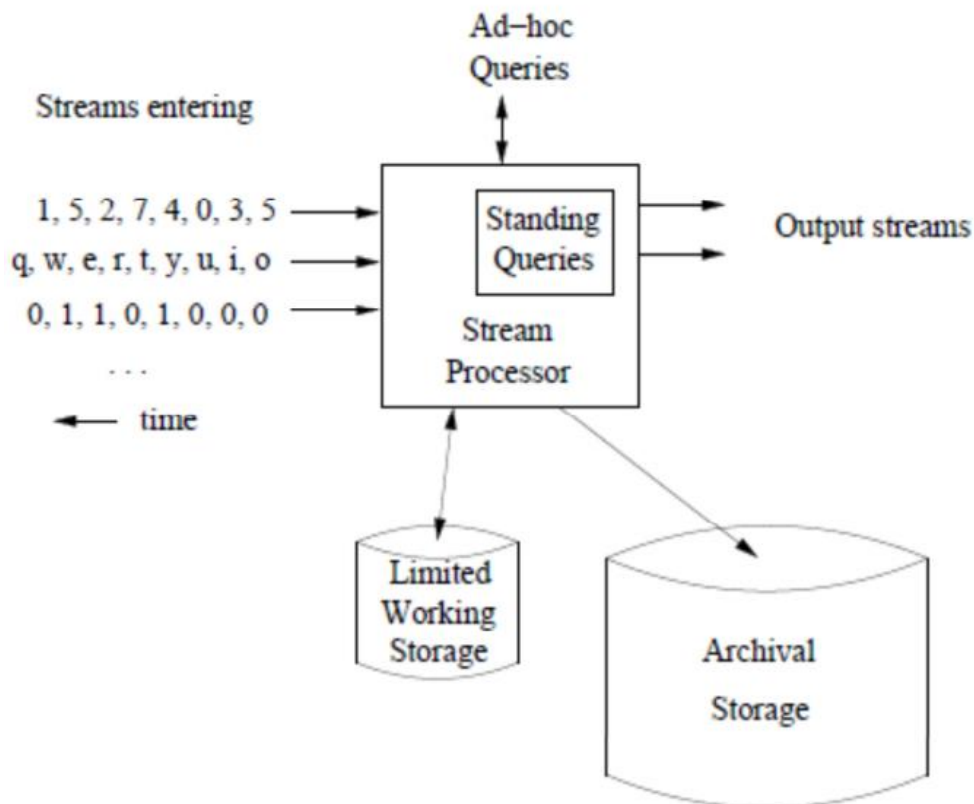


# STREAM COMPUTING

## Streaming of Data

The algorithms for processing streams each involve summarization of the stream in some way. We shall start by considering how to make a useful sample of a stream and how to filter a stream to eliminate most of the “undesirable” elements. We then show how to estimate the number of different elements in a stream using much less storage than would be required if we listed all the elements we have seen. Another approach to summarizing a stream is to look at only a fixed-length “window” consisting of the last  $n$  elements for some (typically large)  $n$ . We then query the window as if it were a relation in a database. If there are many streams and/or  $n$  is large, we may not be able to store the entire window for every stream, so we need to summarize even the windows. We address the fundamental problem of maintaining an approximate count on the number of 1’s in the window of a bit stream, while using much less space than would be needed to store the entire window itself. This technique generalizes to approximating various kinds of sums.

## Data Stream Management System



In analogy to a database-management system, we can view a stream processor as a kind of data-management system, the high-level organization of which is suggested in the above figure. Any number of streams can enter the system. Each stream can provide elements at its own schedule; they need not have the same data rates or data types, and the time between elements of one stream need not be uniform. The fact that the rate of arrival of stream elements is not under the control of the system distinguishes stream processing from the processing of data that goes on within a database-management system. The latter system controls the rate at which data is read from the disk, and therefore never has to worry about data getting lost as it attempts to execute queries. Streams may be archived in a large archival store, but we assume it is not possible to answer queries from the archival store. It could be examined only under special circumstances using time-consuming retrieval processes. There is also a working store, into which summaries or parts of streams may be placed, and which can be used for answering queries. The working store might be disk, or it might be main memory, depending on how fast we need to process queries. But either way, it is of sufficiently limited capacity that it cannot store all the data from all the streams.

## **Sources of Streaming Data**

### **Sensor Data**

Imagine a temperature sensor bobbing about in the ocean, sending back to a base station a reading of the surface temperature each hour. The data produced by this sensor is a stream of real numbers. It is not a very interesting stream, since the data rate is so low. It would not stress modern technology, and the entire stream could be kept in main memory, essentially forever. Now, give the sensor a GPS unit, and let it report surface height instead of temperature. The surface height varies quite rapidly compared with temperature, so we might have the sensor send back a reading every tenth of a second. If it sends a 4-byte real number each time, then it produces 3.5 megabytes per day. It will still take some time to fill up main memory, let alone a single disk. But one sensor might not be that interesting. To learn something about ocean behavior, we might want to deploy a million sensors, each sending back a stream, at the rate of ten per second. A million sensors isn't very many; there would be one for every 150 square miles of ocean. Now we have 3.5 terabytes arriving every day, and we definitely need to think about what can be kept in working storage and what can only be archived.

### **Image Data**

Satellites often send down to earth streams consisting of many terabytes of images per day. Surveillance cameras produce images with lower resolution than satellites, but there can be many of them, each producing a stream of images at intervals like one second. London is said to have six million such cameras, each producing a stream.

## Internet and Web Traffic

A switching node in the middle of the Internet receives streams of IP packets from many inputs and routes them to its outputs. Normally, the job of the switch is to transmit data and not to retain it or query it. But there is a tendency to put more capability into the switch, e.g., the ability to detect denial-of-service attacks or the ability to reroute packets based on information about congestion in the network.

Web sites receive streams of various types. For example, Google receives several hundred million search queries per day. Yahoo! accepts billions of “clicks” per day on its various sites. Many interesting things can be learned from these streams. For example, an increase in queries like “sore throat” enables us to track the spread of viruses. A sudden increase in the click rate for a link could indicate some news connected to that page, or it could mean that the link is broken and needs to be repaired.

## Stream Queries

There are two ways that queries get asked about streams. A place within the processor where standing queries are stored. These queries are, in a sense, permanently executing, and produce outputs at appropriate times.

### Example 1:

The stream produced by the ocean-surface-temperature sensor mentioned at the beginning of Section 4.1.2 might have a standing query to output an alert whenever the temperature exceeds 25 degrees centigrade. This query is easily answered, since it depends only on the most recent stream element. Alternatively, we might have a standing query that, each time a new reading arrives, produces the average of the 24 most recent readings. That query also can be answered easily, if we store the 24 most recent stream elements. When a new stream element arrives, we can drop from the working store the 25th most recent element, since it will never again be needed (unless there is some other standing query that requires it).

Another query we might ask is the maximum temperature ever recorded by that sensor. We can answer this query by retaining a simple summary: the maximum of all stream elements ever seen. It is not necessary to record the entire stream. When a new stream element arrives, we compare it with the stored maximum, and set the maximum to whichever is larger. We can then answer the query by producing the current value of the maximum. Similarly, if we want the average temperature over all time, we have only to record two values: the number of readings ever sent in the stream and the sum of those readings. We can adjust these values easily each time a new reading arrives, and we can produce their quotient as the answer to the query. 2

The other form of query is ad-hoc, a question asked once about the current state of a stream or streams. If we do not store all streams in their entirety, as normally we can not, then we cannot expect to answer arbitrary queries about streams. If we have some idea what kind of queries will be asked through the ad-hoc query interface, then we can prepare for them by storing appropriate parts or summaries of streams as in Example.

If we want the facility to ask a wide variety of ad-hoc queries, a common approach is to store a sliding window of each stream in the working store. A sliding window can be the most recent  $n$  elements of a stream, for some  $n$ , or it can be all the elements that arrived within the last  $t$  time units, e.g., one day. If we regard each stream element as a tuple, we can treat the window as a relation and query it with any SQL query. Of course the stream-management system must keep the window fresh, deleting the oldest elements as new ones come in.

### **Example 2:**

Web sites often like to report the number of unique users over the past month. If we think of each login as a stream element, we can maintain a window that is all logins in the most recent month. We must associate the arrival time with each login, so we know when it no longer belongs to the window. If we think of the window as a relation  $\text{Logins}(\text{name}, \text{time})$ , then it is simple to get the number of unique users over the past month. The SQL query is:

```
SELECT COUNT(DISTINCT(name))
FROM Logins
WHERE time >= t;
```

Here,  $t$  is a constant that represents the time one month before the current time. Note that we must be able to maintain the entire stream of logins for the past month in working storage. However, for even the largest sites, that data is not more than a few terabytes, and so surely can be stored on disk. 2

### **Issues in Stream Processing**

Before proceeding to discuss algorithms, let us consider the constraints under which we work when dealing with streams. First, streams often deliver elements very rapidly. We must process elements in real time, or we lose the opportunity to process them at all, without accessing the archival storage. Thus, it often is important that the stream-processing algorithm is executed in main memory, without access to secondary storage or with only rare accesses to secondary storage. Moreover, even when streams are “slow,” as in the sensor-data example, there may be many such streams. Even if each stream by itself can be processed using a small amount of main memory, the requirements of all the streams together can easily exceed the amount of available main memory.

## Sampling Data in a Stream

As our first example of managing streaming data, we shall look at extracting reliable samples from a stream. As with many stream algorithms, the “trick” involves using hashing in a somewhat unusual way

### A Motivating Example

The general problem we shall address is selecting a subset of a stream so that we can ask queries about the selected subset and have the answers be statistically representative of the stream as a whole. If we know what queries are to be asked, then there are a number of methods that might work, but we are looking for a technique that will allow ad-hoc queries on the sample. We shall look at a particular problem, from which the general idea will emerge. Our running example is the following. A search engine receives a stream of queries, and it would like to study the behavior of typical users.<sup>1</sup> We assume the stream consists of tuples (user, query, time). Suppose that we want to answer queries such as “What fraction of the typical user’s queries were repeated over the past month?” Assume also that we wish to store only 1/10th of the stream elements.

The obvious approach would be to generate a random number, say an integer from 0 to 9, in response to each search query. Store the tuple if and only if the random number is 0. If we do so, each user has, on average, 1/10th of their queries stored. Statistical fluctuations will introduce some noise into the data, but if users issue many queries, the law of large numbers will assure us that most users will have a fraction quite close to 1/10th of their queries stored. However, this scheme gives us the wrong answer to the query asking for the average number of duplicate queries for a user. Suppose a user has issued  $s$  search queries one time in the past month,  $d$  search queries twice, and no search queries more than twice.

If we have a 1/10th sample, of queries, we shall see in the sample for that user an expected  $s/10$  of the search queries issued once. Of the  $d$  search queries issued twice, only  $d/100$  will appear twice in the sample; that fraction is  $d$  times the probability that both occurrences of the query will be in the 1/10th sample. Of the queries that appear twice in the full stream,  $18d/100$  will appear exactly once. To see why, note that  $18/100$  is the probability that one of the two occurrences will be in the 1/10th of the stream that is selected, while the other is in the 9/10th that is not selected.

The correct answer to the query about the fraction of repeated searches is  $d/(s+d)$ . However, the answer we shall obtain from the sample is  $d/(10s+19d)$ . To derive the latter formula, note that  $d/100$  appear twice, while  $s/10+18d/100$  appear once. Thus, the fraction appearing twice in the sample is  $d/100$  divided 1While we shall refer to “users,” the search engine really receives IP addresses from which the search query was issued. We shall assume that these IP addresses identify unique users, which is approximately true, but not exactly true. by  $d/100+ s/10 + 18d/100$ . This ratio is  $d/(10s+ 19d)$ . For no positive values of  $s$  and  $d$  is  $d/(s + d) = d/(10s + 19d)$ .

### Obtaining a Representative Sample

The query of example 1, like many queries about the statistics of typical users, cannot be answered by taking a sample of each user's search queries. Thus, we must strive to pick 1/10th of the users, and take all their searches for the sample, while taking none of the searches from other users. If we can store a list of all users, and whether or not they are in the sample, then we could do the following. Each time a search query arrives in the stream, we look up the user to see whether or not they are in the sample. If so, we add this search query to the sample, and if not, then not. However, if we have no record of ever having seen this user before, then we generate a random integer between 0 and 9. If the number is 0, we add this user to our list with value "in," and if the number is other than 0, we add the user with the value "out."

That method works as long as we can afford to keep the list of all users and their in/out decision in main memory, because there isn't time to go to disk for every search that arrives. By using a hash function, one can avoid keeping

the list of users. That is, we hash each user name to one of ten buckets, 0 through 9. If the user hashes to bucket 0, then accept this search query for the sample, and if not, then not.

Note we do not actually store the user in the bucket; in fact, there is no data in the buckets at all. Effectively, we use the hash function as a random number generator, with the important property that, when applied to the same user several times, we always get the same "random" number. That is, without storing the in/out decision for any user, we can reconstruct that decision any time a search query by that user arrives. More generally, we can obtain a sample consisting of any rational fraction  $a/b$  of the users by hashing user names to  $b$  buckets, 0 through  $b - 1$ . Add the search query to the sample if the hash value is less than  $a$ .

### The General Sampling Problem

The running example is typical of the following general problem. Our stream consists of tuples with  $n$  components. A subset of the components are the key components, on which the selection of the sample will be based. In our running example, there are three components – user, query, and time – of which only user is in the key. However, we could also take a sample of queries by making query be the key, or even take a sample of user-query pairs by making both those components form the key.

To take a sample of size  $a/b$ , we hash the key value for each tuple to  $b$  buckets, and accept the tuple for the sample if the hash value is less than  $a$ . If the key consists of more than one component, the hash function needs to combine the values for those components to make a single hash-value. The result will be a sample consisting of all tuples with certain key values. The selected key values will be approximately  $a/b$  of all the key values appearing in the stream.

### **Varying the Sample Size**

Often, the sample will grow as more of the stream enters the system. In our running example, we retain all the search queries of the selected 1/10th of the users, forever. As time goes on, more searches for the same users will be accumulated, and new users that are selected for the sample will appear in the stream. If we have a budget for how many tuples from the stream can be stored as the sample, then the fraction of key values must vary, lowering as time goes on. In order to assure that at all times, the sample consists of all tuples from a subset of the key values, we choose a hash function  $h$  from key values to a very large number of values  $0, 1, \dots, B-1$ . We maintain a threshold  $t$ , which initially can be the largest bucket number,  $B - 1$ . At all times, the sample consists of those tuples whose key  $K$  satisfies  $h(K) \leq t$ . New tuples from the stream are added to the sample if and only if they satisfy the same condition.

If the number of stored tuples of the sample exceeds the allotted space, we lower  $t$  to  $t-1$  and remove from the sample all those tuples whose key  $K$  hashes to  $t$ . For efficiency, we can lower  $t$  by more than 1, and remove the tuples with several of the highest hash values, whenever we need to throw some key values out of the sample. Further efficiency is obtained by maintaining an index on the hash value, so we can find all those tuples whose keys hash to a particular value quickly.

### **Filtering Streams**

Another common process on streams is selection, or filtering. We want to accept those tuples in the stream that meet a criterion. Accepted tuples are passed to another process as a stream, while other tuples are dropped. If the selection criterion is a property of the tuple that can be calculated (e.g., the first component is less than 10), then the selection is easy to do. The problem becomes harder when the criterion involves lookup for membership in a set. It is especially hard, when that set is too large to store in main memory. In this section, we shall discuss the technique known as “Bloom filtering” as a way to eliminate most of the tuples that do not meet the criterion.

### **Motivating Example**

Again let us start with a running example that illustrates the problem and what we can do about it. Suppose we have a set  $S$  of one billion allowed email addresses – those that we will allow through because we believe them not to be spam. The stream consists of pairs: an email address and the email itself. Since the typical email address is 20 bytes or more, it is not reasonable to store  $S$  in main memory. Thus, we can either use disk accesses to determine whether or not to let through any given stream element, or we can devise a method that requires no more main memory than we have available, and yet will filter most of the undesired stream elements.

Suppose for argument's sake that we have one gigabyte of available main memory. In the technique known as Bloom filtering, we use that main memory as a bit array. In this case, we have room for eight billion bits, since one byte equals eight bits. Devise a hash function  $h$  from email addresses to eight billion buckets. Hash each member of  $S$  to a bit, and set that bit to 1. All other bits of the array remain 0.

Since there are one billion members of  $S$ , approximately  $1/8$ th of the bits will be 1. The exact fraction of bits set to 1 will be slightly less than  $1/8$ th, because it is possible that two members of  $S$  hash to the same bit. We shall discuss the exact fraction of 1's in Section 4.3.3. When a stream element arrives, we hash its email address. If the bit to which that email address hashes is 1, then we let the email through. But if the email address hashes to a 0, we are certain that the address is not in  $S$ , so we can drop this stream element.

Unfortunately, some spam email will get through. Approximately  $1/8$ th of the stream elements whose email address is not in  $S$  will happen to hash to a bit whose value is 1 and will be let through. Nevertheless, since the majority of emails are spam (about 80% according to some reports), eliminating  $7/8$ th of the spam is a significant benefit. Moreover, if we want to eliminate every spam, we need only check for membership in  $S$  those good and bad emails that get through the filter. Those checks will require the use of secondary memory to access  $S$  itself. There are also other options, as we shall see when we study the general Bloom-filtering technique. As a simple example, we could use a cascade of filters, each of which would eliminate  $7/8$ th of the remaining spam.

## The Bloom Filter

A Bloom filter consists of:

1. An array of  $n$  bits, initially all 0's.
2. A collection of hash functions  $h_1, h_2, \dots, h_k$ . Each hash function maps "key" values to  $n$  buckets, corresponding to the  $n$  bits of the bit-array.
3. A set  $S$  of  $m$  key values.

The purpose of the Bloom filter is to allow through all stream elements whose keys are in  $S$ , while rejecting most of the stream elements whose keys are not in  $S$ . To initialize the bit array, begin with all bits 0. Take each key value in  $S$  and hash it using each of the  $k$  hash functions. Set to 1 each bit that is  $h_i(K)$  for some hash function  $h_i$  and some key value  $K$  in  $S$ . To test a key  $K$  that arrives in the stream, check that all of  $h_1(K), h_2(K), \dots, h_k(K)$  are 1's in the bit-array. If all are 1's, then let the stream element through. If one or more of these bits are 0, then  $K$  could not be in  $S$ , so reject the stream element.

## Analysis of Bloom Filtering

If a key value is in  $S$ , then the element will surely pass through the Bloom filter. However, if the key value is not in  $S$ , it might still pass. We need to understand how to calculate the probability of a false positive, as a function of  $n$ , the bit-array length,  $m$  the number of members of  $S$ , and  $k$ , the number of hash functions. The model to use is throwing darts at targets. Suppose we have  $x$  targets and  $y$  darts. Any dart is equally likely to hit any target. After throwing the darts, how many targets can we expect to be hit at least once? The analysis is similar to the analysis in Section 3.4.2, and goes as follows:

- The probability that a given dart will not hit a given target is  $(x - 1)/x$ .
- The probability that none of the  $y$  darts will hit a given target is  $(\frac{x-1}{x})^y$ . We can write this expression as  $(1 - \frac{1}{x})^{x(\frac{y}{x})}$ .
- Using the approximation  $(1 - \epsilon)^{1/\epsilon} = 1/e$  for small  $\epsilon$  (recall Section 1.3.5), we conclude that the probability that none of the  $y$  darts hit a given target is  $e^{-y/x}$ .

**Example 3 :** Consider the running example of Section 4.3.1. We can use the above calculation to get the true expected number of 1's in the bit array. Think of each bit as a target, and each member of  $S$  as a dart. Then the probability that a given bit will be 1 is the probability that the corresponding target will be hit by one or more darts. Since there are one billion members of  $S$ , we have  $y = 10^9$  darts. As there are eight billion bits, there are  $x = 8 \times 10^9$  targets. Thus, the probability that a given target is not hit is  $e^{-y/x} = e^{-1/8}$  and the probability that it is hit is  $1 - e^{-1/8}$ . That quantity is about 0.1175. we suggested that  $1/8 = 0.125$  is a good approximation, which it is, but now we have the exact calculation. 2

We can apply the rule to the more general situation, where set  $S$  has  $m$  members, the array has  $n$  bits, and there are  $k$  hash functions. The number of targets is  $x = n$ , and the number of darts is  $y = km$ . Thus, the probability that a bit remains 0 is  $e^{-km/n}$ . We want the fraction of 0 bits to be fairly large, or else the probability that a nonmember of  $S$  will hash at least once to a 0 becomes too small, and there are too many false positives. For example, we might choose  $k$ , the number of hash functions to be  $n/m$  or less. Then the probability of a 0 is at least  $e^{-1}$  or 37%. In general, the probability of a false positive is the probability of a 1 bit, which is  $1 - e^{-km/n}$ , raised to the  $k$ th power, i.e.,  $(1 - e^{-km/n})^k$ .

**Example 4 :** In Example 3 we found that the fraction of 1's in the array of our running example is 0.1175, and this fraction is also the probability of a false positive. That is, a nonmember of  $S$  will pass through the filter if it hashes to a 1, and the probability of it doing so is 0.1175. Suppose we used the same  $S$  and the same array, but used two different hash functions. This situation corresponds to throwing two billion darts at eight billion targets, and the probability that a bit remains 0 is  $e^{-1/4}$ .

In order to be a false positive, a nonmember of  $S$  must hash twice to bits that are 1, and this probability is  $(1 - e^{-1/4})^2$ , or approximately 0.0493. Thus, adding a second hash function for our running example is an improvement, reducing the false-positive rate from 0.1175 to 0.0493.

### The Count-Distinct Problem

Suppose stream elements are chosen from some universal set. We would like to know how many different elements have appeared in the stream, counting either from the beginning of the stream or from some known time in the past.

**Example 5 :** As a useful example of this problem, consider a Web site gathering statistics on how many unique users it has seen in each given month. The universal set is the set of logins for that site, and a stream element is generated each time someone logs in. This measure is appropriate for a site like Amazon, where the typical user logs in with their unique login name. A similar problem is a Web site like Google that does not require login to issue a search query, and may be able to identify users only by the IP address from which they send the query. There are about 4 billion IP addresses,<sup>2</sup> sequences of four 8-bit bytes will serve as the universal set in this case. 2

The obvious way to solve the problem is to keep in main memory a list of all the elements seen so far in the stream. Keep them in an efficient search structure such as a hash table or search tree, so one can quickly add new elements and check whether or not the element that just arrived on the stream was already seen. As long as the number of distinct elements is not too great, this structure can fit in main memory and there is little problem obtaining an exact answer to the question how many distinct elements appear in the stream.

However, if the number of distinct elements is too great, or if there are too many streams that need to be processed at once (e.g., Yahoo! wants to count the number of unique users viewing each of its pages in a month), then we cannot store the needed data in main memory. There are several options. We could use more machines, each machine handling only one or several of the streams. We could store most of the data structure in secondary memory and batch stream elements so whenever we brought a disk block to main memory there would be many tests and updates to be performed on the data in that block.

## The Flajolet-Martin Algorithm

It is possible to estimate the number of distinct elements by hashing the elements of the universal set to a bit-string that is sufficiently long. The length of the bit-string must be sufficient that there are more possible results of the hash function than there are elements of the universal set. For example, 64 bits is sufficient to hash URL's. We shall pick many different hash functions and hash each element of the stream using these hash functions. The important property of a hash function is that when applied to the same element, it always produces the same result.

The idea behind the Flajolet-Martin Algorithm is that the more different elements we see in the stream, the more different hash-values we shall see. As we see more different hash-values, it becomes more likely that one of these values will be "unusual." The particular unusual property we shall exploit is that the value ends in many 0's, although many other options exist. Whenever we apply a hash function  $h$  to a stream element  $a$ , the bit string  $h(a)$  will end in some number of 0's, possibly none. Call this number the tail length for  $a$  and  $h$ . Let  $R$  be the maximum tail length of any  $a$  seen so far in the stream. Then we shall use estimate  $2R$  for the number of distinct elements seen in the stream. This estimate makes intuitive sense.

The probability that a given stream element  $a$  has  $h(a)$  ending in at least  $r$  0's is  $2^{-r}$ . Suppose there are  $m$  distinct elements in the stream. Then the probability that none of them has tail length at least  $r$  is  $(1 - 2^{-r})^m$ . This sort of expression should be familiar by now.

We can rewrite it as  $((1 - 2^{-r})^{2^r})^{m/2^r}$ . Assuming  $r$  is reasonably large, the inner expression is of the form  $(1 - \epsilon)^{1/\epsilon}$ , which is approximately  $1/e$ . Thus, the probability of not finding a stream element with as many as  $r$  0's at the end of its hash value is  $e^{-m2^{-r}}$ . We can conclude:

1. If  $m$  is much larger than  $2^r$ , then the probability that we shall find a tail of length at least  $r$  approaches 1.
2. If  $m$  is much less than  $2^r$ , then the probability of finding a tail length at least  $r$  approaches 0.

We conclude from these two points that the proposed estimate of  $m$ , which is  $2_R$  (recall  $R$  is the largest tail length for any stream element) is unlikely to be either much too high or much too low.

### Combining Estimates

Unfortunately, there is a trap regarding the strategy for combining the estimates of  $m$ , the number of distinct elements, that we obtain by using many different hash functions. Our first assumption would be that if we take the average of the values  $2_R$  that we get from each hash function, we shall get a value that approaches the true  $m$ , the more hash functions we use. However, that is not the case, and the reason has to do with the influence an overestimate has on the average.

Consider a value of  $r$  such that  $2_r$  is much larger than  $m$ . There is some probability  $p$  that we shall discover  $r$  to be the largest number of 0's at the end of the hash value for any of the  $m$  stream elements. Then the probability of finding  $r+1$  to be the largest number of 0's instead is at least  $p/2$ . However, if we do increase by 1 the number of 0's at the end of a hash value, the value of  $2_R$  doubles. Consequently, the contribution from each possible large  $R$  to the expected value of  $2_R$  grows as  $R$  grows, and the expected value of  $2_R$  is actually infinite.<sup>3</sup>

Another way to combine estimates is to take the median of all estimates. The median is not affected by the occasional outsized value of  $2_R$ , so the worry described above for the average should not carry over to the median. Unfortunately, the median suffers from another defect: it is always a power of 2. Thus, no matter how many hash functions we use, should the correct value of  $m$  be between two powers of 2, say 400, then it will be impossible to obtain a close estimate.

There is a solution to the problem, however. We can combine the two methods. First, group the hash functions into small groups, and take their average. Then, take the median of the averages. It is true that an occasional oversized  $2R$  will bias some of the groups and make them too large. However, taking the median of group averages will reduce the influence of this effect almost to nothing. Moreover, if the groups themselves are large enough, then the averages can be essentially any number, which enables us to approach the true value  $m$  as long as we use enough hash functions. In order to guarantee that any possible average can be obtained, groups should be of size at least a small multiple of  $\log_2 m$ .

## Space Requirements

Observe that as we read the stream it is not necessary to store the elements seen. The only thing we need to keep in main memory is one integer per hash function; this integer records the largest tail length seen so far for that hash function and any stream element. If we are processing only one stream, we could use millions of hash functions, which is far more than we need to get a 3. Technically, since the hash value is a bit-string of finite length, there is no contribution to  $2R$  for  $R$ 's that are larger than the length of the hash value. However, this effect is not enough to avoid the conclusion that the expected value of  $2R$  is much too large.

## ESTIMATING MOMENTS

### Definition of Moments:

Suppose a stream consists of elements chosen from a universal set. Assume the universal set is ordered so we can speak of the  $i$ th element for any  $i$ . Let  $m_i$  be the number of occurrences of the  $i$ th element for any  $i$ . Then the  $k$ th-order moment (or just  $k$ th moment) of the stream is the sum over all  $i$  of  $(m_i)^k$ .

**Example 6 :** The 0th moment is the sum of 1 for each  $m_i$  that is greater than 0. That is, the 0th moment is a count of the number of distinct elements in the stream. We can use the method of Section 4.4 to estimate the 0th moment of a stream. Technically, since  $m_i$  could be 0 for some elements in the universal set, we need to make explicit in the definition of "moment" that 0<sup>0</sup> is taken to be 0. For moments 1 and above, the contribution of  $m_i$ 's that are 0 is surely 0. The 1st moment is the sum of the  $m_i$ 's, which must be the length of the stream. Thus, first moments are especially easy to compute; just count the length of the stream seen so far. The second moment is the sum of the squares of the  $m_i$ 's. It is sometimes called the surprise number, since it measures how uneven the distribution of elements in the stream is. To see the distinction, suppose we have a stream of length 100, in which eleven different elements appear.

The most even distribution of these eleven elements would have one appearing 10 times and the other ten appearing 9 times each. In this case, the surprise number is  $10 \times 10 + 10 \times 9 = 190$ . At the other extreme, one of the eleven elements could appear 90 times and the other ten appear 1 time each. Then, the surprise number would be  $10 \times 10 + 10 \times 1 = 110$ .

There is no problem computing moments of any order if we can afford to keep in main memory a count for each element that appears in the stream. However, also as in that section, if we cannot afford to use that much memory, then we need to estimate the  $k$ th moment by keeping a limited number of values in main memory and computing an estimate from these values. For the case of distinct elements, each of these values were counts of the longest tail produced by a single hash function. We shall see another form of value that is useful for second and higher moments.

### The Alon-Matias-Szegedy Algorithm for Second Moments

For now, let us assume that a stream has a particular length  $n$ . We shall show how to deal with growing streams in the next section. Suppose we do not have enough space to count all the  $m_i$ 's for all the elements of the stream. We can still estimate the second moment of the stream using a limited amount of space; the more space we use, the more accurate the estimate will be. We compute some number of variables. For each variable  $X$ , we store:

1. A particular element of the universal set, which we refer to as  $X.\text{element}$ ,  
and
2. An integer  $X.\text{value}$ , which is the value of the variable.

To determine the value of a variable  $X$ , we choose a position in the stream between 1 and  $n$ , uniformly and at random. Set  $X.\text{element}$  to be the element found there, and initialize  $X.\text{value}$  to 1. As we read the stream, add 1 to  $X.\text{value}$  each time we encounter another occurrence of  $X.\text{element}$ .

**Example 7 :** Suppose the stream is a, b, c, b, d, a, c, d, a, b, d, c, a, a, b. The length of the stream is  $n = 15$ . Since a appears 5 times, b appears 4 times, and c and d appear three times each, the second moment for the stream is  $5^2 + 4^2 + 3^2 + 3^2 = 59$ . Suppose we keep three variables,  $X_1$ ,  $X_2$ , and  $X_3$ . Also, assume that at "random" we pick the 3rd, 8th, and 13th positions to define these three variables. When we reach position 3, we find element c, so we set  $X_1.\text{element} = c$  and  $X_1.\text{value} = 1$ . Position 4 holds b, so we do not change  $X_1$ .

Likewise, nothing happens at positions 5 or 6. At position 7, we see c again, so we set  $X1.value = 2$ . At position 8 we find d, and so set  $X2.element = d$  and  $X2.value = 1$ . Positions 9 and 10 hold a and b, so they do not affect X1 or X2. Position 11 holds d so we set  $X2.value = 2$ , and position 12 holds c so we set  $X1.value = 3$ . At position 13, we find element a, and so set  $X3.element = a$  and  $X3.value = 1$ . Then, at position 14 we see another a and so set  $X3.value = 2$ . Position 15, with element b does not affect any of the variables, so we are done, with final values  $X1.value = 3$  and  $X2.value = X3.value = 2$ . We can derive an estimate of the second moment from any variable X. This estimate is  $n(2X.value - 1)$ .

### Higher-Order Moments

We estimate kth moments, for  $k > 2$ , in essentially the same way as we estimate second moments. The only thing that changes is the way we derive an estimate from a variable. In Section 4.5.2 we used the formula  $n(2v - 1)$  to turn a value v, the count of the number of occurrences of some particular stream element a, into an estimate of the second moment. We saw why this formula works: the terms  $2v - 1$ , for  $v = 1, 2, \dots, m$  sum to  $m^2$ , where m is the number of times a appears in the stream. Notice that  $2v - 1$  is the difference between  $v^2$  and  $(v - 1)^2$ . Suppose we wanted the third moment rather than the second. Then all we have to do is replace  $2v - 1$  by  $v^3 - (v - 1)^3 = 3v^2 - 3v + 1$ . Then  $\sum_{v=1}^m (3v^2 - 3v + 1) = m^3$ , so we can use as our estimate of the third moment the formula  $n(3v^2 - 3v + 1)$ , where  $v = X.value$  is the value associated with some variable X. More generally, we can estimate kth moments for any  $k \geq 2$  by turning value  $v = X.value$  into  $n(v^k - (v - 1)^k)$ .

### Dealing With Infinite Streams

Technically, the estimate we used for second and higher moments assumes that n, the stream length, is a constant. In practice, n grows with time. That fact, by itself, doesn't cause problems, since we store only the values of variables and multiply some function of that value by n when it is time to estimate the moment. If we count the number of stream elements seen and store this value, which only requires  $\log n$  bits, then we have n available whenever we need it.

A more serious problem is that we must be careful how we select the positions for the variables. If we do this selection once and for all, then as the stream gets longer, we are biased in favor of early positions, and the estimate of the moment will be too large. On the other hand, if we wait too long to pick positions, then early in the stream we do not have many variables and so will get an unreliable estimate.

The proper technique is to maintain as many variables as we can store at all times, and to throw some out as the stream grows. The discarded variables are replaced by new ones, in such a way that at all times, the probability of picking any one position for a variable is the same as that of picking any other position. Suppose we have space to store  $s$  variables. Then the first  $s$  positions of the stream are each picked as the position of one of the  $s$  variables.

Inductively, suppose we have seen  $n$  stream elements, and the probability of any particular position being the position of a variable is uniform, that is  $s/n$ . When the  $(n+1)$ st element arrives, pick that position with probability  $s/(n+1)$ . If not picked, then the  $s$  variables keep their same positions. However, if the  $(n+1)$ st position is picked, then throw out one of the current  $s$  variables, with equal probability. Replace the one discarded by a new variable whose element is the one at position  $n + 1$  and whose value is 1. Surely, the probability that position  $n + 1$  is selected for a variable is what it should be:  $s/(n + 1)$ . However, the probability of every other position also is  $s/(n + 1)$ , as we can prove by induction on  $n$ . By the inductive hypothesis, before the arrival of the  $(n + 1)$ st stream element, this probability was  $s/n$ . With probability  $1 - s/(n + 1)$  the  $(n + 1)$ st position will not be selected, and the probability of each of the first  $n$  positions remains  $s/n$ . However, with probability  $s/(n + 1)$ , the  $(n + 1)$ st position is picked, and the probability for each of the first  $n$  positions is reduced by factor  $(s-1)/s$ .

Considering the two cases, the probability of selecting each of the first  $n$  positions is

$$\left(1 - \frac{s}{n+1}\right)\left(\frac{s}{n}\right) + \left(\frac{s}{n+1}\right)\left(\frac{s-1}{s}\right)\left(\frac{s}{n}\right)$$

This expression simplifies to

$$\left(1 - \frac{s}{n+1}\right)\left(\frac{s}{n}\right) + \left(\frac{s-1}{n+1}\right)\left(\frac{s}{n}\right)$$

and then to

$$\left(\left(1 - \frac{s}{n+1}\right) + \left(\frac{s-1}{n+1}\right)\right)\left(\frac{s}{n}\right)$$

which in turn simplifies to

$$\left(\frac{n}{n+1}\right)\left(\frac{s}{n}\right) = \frac{s}{n+1}$$

Thus, we have shown by induction on the stream length  $n$  that all positions have equal probability  $s/n$  of being chosen as the position of a variable.