ .

**Exercise 9.** Given a permutation  $\sigma = \sigma_1, \dots, \sigma_n$  of  $[n]$ , let  $\text{inv}(\sigma)$  be the number of distinct pairs  $i < j$  such that  $\sigma_i > \sigma_j$ . I.e.  $\text{inv}(\sigma)$  is the number of pairs of integers that are inverted in  $\sigma$ . If  $\sigma$  is a uniformly random permutation of  $[n]$ , then what is  $\mathbb{E}[\text{inv}(\sigma)]$ ? What is  $\text{Var}(\text{inv}(\sigma))$ ?

HINT: Use the indicator random variables  $X_{ij}$ , defined for any  $1 \leq i < j \leq n$ , where  $X_{ij} = 1$  if  $\sigma_i > \sigma_j$ , and  $X_{ij} = 0$  otherwise. Compute  $\mathbb{E}[X_{ij}]$ , and use this to calculate  $\mathbb{E}[\text{inv}(\sigma)]$ . For  $\text{Var}(\text{inv}(\sigma))$ , expand  $\mathbb{E}[\text{inv}(\sigma)^2]$  and compute the quantities

1.  $\mathbb{E}[X_{ij}X_{k\ell}]$  where  $i \neq k$  and  $j \neq \ell$ ;
2.  $\mathbb{E}[X_{ij}X_{i\ell}]$  where  $j \neq \ell$ ;
3.  $\mathbb{E}[X_{ij}X_{j\ell}]$ ;
4.  $\mathbb{E}[X_{ij}^2]$ .

Use these expectations to compute  $\mathbb{E}[\text{inv}(\sigma)^2]$  and  $\text{Var}(\text{inv}(\sigma))$ .

We are now ready to analyze DISTINCT-SIMPLE.

**Theorem 6.** If  $\tilde{F}_0$  satisfies (1), then with probability at least  $1/2$ , the estimate  $\hat{F}_0$  output by  $\text{DISTINCT-SIMPLE}(\sigma, \tilde{F}_0, k)$  satisfies

$$\left(1 - \frac{\sqrt{8}}{\sqrt{k}}\right)F_0 \leq \hat{F}_0 \leq \left(1 + \frac{\sqrt{8}}{\sqrt{k}}\right)F_0.$$

Moreover, the algorithm uses at most  $O(k)$  words of memory in expectation.

*Proof.* Let  $D$  be the set of those  $i \in [n]$  that appear in the stream, i.e.  $D = \{i : f_i > 0\}$ . By this definition,  $F_0 = |D|$ . For any  $i \in D$ , let  $X_i$  be the indicator random variable which equals 1

if  $i \in S$ , and 0 otherwise. We have  $\mathbb{E}[X_i] = \mathbb{P}(i \in S) = 2^{-d}$  because  $h(i)$  is a uniformly random string in  $\{0, 1\}^L$  and exactly  $2^{L-d}$  out of all the  $2^L$  such strings have their  $d$  leftmost bits set to 0. It follows that

$$\mathbb{E}[|S|] = \mathbb{E}\left[\sum_{i \in D} X_i\right] = \sum_{i \in D} \mathbb{P}(i \in S) = 2^{-d} F_0.$$

Since  $F_0 \leq \tilde{F}_0$ , and we chose  $d$  so that  $2^{-d} \tilde{F}_0 \leq k$ , we have  $\mathbb{E}[|S|] = 2^{-d} F_0 \leq 2^{-d} \tilde{F}_0 \leq k$ . Since storing  $S$  dominates the space usage of the algorithm, it follows that the algorithm uses at most  $O(k)$  words of memory in expectation. It remains to analyze how close  $\hat{F}_0$  is to  $F_0$ .

It is easy to see that  $\hat{F}_0$  equals  $F_0$  in expectation:  $\mathbb{E}[\hat{F}_0] = \mathbb{E}[2^d |S|] = F_0$ . However, this is not enough to show that  $\hat{F}_0$  is close to  $F_0$ : a random variable could be very far from its expectation with high probability. This is where Chebyshev's inequality comes to the rescue: if we know that the variance of  $\hat{F}_0$  is small, then we know it is unlikely to be too far from its expectation.

Recall that we assumed  $h(1), \dots, h(n)$  are all independent, and, therefore, the random variables  $X_i$  defined for each  $i \in D$  are also independent. By Proposition 3, it follows that

$$\text{Var}(|S|) = \text{Var}\left(\sum_{i \in D} X_i\right) = \sum_{i \in D} \text{Var}(X_i).$$

For any  $i \in D$ ,

$$\text{Var}(X_i) = \mathbb{E}[X_i^2] - \mathbb{E}[X_i]^2 \leq \mathbb{E}[X_i^2] = \mathbb{E}[X_i] = 2^{-d}.$$

Here, we used the fact that  $\mathbb{E}[X_i^2] \geq 0$ , and that  $X_i^2 = X_i$  because  $X_i \in \{0, 1\}$ . Plugging this into the equation for  $\text{Var}(|S|)$ , we get that  $\text{Var}(|S|) \leq 2^{-d} F_0 = \mathbb{E}[|S|]$ .

We are ready to finish the proof. Let  $\varepsilon = \frac{\sqrt{8}}{\sqrt{k}}$ . By Chebyshev's inequality,

$$\begin{aligned} \mathbb{P}(|\hat{F}_0 - F_0| \geq \varepsilon F_0) &= \mathbb{P}(|\hat{F}_0 - \mathbb{E}[\hat{F}_0]| \geq \varepsilon \mathbb{E}[\hat{F}_0]) \\ &= \mathbb{P}(|S| - \mathbb{E}[|S|] \geq \varepsilon \mathbb{E}[|S|]) \\ &\leq \frac{\text{Var}(|S|)}{\varepsilon^2 \mathbb{E}[|S|]^2} \leq \frac{1}{\varepsilon^2 \mathbb{E}[|S|]}. \end{aligned}$$

I.e. if the expected size of  $S$  is not too small, we have a large probability of an accurate estimate. (This makes intuitive sense: a large sample gives a better estimate.) By the choice of  $d$  we know that  $2^d \leq 2\tilde{F}_0/k$ , which, after re-arranging, gives  $2^{-d} \tilde{F}_0 \geq k/2$ . Therefore, by (1),

$$\mathbb{E}[|S|] = 2^{-d} F_0 \geq 2^{-d-1} \tilde{F}_0 \geq \frac{k}{4}.$$

Substituting into the expression we got from Chebyshev's inequality, we have

$$\mathbb{P}(|\hat{F}_0 - F_0| \geq \varepsilon F_0) \leq \frac{4}{\varepsilon^2 k} = \frac{1}{2},$$

as we wanted. □

**Exercise 10.** Use Exercise 8 to analyze the following algorithm: given input array  $A[1..n]$  of distinct integers, the algorithm samples  $k$  indexes  $i_1, \dots, i_k$  independently and uniformly at random with replacement from  $\{1, \dots, n\}$ . Then it computes and outputs the median  $X$  of  $A[i_1], \dots, A[i_k]$ . Show that if  $k \geq \frac{4}{\varepsilon^2}$ , then the output  $X$  of the of the algorithm has rank between  $(\frac{1}{2} - \varepsilon)n$  and  $(\frac{1}{2} + \varepsilon)n$  in  $A$  with probability at least  $\frac{1}{2}$ .



### 3.2 Adaptive Sampling

Next we describe a single pass streaming algorithm for distinct counts which does not assume we know an estimate  $\hat{F}_0$ . We describe one of the simpler solutions for this problem, known as Adaptive Sampling, as the procedure DISTINCT.

```

DISTINCT( $\sigma, k$ )
1   $S = \emptyset$ 
2   $d = 0$ 
3   $L = \lceil \log_2 n \rceil$ 
4  Pick a hash function  $h : [n] \rightarrow \{0, 1\}^L$ 
5  for  $t = 1$  to  $m$ 
6      if  $h(i_t) \in 0^d \{0, 1\}^{L-d}$  and  $i_t \notin S$ 
7           $S = S \cup \{i_t\}$ 
8      while  $|S| > k$ 
9           $d = d + 1$ 
10          $T = \emptyset$ 
11         for  $j \in S$ 
12             if  $h(j) \in 0^d \{0, 1\}^{L-d}$ 
13                  $T = T \cup \{j\}$ 
14              $S = T$ 
15 return  $\hat{F}_0 = 2^d \cdot |S|$ 

```

This algorithm is a variant of DISTINCT-SIMPLE which sets the value  $d$ , and, consequently, the sampling rate, adaptively. The algorithm keeps the invariant that  $|S| \leq k$ , and, whenever  $S$  exceeds size  $k$ , we increase  $d$  by 1, and drop, in expectation, half of its elements. At the end, we expect  $S$  to hold about  $2^{-d}F_0$  elements, so we output  $2^d|S|$  as our estimate of  $F_0$ . Notice that the algorithm is guaranteed to use  $O(k)$  words of memory.

We are now ready to analyze DISTINCT. Warning: this is one of the more involved probabilistic analyses you will see in this course, and it is somewhat heavy in calculations.

**Theorem 7.** *Let  $k \geq 144$ . With probability at least  $1/2$ , the estimate  $\hat{F}_0$  output by DISTINCT( $\sigma, k$ ) satisfies*

$$\left(1 - \frac{4}{\sqrt{k}}\right)F_0 \leq \hat{F}_0 \leq \left(1 + \frac{4}{\sqrt{k}}\right)F_0.$$

*Proof.* Let  $D$  be the set of those  $i \in [n]$  that appear in the stream, i.e.  $D = \{i : f_i > 0\}$ . By this definition,  $F_0 = |D|$ . Let us define the sets  $S_0, S_1, \dots, S_L$  by  $S_\ell = \{i \in D : h(i) \in 0^\ell \{0, 1\}^{L-\ell}\}$ , and observe that  $S_L \subseteq S_{L-1} \subseteq \dots \subseteq S_0 = D$ . It is easy to show that the final value of  $d$  computed by DISTINCT equals  $\min\{\ell : |S_\ell| \leq k\}$ , and the estimate output by the algorithm equals  $2^d|S_d|$ .

Using a calculation analogous to the one in Theorem 6, we get

$$\begin{aligned} \mathbb{E}[|S_\ell|] &= 2^{-\ell}F_0, \\ \text{Var}(|S_\ell|) &\leq 2^{-\ell}F_0 = \mathbb{E}[|S_\ell|]. \end{aligned}$$

Let  $\varepsilon = \frac{4}{\sqrt{k}} \leq \frac{1}{3}$ , and let  $a$  be the largest integer such that  $(1 - \varepsilon)2^{-a}F_0 > k$ , and let  $b$  be the smallest integer such that  $(1 + \varepsilon)2^{-b}F_0 \leq k$ . Verify that  $b \geq a$ , and  $b - a \leq 1 + \log_2 \frac{1+\varepsilon}{1-\varepsilon} \leq 2$ , i.e.  $b \in \{a, a + 1, a + 2\}$ .

Since  $\mathbb{E}[|S_a|] = 2^{-a}F_0$ , we have, by Chebyshev's inequality

$$\begin{aligned} \mathbb{P}(|S_a| \leq k) &\leq \mathbb{P}(|S_a| < (1 - \varepsilon)2^{-a}F_0) \\ &= \mathbb{P}(\mathbb{E}[|S_a|] - |S_a| > \varepsilon\mathbb{E}[|S_a|]) \\ &< \frac{\text{Var}(|S_a|)}{\varepsilon^2\mathbb{E}[|S_a|]^2} \leq \frac{1}{\varepsilon^2\mathbb{E}[|S_a|]}. \end{aligned}$$

Since  $\mathbb{E}[|S_a|] > k$ , and  $\varepsilon^2 = \frac{16}{k}$ , the right hand side is at most  $1/16$ . By an analogous calculation:

$$\mathbb{P}(|S_b| - \mathbb{E}[|S_b|]| > \varepsilon 2^{-b}F_0) = \mathbb{P}(|S_b| - \mathbb{E}[|S_b|]| > \varepsilon\mathbb{E}[|S_b|]) < \frac{1}{\varepsilon^2\mathbb{E}[|S_b|]}.$$

By the choice of  $b$ ,  $(1 + \varepsilon)2^{-b}F_0 \geq k/2$  (or we would've chosen a smaller  $b$ ), so  $\mathbb{E}[|S_b|] = 2^{-b}F_0 \geq 3k/8$ . Since  $\varepsilon^2 = \frac{16}{k}$ , the right hand side of the inequality above is at most  $1/6$ . Another analogous calculation shows that, if  $b = a + 2$ , then

$$\mathbb{P}(|S_{b-1}| - \mathbb{E}[|S_{b-1}|]| > \varepsilon 2^{-b+1}F_0) < \frac{1}{6}.$$

By the union bound, with probability at least  $1 - 1/6 - 1/6 - 1/16 > 1/2$ , we have

$$\begin{aligned} |S_a| &> k, \\ ||S_{b-1}| - \mathbb{E}[|S_{b-1}|]| &\leq \varepsilon 2^{-b+1}F_0, \\ ||S_b| - \mathbb{E}[|S_b|]| &\leq \varepsilon 2^{-b}F_0. \end{aligned}$$

The last inequality and the choice of  $b$  imply that  $|S_b| \leq k$ . Therefore, at the end of the execution of the algorithm,  $d = b$  or  $d = b - 1$ , and the algorithm outputs the estimate  $\hat{F}_0 \in \{2^{b-1}|S_{b-1}|, 2^b|S_b|\}$ . If  $\hat{F}_0 = 2^{b-1}|S_{b-1}|$ , we have

$$|\hat{F}_0 - F_0| = 2^{b-1}||S_{b-1}| - \mathbb{E}[|S_{b-1}|]| \leq \varepsilon F_0.$$

Analogously, if  $\hat{F}_0 = 2^b|S_b|$ , we have

$$|\hat{F}_0 - F_0| = 2^b||S_b| - \mathbb{E}[|S_b|]| \leq \varepsilon F_0.$$

This completes the proof of the theorem. □

The constants in the proof are not chosen to be the tightest possible, but rather to make the calculations relatively painless. Another way to state the result we proved is that, using `DISTINCT`, we can approximate  $F_0$  up to a factor  $1 \pm \varepsilon$  using  $O(\frac{1}{\varepsilon^2})$  words of memory, or, equivalently,  $O(\frac{\log n}{\varepsilon^2})$  bits of memory. More sophisticated algorithms are known which use  $O(\frac{1}{\varepsilon^2} + \log n)$  bits of memory; this latter bound is the best possible in the worst case.

**Exercise 11.** *The theorem above does not show any guarantee when  $k = 1$ . Prove that there exists a constant  $C$ , independent of  $n$ ,  $m$  or the stream, so that if  $\hat{F}_0$  is the estimate output by `DISTINCT`( $\sigma, 1$ ), then, with probability at least  $1/2$ ,*

$$\frac{1}{C}F_0 \leq \hat{F}_0 \leq CF_0$$

Let us finally make a remark about making the assumptions on the hash function in the algorithm more realistic. Carefully checking the calculation in Proposition 3 shows that it is enough if the random variables  $X_1, \dots, X_n$  are *pairwise independent*, i.e. if for every  $i \neq j$  and every two values  $x, x'$ , we have  $\mathbb{P}(X_i = x, X_j = x') = \mathbb{P}(X_i = x)\mathbb{P}(X_j = x')$ . (This is not the same as full independence: take for example  $X_1$  and  $X_2$  to be independent and uniform in  $\{0, 1\}$ , and take  $X_3$  to be the XOR of  $X_1$  and  $X_2$ .)

This observation inspires the following definition, which is closely related to the definition of a universal hashing family.

**Definition 8.** *A family  $\mathcal{H}$  of functions from  $[n]$  to  $\{0, 1\}^L$  is a pairwise independent hash family if for all  $i, j \in [n]$ ,  $i \neq j$ , and for any two bit-strings  $x, x' \in \{0, 1\}^L$ , we have*

$$\mathbb{P}(h(i) = x, h(j) = x') = \frac{1}{2^{2L}},$$

where the probability is taken over picking a uniformly random  $h \in \mathcal{H}$ .

A pairwise independent hash family is also a universal family, but not necessarily the other way around. Using a little bit of algebra, it is easy to construct pairwise independent hash family of size  $2^{2L} = O(n^2)$ . Picking a random function from this family requires only picking  $2L = O(\log n)$  random bits, and, moreover, the value of the hash function on any element of  $[n]$  can be quickly computed from the random bits. So, in DISTINCT we will take  $h$  to be a random function from such a pairwise independent hash family. This will only increase our space complexity by a constant number of words, and will not affect the correctness guarantees of the algorithm, because in the analysis we only used pairwise independence.

## 4 Parallel Algorithms

### 4.1 The Model

In this section we explore another computational model which addresses the challenges of processing big data. We still consider a setting in which the data we want to process is too large to fit on a machine. However, we are going to assume that we possess a cluster of many machines: enough so that the total space on them fits all our data, and potentially slightly more than that. For concreteness, let us say that the input data is a set of  $n$  integers, we have  $\sqrt{n}$  machines, and each of them can keep  $O(\sqrt{n})$  integers in its local storage. Assume that the input starts out partitioned arbitrarily among the machines. Each machine can compute on the part of the input in its local storage. Since no machine sees the entire input, we cannot hope to solve even very simple problems, unless we allow the machines to exchange messages. However, communication between the machines is usually slow: sending data over a network is much slower than internal computation. For this reason, we will try to design algorithms that process all the data with as little communication between the machines as possible.

Let us describe the model more formally. Assume that the size of the input is  $n$  (measured, for example, in words of memory). We have  $m$  machines, and each of them has local storage of  $s$  words of memory. We assume that  $ms = \Omega(n)$ , and sometimes even allow for  $ms = \Omega(n^{1+\varepsilon})$  for some relatively small constant  $\varepsilon > 0$ . The input starts out partitioned arbitrarily between the

$m$  machines. The computation proceeds in *rounds*. In a single round, each machine can execute any polynomial time algorithm on its local storage. Once all machines have finished their local computation, they are allowed to exchange messages. All machines simultaneously send messages to the other machines (assume the machines are numbered). The total size of the messages sent or received by a machine must not exceed  $s$ . Once the messages are exchanged, the round is complete. In the next round the local storage of each machine contains the union of the messages it received in the previous round together with the contents of its local storage from the previous round. We have the constraint that the local storage of any machine must not exceed  $s$  at any time. To satisfy this constraint, a machine can discard information from its local storage during the algorithm. Our main goal is to complete the entire computation in as few rounds as possible, preferably a constant number of rounds.

This model is an abstraction and simplification of real-world systems like MapReduce and Hadoop. The point here is not to perfectly model the actual systems, just like the point of the RAM model is not to be a perfect model of a physical computer. Our goal instead is to have a model which is simple enough to allow us to analyze our algorithms, and yet captures some of the algorithmic challenges of designing algorithms for parallel systems. We should note that there are other theoretical models of parallel computation, like the various flavors of the PRAM, and the BSP model of Valiant. The model above can be seen as a simplification of BSP.

## 4.2 Simple Algorithms

First a warm-up exercise.

**Exercise 12.** *Give a parallel algorithm which finds the maximum and the sum of  $n$  integers in a constant number of rounds when  $s = m = \Theta(\sqrt{n})$ . At the end, one specially designated machine must have the output written in its local memory.*

Now let us look at something a bit more interesting: the prefix sums problem. Our input is  $n$  pairs of integers,  $(1, x_1), \dots, (n, x_n)$ , and our goal is to compute, for all  $1 \leq i \leq n$ , the sum  $x_1 + \dots + x_i$ . We assume we have  $s = m = \Theta(\sqrt{n})$ . Notice that no machine has enough space for the entire output, so we will be satisfied having the different pieces of the output stored on different machines, as long as some prefix sum is stored by some machine.

A standard algorithm is as follows:

1. Exchange messages so that  $(1, x_1), \dots, (\sqrt{n}, x_{\sqrt{n}})$  are stored on the first machine,  $(\sqrt{n} + 1, x_{\sqrt{n}+1}), \dots, (2\sqrt{n}, x_{2\sqrt{n}})$  on the second machine, etc.
2. In the next round, the  $i$ -th machine computes  $x_{(i-1)\sqrt{n}+1} + \dots + x_j$  for all  $(i-1)\sqrt{n} + 1 \leq j \leq i\sqrt{n}$ . Moreover, the  $i$ -th machine sends to all machines numbered higher than itself the sum of all integers in its local memory.
3. In the next and final round, the local storage of the  $i$ -th machine contains the sums of the integers on machines  $1, \dots, i-1$ . It can add up these sums to get  $x_1 + \dots + x_{(i-1)\sqrt{n}}$ . The local storage of the  $i$ -th machine also contains the sums  $x_{(i-1)\sqrt{n}+1} + \dots + x_j$  for all  $(i-1)\sqrt{n} + 1 \leq j \leq i\sqrt{n}$  from the previous round. From these, the  $i$ -th machine can compute the prefix sums  $x_1 + \dots + x_j$  for all  $(i-1)\sqrt{n} + 1 \leq j \leq i\sqrt{n}$ .

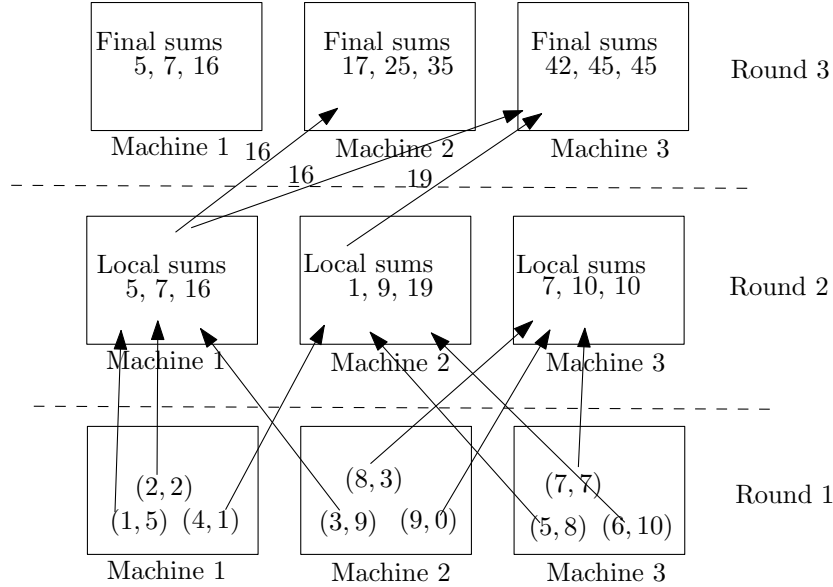


Figure 1: An example of the prefix sums algorithm.

An illustrative example is shown in Figure 1.

**Exercise 13.** Suppose the input is a weighted graph  $G$  on  $n$  nodes and  $\Theta(n^2)$  edges. Assume that  $s = \Theta(n^{3/2})$  and  $m = \Theta(n)$ . Initially, each machine has  $n$  edges of the graph in its local storage, where each edge is represented as a pair of vertices with a weight. Give a parallel algorithm that computes the minimum spanning tree of  $G$  in a constant number of rounds.

### 4.3 Using Streaming Algorithms in Parallel

Many streaming algorithms can be turned into parallel algorithms without much effort. For example, suppose that our input is a sequence  $\sigma = (\sigma_1, \dots, \sigma_n)$  of  $n$  integers, which could possibly repeat. Each of our  $m = \Theta(\sqrt{n})$  machines has  $\Theta(\sqrt{n})$  integers from  $\sigma$ . We want to estimate the number of distinct integers in  $\sigma$ .

We show how to use the adaptive sampling algorithm DISTINCT to achieve this goal. First one machine samples a random hash function and sends it to all other machines: this way all machines use the same hash function. **(This can be done in two rounds: see how?)** Then, we run DISTINCT on each machine with parameter  $k = \Theta(\frac{1}{\varepsilon^2})$ , and once machine  $i$  has processed its input, it can send the set of at most  $k$  elements  $S^{(i)}$  computed by DISTINCT on its local input, together with the final value  $d^{(i)}$  of  $d$  used to compute  $S^{(i)}$ , to a specially designated machine, say the first one. (To make sure the notation is clear, we point out that  $S^{(i)}$  is the set of integers on the  $i$ -th machine whose hash value given by  $h$  starts with  $d^{(i)}$  0's, and  $d^{(i)}$  is chosen so that  $|S^{(i)}| \leq k$ .) As long as  $\varepsilon$  is a constant, the total number of words sent to the first machine is  $O(\varepsilon^{-2}m) = O(\sqrt{n})$ , and does not exceed the storage requirement.

It remains for the first machine to finish the computation. We claim that it can exactly simulate the DISTINCT algorithm when run on the entire sequence  $\sigma$ . Recall from the proof of Theorem 7 the notation  $D$  for the set of distinct integers in  $\sigma$ , and the sets  $S_\ell = \{i : h(i) \in 0^\ell\{0, 1\}^{L-\ell}\}$ .

Recall further that the final value of  $d$  computed by DISTINCT equals  $\min\{\ell : |S_\ell| \leq k\}$ , and the estimate output by the algorithm equals  $2^d |S_d|$ . Observe that for any  $\ell \geq \max\{d^{(1)}, \dots, d^{(m)}\}$ ,  $S_\ell \subseteq S^{(1)} \cup \dots \cup S^{(m)}$ , and, moreover, that  $d \geq \max\{d^{(1)}, \dots, d^{(m)}\}$ . So, the first machine can just go over the different possible values for  $d$ , from  $\max\{d^{(1)}, \dots, d^{(m)}\}$  to  $L$ , and check for each one whether  $|S_d| \leq k$ . The smallest  $d$  for which this is true gives us our desired output  $2^d |S_d|$ . Since we already proved in Theorem 7 that DISTINCT computed an estimate which is within a  $(1 \pm \varepsilon)$  multiplicative factor from the true number of distinct elements, we now also have a parallel algorithm with the same guarantee.

**Exercise 14.** *Show how to use a constant number of additional rounds of computation to refine the approximation factor above from  $(1 \pm \varepsilon)$  to  $(1 \pm \frac{C}{\sqrt{s}})$ , for a constant  $C$ .*