

# Randomised Algorithms

## Lecture 13: Streaming Algorithms

Thomas Sauerwald ([tms41@cam.ac.uk](mailto:tms41@cam.ac.uk))

# Outline

---

Introduction

Approximate Counting

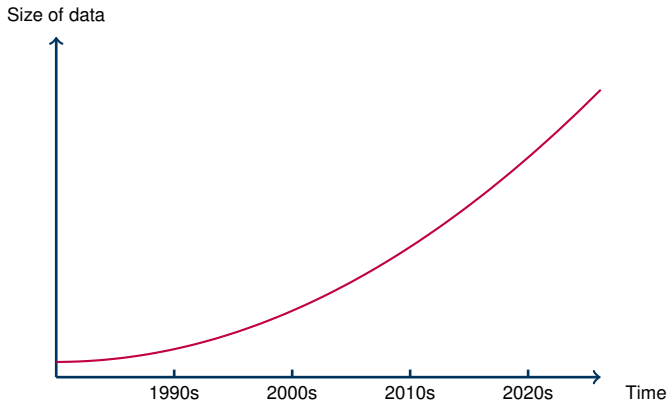
Distinct Elements and Frequency Moments

Extra Material (non-examinable): An Algorithm for  $F_0$  in the Turnstile Model

## Background of Streaming Algorithms

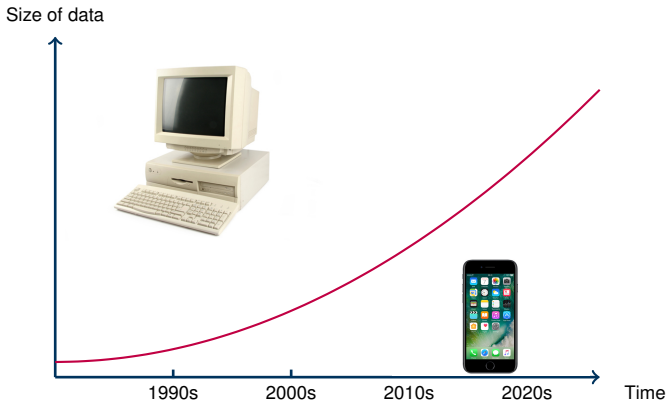
---

- The amount of data has been increased exponentially over the last years



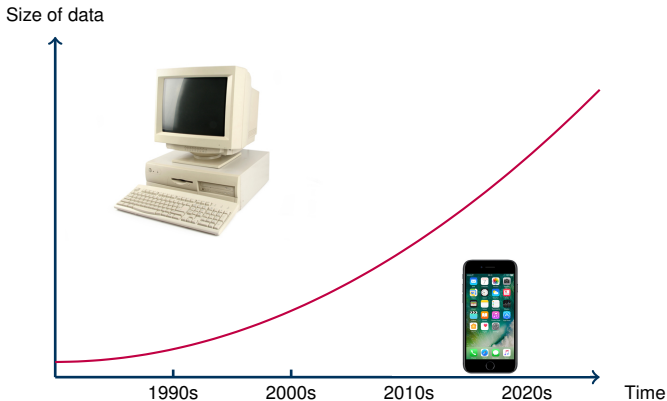
## Background of Streaming Algorithms

- The amount of data has been increased exponentially over the last years
- For many applications computational devices' memories are limited



## Background of Streaming Algorithms

- The amount of data has been increased exponentially over the last years
- For many applications computational devices' memories are limited
- We need to find good (approximate) solutions without storing the entire input!



# Motivation: Analysing Search Engine Queries

---



# Motivation: Analysing Search Engine Queries



IP: 54.73.136.89										
Time:										
Text:										

.....

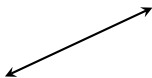
# Motivation: Analysing Search Engine Queries



IP: 54.73.136.89	IP: 102.58.22.231								.....
Time:	Time:								
Text:	Text:								



# Motivation: Analysing Search Engine Queries



IP: 54.73.136.89	IP: 102.58.22.231	IP: 54.73.136.89							.....
Time:	Time:	Time:							
Text:	Text:	Text:							

# Motivation: Analysing Search Engine Queries



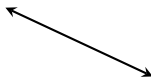
IP: 54.73.136.89	IP: 102.58.22.231	IP: 54.73.136.89	IP: 170.9.103.244					.....
Time:	Time:	Time:	Time:					
Text:	Text:	Text:	Text:					

# Motivation: Analysing Search Engine Queries



IP: 54.73.136.89	IP: 102.58.22.231	IP: 54.73.136.89	IP: 170.9.103.244	IP: 189.105.32.75				.....
Time:	Time:	Time:	Time:	Time:				
Text:	Text:	Text:	Text:	Text:				

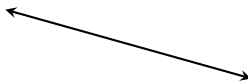
# Motivation: Analysing Search Engine Queries



IP: 54.73.136.89	IP: 102.58.22.231	IP: 54.73.136.89	IP: 170.9.103.244	IP: 189.105.32.75	IP: 54.73.136.89		
Time:	Time:	Time:	Time:	Time:	Time:		
Text:	Text:	Text:	Text:	Text:	Text:		

.....

# Motivation: Analysing Search Engine Queries



IP: 54.73.136.89	IP: 102.58.22.231	IP: 54.73.136.89	IP: 170.9.103.244	IP: 189.105.32.75	IP: 54.73.136.89	IP: 144.66.18.240	
Time:	Time:	Time:	Time:	Time:	Time:	Time:	
Text:	Text:	Text:	Text:	Text:	Text:	Text:	

.....

# Motivation: Analysing Search Engine Queries



IP: 54.73.136.89	IP: 102.58.22.231	IP: 54.73.136.89	IP: 170.9.103.244	IP: 189.105.32.75	IP: 54.73.136.89	IP: 144.66.18.240	IP: 102.58.22.231	.....
Time:	Time:	Time:	Time:	Time:	Time:	Time:	Time:	
Text:	Text:	Text:	Text:	Text:	Text:	Text:	Text:	

# Motivation: Analysing Search Engine Queries



IP: 54.73.136.89	IP: 102.58.22.231	IP: 54.73.136.89	IP: 170.9.103.244	IP: 189.105.32.75	IP: 54.73.136.89	IP: 144.66.18.240	IP: 102.58.22.231
Time:	Time:	Time:	Time:	Time:	Time:	Time:	Time:
Text:	Text:	Text:	Text:	Text:	Text:	Text:	Text:

.....

# Motivation: Analysing Search Engine Queries



- memory is much smaller than needed to store entire data stream

IP: 54.73.136.89	IP: 102.58.22.231	IP: 54.73.136.89	IP: 170.9.103.244	IP: 189.105.32.75	IP: 54.73.136.89	IP: 144.66.18.240	IP: 102.58.22.231	.....
Time:	Time:	Time:	Time:	Time:	Time:	Time:	Time:	
Text:	Text:	Text:	Text:	Text:	Text:	Text:	Text:	



## Motivation: Analysing Search Engine Queries



- memory is much smaller than needed to store entire data stream  
⇒ We can only read each data item **once** and in sequential order

IP: 54.73.136.89	IP: 102.58.22.231	IP: 54.73.136.89	IP: 170.9.103.244	IP: 189.105.32.75	IP: 54.73.136.89	IP: 144.66.18.240	IP: 102.58.22.231	.....
Time:	Time:	Time:	Time:	Time:	Time:	Time:	Time:	
Text:	Text:	Text:	Text:	Text:	Text:	Text:	Text:	

## Motivation: Analysing Search Engine Queries

- What is the total number queries?
- What is the total number of different IP addresses?
- Extension 1: only consider queries within a certain interval (sliding window)
- Extension 2: also allow the cancellation/removal of a query (turnstile model)
- Extension 3: What if we have different data centers? (distributed streaming)
- ⋮

- **memory is much smaller** than needed to store entire data stream  
⇒ We can only read each data item **once** and in sequential order

IP: 54.73.136.89	IP: 102.58.22.231	IP: 54.73.136.89	IP: 170.9.103.244	IP: 189.105.32.75	IP: 54.73.136.89	IP: 144.66.18.240	IP: 102.58.22.231	.....
Time:	Time:	Time:	Time:	Time:	Time:	Time:	Time:	
Text:	Text:	Text:	Text:	Text:	Text:	Text:	Text:	

## Motivation: Analysing Search Engine Queries

- What is the total number queries?
- What is the total number of different IP addresses?
- Extension 1: only consider queries within a certain interval (sliding window)
- Extension 2: also allow the cancellation/removal of a query (turnstile model)
- Extension 3: What if we have different data centers? (distributed streaming)
- ⋮

- **memory** is **much smaller** than needed to store entire data stream  
⇒ We can only read each data item **once** and in sequential order

IP:  
54.73.1  
Time:  
Text:

### Other Applications:

- Monitoring Financial Transactions
- Analysing Buying Histories of Users

## Motivation: Analysing Search Engine Queries

- What is the total number queries?
- What is the total number of different IP addresses?
- Extension 1: only consider queries within a certain interval (sliding window)
- Extension 2: also allow the cancellation/removal of a query (turnstile model)
- Extension 3: What if we have different data centers? (distributed streaming)
- ⋮

- memory is much smaller than needed to store entire data stream  
⇒ We can only read each data item once and in sequential order

IP:  
54.73.1  
Time:  
Text:

### Other Applications:

- Monitoring Financial Transactions
- Analysing Buying Histories of Users

- The **input** of a streaming algorithm is given as a **data stream**, which is a sequence of data

$$\mathcal{S} = s_1, s_2, \dots, s_i, \dots$$

and every  $s_i$  belongs to the universe  $U$ .

## Streaming algorithms

---

- The **input** of a streaming algorithm is given as a **data stream**, which is a sequence of data

$$\mathcal{S} = s_1, s_2, \dots, s_i, \dots$$

and every  $s_i$  belongs to the universe  $U$ .

- **Constraints for streaming algorithms:** the space complexity should be **sublinear** in  $|U|$  and  $|\mathcal{S}|$ .

- The **input** of a streaming algorithm is given as a **data stream**, which is a sequence of data

$$\mathcal{S} = s_1, s_2, \dots, s_i, \dots$$

and every  $s_i$  belongs to the universe  $U$ .

- **Constraints for streaming algorithms:** the space complexity should be sublinear in  $|U|$  and  $|\mathcal{S}|$ .
- **Quality of the output:** The algorithm needs to give a good **approximate** value with **high probability**.

## Streaming algorithms

- The **input** of a streaming algorithm is given as a **data stream**, which is a sequence of data

$$\mathcal{S} = s_1, s_2, \dots, s_i, \dots$$

and every  $s_i$  belongs to the universe  $U$ .

- Constraints for streaming algorithms:** the space complexity should be **sublinear** in  $|U|$  and  $|\mathcal{S}|$ .
- Quality of the output:** The algorithm needs to give a good **approximate** value with **high probability**.

$(\epsilon, \delta)$ -approximation

For confidence parameter  $\delta$  and approximation parameter  $\epsilon$ , the algorithm's output **Output** and the exact answer **Exact** satisfies

$$\mathbf{P}[\mathbf{Output} \in (1 - \epsilon, 1 + \epsilon) \cdot \mathbf{Exact}] \geq 1 - \delta.$$



# Outline

---

Introduction

Approximate Counting

Distinct Elements and Frequency Moments

Extra Material (non-examinable): An Algorithm for  $F_0$  in the Turnstile Model

## Approximate Counting and Morris Algorithm

---

### Approximate Counting

An approximate counting algorithm must monitor a sequence of events. At any given time, the algorithm must output an **estimate** of the **number of events**.

## Approximate Counting and Morris Algorithm

This could be also described as a data structure maintaining an integer  $n$  and supporting two operations:

- `update()`: increment  $n$  by 1
- `query()`: output  $n$

### Approximate Counting

An approximate counting algorithm must monitor a sequence of events. At any given time, the algorithm must output an **estimate** of the **number of events**.

## Approximate Counting and Morris Algorithm

This could be also described as a data structure maintaining an integer  $n$  and supporting two operations:

- `update()`: increment  $n$  by 1
- `query()`: output  $n$

### Approximate Counting

An approximate counting algorithm must monitor a sequence of events. At any given time, the algorithm must output an **estimate** of the **number of events**.

Trivial (and exact) solution uses  $\log_2 n$  space. Can we do better?

## Approximate Counting and Morris Algorithm

This could be also described as a data structure maintaining an integer  $n$  and supporting two operations:

- `update()`: increment  $n$  by 1
- `query()`: output  $n$

### Approximate Counting

An approximate counting algorithm must monitor a sequence of events. At any given time, the algorithm must output an **estimate** of the **number of events**.

Trivial (and exact) solution uses  $\log_2 n$  space. Can we do better?

### MORRIS ALGORITHM

- 1:  $X \leftarrow 0$
- 2: **While** `update()`
- 3:     With probability  $2^{-X}$  set  $X \leftarrow X + 1$
- 4: **Return**  $2^X - 1$

## Approximate Counting and Morris Algorithm

This could be also described as a data structure maintaining an integer  $n$  and supporting two operations:

- `update()`: increment  $n$  by 1
- `query()`: output  $n$

### Approximate Counting

An approximate counting algorithm must monitor a sequence of events. At any given time, the algorithm must output an **estimate** of the **number of events**.

Trivial (and exact) solution uses  $\log_2 n$  space. Can we do better?

### MORRIS ALGORITHM

- 1:  $X \leftarrow 0$
- 2: **While** `update()`
- 3:     With probability  $2^{-X}$  set  $X \leftarrow X + 1$
- 4: **Return**  $2^X - 1$

**Intuition:**  $X$  will be an approximation of  $\log_2 n$  (that is, we try to approximate the number of bits of  $n$  in binary)

## Analysis (1/3)

---

### Lemma (Expectation Analysis)

Let  $X_n$  denote the value of  $X$  after  $n$  updates. For every  $n \geq 0$ ,

$$\mathbf{E} \left[ 2^{X_n} \right] = n + 1.$$

## Analysis (1/3)

---

### Lemma (Expectation Analysis)

Let  $X_n$  denote the value of  $X$  after  $n$  updates. For every  $n \geq 0$ ,

$$\mathbf{E} \left[ 2^{X_n} \right] = n + 1.$$

Proof:



## Analysis (1/3)

---

### Lemma (Expectation Analysis)

Let  $X_n$  denote the value of  $X$  after  $n$  updates. For every  $n \geq 0$ ,

$$\mathbf{E} \left[ 2^{X_n} \right] = n + 1.$$

Proof:

- Base case: For  $n = 0$ , we have  $X_n = X_0 = 0$

## Analysis (1/3)

---

### Lemma (Expectation Analysis)

Let  $X_n$  denote the value of  $X$  after  $n$  updates. For every  $n \geq 0$ ,

$$\mathbf{E} \left[ 2^{X_n} \right] = n + 1.$$

Proof:

- Base case: For  $n = 0$ , we have  $X_n = X_0 = 0 \checkmark$

## Analysis (1/3)

### Lemma (Expectation Analysis)

Let  $X_n$  denote the value of  $X$  after  $n$  updates. For every  $n \geq 0$ ,

$$\mathbf{E} \left[ 2^{X_n} \right] = n + 1.$$

Proof:

- **Base case:** For  $n = 0$ , we have  $X_n = X_0 = 0 \checkmark$
- **Induction step:**  $n \rightarrow n + 1$ : By conditioning on  $X_n$ ,

$$\mathbf{E} \left[ 2^{X_{n+1}} \right] = \sum_{j=0}^{\infty} \mathbf{P} [ X_n = j ] \cdot \mathbf{E} \left[ 2^{X_{n+1}} \mid X_n = j \right]$$

## Analysis (1/3)

### Lemma (Expectation Analysis)

Let  $X_n$  denote the value of  $X$  after  $n$  updates. For every  $n \geq 0$ ,

$$\mathbf{E} \left[ 2^{X_n} \right] = n + 1.$$

Proof:

- **Base case:** For  $n = 0$ , we have  $X_n = X_0 = 0$  ✓
- **Induction step:**  $n \rightarrow n + 1$ : By conditioning on  $X_n$ ,

$$\begin{aligned} \mathbf{E} \left[ 2^{X_{n+1}} \right] &= \sum_{j=0}^{\infty} \mathbf{P} [ X_n = j ] \cdot \mathbf{E} \left[ 2^{X_{n+1}} \mid X_n = j \right] \\ &= \sum_{j=0}^{\infty} \mathbf{P} [ X_n = j ] \cdot \left( 2^j \cdot \left( 1 - \frac{1}{2^j} \right) + 2^{j+1} \cdot \frac{1}{2^j} \right) \end{aligned}$$

## Analysis (1/3)

### Lemma (Expectation Analysis)

Let  $X_n$  denote the value of  $X$  after  $n$  updates. For every  $n \geq 0$ ,

$$\mathbf{E} \left[ 2^{X_n} \right] = n + 1.$$

Proof:

- **Base case:** For  $n = 0$ , we have  $X_n = X_0 = 0$  ✓
- **Induction step:**  $n \rightarrow n + 1$ : By conditioning on  $X_n$ ,

$$\begin{aligned} \mathbf{E} \left[ 2^{X_{n+1}} \right] &= \sum_{j=0}^{\infty} \mathbf{P} [ X_n = j ] \cdot \mathbf{E} \left[ 2^{X_{n+1}} \mid X_n = j \right] \\ &= \sum_{j=0}^{\infty} \mathbf{P} [ X_n = j ] \cdot \left( 2^j \cdot \left( 1 - \frac{1}{2^j} \right) + 2^{j+1} \cdot \frac{1}{2^j} \right) \\ &= \sum_{j=0}^{\infty} \mathbf{P} [ X_n = j ] \cdot 2^j + \sum_{j=0}^{\infty} \mathbf{P} [ X_n = j ] \end{aligned}$$

## Analysis (1/3)

### Lemma (Expectation Analysis)

Let  $X_n$  denote the value of  $X$  after  $n$  updates. For every  $n \geq 0$ ,

$$\mathbf{E} \left[ 2^{X_n} \right] = n + 1.$$

Proof:

- **Base case:** For  $n = 0$ , we have  $X_n = X_0 = 0$  ✓
- **Induction step:**  $n \rightarrow n + 1$ : By conditioning on  $X_n$ ,

$$\begin{aligned} \mathbf{E} \left[ 2^{X_{n+1}} \right] &= \sum_{j=0}^{\infty} \mathbf{P} [ X_n = j ] \cdot \mathbf{E} \left[ 2^{X_{n+1}} \mid X_n = j \right] \\ &= \sum_{j=0}^{\infty} \mathbf{P} [ X_n = j ] \cdot \left( 2^j \cdot \left( 1 - \frac{1}{2^j} \right) + 2^{j+1} \cdot \frac{1}{2^j} \right) \\ &= \sum_{j=0}^{\infty} \mathbf{P} [ X_n = j ] \cdot 2^j + \sum_{j=0}^{\infty} \mathbf{P} [ X_n = j ] \\ &= \mathbf{E} \left[ 2^{X_n} \right] + 1 \end{aligned}$$

## Analysis (1/3)

### Lemma (Expectation Analysis)

Let  $X_n$  denote the value of  $X$  after  $n$  updates. For every  $n \geq 0$ ,

$$\mathbf{E} \left[ 2^{X_n} \right] = n + 1.$$

Proof:

- **Base case:** For  $n = 0$ , we have  $X_n = X_0 = 0$  ✓
- **Induction step:**  $n \rightarrow n + 1$ : By conditioning on  $X_n$ ,

$$\begin{aligned} \mathbf{E} \left[ 2^{X_{n+1}} \right] &= \sum_{j=0}^{\infty} \mathbf{P} [ X_n = j ] \cdot \mathbf{E} \left[ 2^{X_{n+1}} \mid X_n = j \right] \\ &= \sum_{j=0}^{\infty} \mathbf{P} [ X_n = j ] \cdot \left( 2^j \cdot \left( 1 - \frac{1}{2^j} \right) + 2^{j+1} \cdot \frac{1}{2^j} \right) \\ &= \sum_{j=0}^{\infty} \mathbf{P} [ X_n = j ] \cdot 2^j + \sum_{j=0}^{\infty} \mathbf{P} [ X_n = j ] \\ &= \mathbf{E} \left[ 2^{X_n} \right] + 1 \\ &= (n + 1) + 1. \end{aligned}$$

## Analysis (1/3)

### Lemma (Expectation Analysis)

Let  $X_n$  denote the value of  $X$  after  $n$  updates. For every  $n \geq 0$ ,

$$\mathbf{E} \left[ 2^{X_n} \right] = n + 1.$$

Proof:

- **Base case:** For  $n = 0$ , we have  $X_n = X_0 = 0 \checkmark$
- **Induction step:**  $n \rightarrow n + 1$ : By conditioning on  $X_n$ ,

$$\begin{aligned} \mathbf{E} \left[ 2^{X_{n+1}} \right] &= \sum_{j=0}^{\infty} \mathbf{P} [ X_n = j ] \cdot \mathbf{E} \left[ 2^{X_{n+1}} \mid X_n = j \right] \\ &= \sum_{j=0}^{\infty} \mathbf{P} [ X_n = j ] \cdot \left( 2^j \cdot \left( 1 - \frac{1}{2^j} \right) + 2^{j+1} \cdot \frac{1}{2^j} \right) \\ &= \sum_{j=0}^{\infty} \mathbf{P} [ X_n = j ] \cdot 2^j + \sum_{j=0}^{\infty} \mathbf{P} [ X_n = j ] \\ &= \mathbf{E} \left[ 2^{X_n} \right] + 1 \\ &= (n + 1) + 1. \end{aligned}$$

By Induction Hypothesis



## Analysis (1/3)

### Lemma (Expectation Analysis)

Let  $X_n$  denote the value of  $X$  after  $n$  updates. For every  $n \geq 0$ ,

$$\mathbf{E} \left[ 2^{X_n} \right] = n + 1.$$

Hence  $\Theta_n := 2^{X_n} - 1$  is an unbiased estimator of  $n$ .

Proof:

- Base case: For  $n = 0$ , we have  $X_n = X_0 = 0$  ✓
- Induction step:  $n \rightarrow n + 1$ : By conditioning on  $X_n$ ,

$$\begin{aligned} \mathbf{E} \left[ 2^{X_{n+1}} \right] &= \sum_{j=0}^{\infty} \mathbf{P} [ X_n = j ] \cdot \mathbf{E} \left[ 2^{X_{n+1}} \mid X_n = j \right] \\ &= \sum_{j=0}^{\infty} \mathbf{P} [ X_n = j ] \cdot \left( 2^j \cdot \left( 1 - \frac{1}{2^j} \right) + 2^{j+1} \cdot \frac{1}{2^j} \right) \\ &= \sum_{j=0}^{\infty} \mathbf{P} [ X_n = j ] \cdot 2^j + \sum_{j=0}^{\infty} \mathbf{P} [ X_n = j ] \\ &= \mathbf{E} \left[ 2^{X_n} \right] + 1 \\ &= (n + 1) + 1. \end{aligned}$$

By Induction Hypothesis

## Analysis (2/3)

### Lemma (Second Moment Analysis)

Let  $X_n$  denote the value of  $X$  after  $n$  updates. For every  $n \geq 0$ ,

$$\mathbf{E} \left[ \left( 2^{X_n} \right)^2 \right] = \mathbf{E} \left[ 2^{2 \cdot X_n} \right] = \frac{3}{2} n^2 + \frac{3}{2} n + 1.$$

## Analysis (2/3)

### Lemma (Second Moment Analysis)

Let  $X_n$  denote the value of  $X$  after  $n$  updates. For every  $n \geq 0$ ,

$$\mathbf{E} \left[ \left( 2^{X_n} \right)^2 \right] = \mathbf{E} \left[ 2^{2 \cdot X_n} \right] = \frac{3}{2} n^2 + \frac{3}{2} n + 1.$$

This is shown similarly to that of the previous Lemma (see supervision sheet)

## Analysis (2/3)

### Lemma (Second Moment Analysis)

Let  $X_n$  denote the value of  $X$  after  $n$  updates. For every  $n \geq 0$ ,

$$\mathbf{E} \left[ \left( 2^{X_n} \right)^2 \right] = \mathbf{E} \left[ 2^{2 \cdot X_n} \right] = \frac{3}{2} n^2 + \frac{3}{2} n + 1.$$

This is shown similarly to that of the previous Lemma (see supervision sheet)

- Recall  $\Theta_n = 2^{X_n} - 1$ .

## Analysis (2/3)

### Lemma (Second Moment Analysis)

Let  $X_n$  denote the value of  $X$  after  $n$  updates. For every  $n \geq 0$ ,

$$\mathbf{E} \left[ \left( 2^{X_n} \right)^2 \right] = \mathbf{E} \left[ 2^{2 \cdot X_n} \right] = \frac{3}{2} n^2 + \frac{3}{2} n + 1.$$

This is shown similarly to that of the previous Lemma (see supervision sheet)

- Recall  $\Theta_n = 2^{X_n} - 1$ .
- Since  $\mathbf{V}[Z] = \mathbf{E}[Z^2] - \mathbf{E}[Z]^2$ ,

$$\mathbf{V}[\Theta_n] = \mathbf{V}[2^{X_n}]$$

## Analysis (2/3)

### Lemma (Second Moment Analysis)

Let  $X_n$  denote the value of  $X$  after  $n$  updates. For every  $n \geq 0$ ,

$$\mathbf{E} \left[ \left( 2^{X_n} \right)^2 \right] = \mathbf{E} \left[ 2^{2 \cdot X_n} \right] = \frac{3}{2} n^2 + \frac{3}{2} n + 1.$$

This is shown similarly to that of the previous Lemma (see supervision sheet)

- Recall  $\Theta_n = 2^{X_n} - 1$ .
- Since  $\mathbf{V}[Z] = \mathbf{E}[Z^2] - \mathbf{E}[Z]^2$ ,

$$\mathbf{V}[\Theta_n] = \mathbf{V}[2^{X_n}] = \mathbf{E}[2^{2 \cdot X_n}] - \left( \mathbf{E}[2^{X_n}] \right)^2$$

## Analysis (2/3)

### Lemma (Second Moment Analysis)

Let  $X_n$  denote the value of  $X$  after  $n$  updates. For every  $n \geq 0$ ,

$$\mathbf{E} \left[ \left( 2^{X_n} \right)^2 \right] = \mathbf{E} \left[ 2^{2 \cdot X_n} \right] = \frac{3}{2} n^2 + \frac{3}{2} n + 1.$$

This is shown similarly to that of the previous Lemma (see supervision sheet)

- Recall  $\Theta_n = 2^{X_n} - 1$ .
- Since  $\mathbf{V}[Z] = \mathbf{E}[Z^2] - \mathbf{E}[Z]^2$ ,

$$\begin{aligned} \mathbf{V}[\Theta_n] &= \mathbf{V}[2^{X_n}] = \mathbf{E}[2^{2 \cdot X_n}] - \left( \mathbf{E}[2^{X_n}] \right)^2 \\ &= \frac{3}{2} n^2 + \frac{3}{2} n + 1 - (n + 1)^2 = \frac{n^2 - n}{2} \end{aligned}$$

## Analysis (2/3)

### Lemma (Second Moment Analysis)

Let  $X_n$  denote the value of  $X$  after  $n$  updates. For every  $n \geq 0$ ,

$$\mathbf{E} \left[ \left( 2^{X_n} \right)^2 \right] = \mathbf{E} \left[ 2^{2 \cdot X_n} \right] = \frac{3}{2} n^2 + \frac{3}{2} n + 1.$$

This is shown similarly to that of the previous Lemma (see supervision sheet)

- Recall  $\Theta_n = 2^{X_n} - 1$ .
- Since  $\mathbf{V}[Z] = \mathbf{E}[Z^2] - \mathbf{E}[Z]^2$ ,

$$\begin{aligned} \mathbf{V}[\Theta_n] &= \mathbf{V}[2^{X_n}] = \mathbf{E}[2^{2 \cdot X_n}] - \left( \mathbf{E}[2^{X_n}] \right)^2 \\ &= \frac{3}{2} n^2 + \frac{3}{2} n + 1 - (n + 1)^2 = \frac{n^2 - n}{2} \end{aligned}$$

- Using Chebysheff's inequality,



## Analysis (2/3)

### Lemma (Second Moment Analysis)

Let  $X_n$  denote the value of  $X$  after  $n$  updates. For every  $n \geq 0$ ,

$$\mathbf{E} \left[ \left( 2^{X_n} \right)^2 \right] = \mathbf{E} \left[ 2^{2 \cdot X_n} \right] = \frac{3}{2} n^2 + \frac{3}{2} n + 1.$$

This is shown similarly to that of the previous Lemma (see supervision sheet)

- Recall  $\Theta_n = 2^{X_n} - 1$ .
- Since  $\mathbf{V}[Z] = \mathbf{E}[Z^2] - \mathbf{E}[Z]^2$ ,

$$\begin{aligned} \mathbf{V}[\Theta_n] &= \mathbf{V}[2^{X_n}] = \mathbf{E}[2^{2 \cdot X_n}] - \left( \mathbf{E}[2^{X_n}] \right)^2 \\ &= \frac{3}{2} n^2 + \frac{3}{2} n + 1 - (n + 1)^2 = \frac{n^2 - n}{2} \end{aligned}$$

- Using Chebysheff's inequality,

$$\mathbf{P}[|\Theta_n - n| \geq \epsilon \cdot n] \leq \frac{\mathbf{V}[\Theta_n]}{\epsilon^2 \cdot n^2}$$

## Analysis (2/3)

### Lemma (Second Moment Analysis)

Let  $X_n$  denote the value of  $X$  after  $n$  updates. For every  $n \geq 0$ ,

$$\mathbf{E} \left[ \left( 2^{X_n} \right)^2 \right] = \mathbf{E} \left[ 2^{2 \cdot X_n} \right] = \frac{3}{2} n^2 + \frac{3}{2} n + 1.$$

This is shown similarly to that of the previous Lemma (see supervision sheet)

- Recall  $\Theta_n = 2^{X_n} - 1$ .
- Since  $\mathbf{V}[Z] = \mathbf{E}[Z^2] - \mathbf{E}[Z]^2$ ,

$$\begin{aligned} \mathbf{V}[\Theta_n] &= \mathbf{V}[2^{X_n}] = \mathbf{E}[2^{2 \cdot X_n}] - \left( \mathbf{E}[2^{X_n}] \right)^2 \\ &= \frac{3}{2} n^2 + \frac{3}{2} n + 1 - (n+1)^2 = \frac{n^2 - n}{2} \end{aligned}$$

- Using Chebysheff's inequality,

$$\mathbf{P}[|\Theta_n - n| \geq \epsilon \cdot n] \leq \frac{\mathbf{V}[\Theta_n]}{\epsilon^2 \cdot n^2} \leq \frac{\frac{n^2}{2}}{\epsilon^2 \cdot n^2} = \frac{1}{2\epsilon^2}.$$

## Analysis (2/3)

### Lemma (Second Moment Analysis)

Let  $X_n$  denote the value of  $X$  after  $n$  updates. For every  $n \geq 0$ ,

$$\mathbf{E} \left[ \left( 2^{X_n} \right)^2 \right] = \mathbf{E} \left[ 2^{2 \cdot X_n} \right] = \frac{3}{2} n^2 + \frac{3}{2} n + 1.$$

This is shown similarly to that of the previous Lemma (see supervision sheet)

- Recall  $\Theta_n = 2^{X_n} - 1$ .
- Since  $\mathbf{V}[Z] = \mathbf{E}[Z^2] - \mathbf{E}[Z]^2$ ,

$$\begin{aligned} \mathbf{V}[\Theta_n] &= \mathbf{V}[2^{X_n}] = \mathbf{E}[2^{2 \cdot X_n}] - \left( \mathbf{E}[2^{X_n}] \right)^2 \\ &= \frac{3}{2} n^2 + \frac{3}{2} n + 1 - (n+1)^2 = \frac{n^2 - n}{2} \end{aligned}$$

- Using Chebysheff's inequality,

This failure probability (estimate) is at least  $\frac{1}{2}$  ☹️

$$\mathbf{P}[|\Theta_n - n| \geq \epsilon \cdot n] \leq \frac{\mathbf{V}[\Theta_n]}{\epsilon^2 \cdot n^2} \leq \frac{\frac{n^2}{2}}{\epsilon^2 \cdot n^2} = \frac{1}{2\epsilon^2}.$$

## Analysis (3/3)

---

**Idea:** Reduce Variance by Running Independent Instances and Taking Average.

## Analysis (3/3)

**Idea:** Reduce Variance by Running Independent Instances and Taking Average.

IMPROVED MORRIS ALGORITHM(G)

- 1: Let  $\Theta^1, \Theta^2, \dots, \Theta^k$  be  $k$  independent instances of MORRIS
- 2: **Return**  $\bar{\Theta} := \frac{1}{k} \sum_{i=1}^k \Theta^i$

**Idea:** Reduce Variance by Running Independent Instances and Taking Average.

IMPROVED MORRIS ALGORITHM( $G$ )

1: Let  $\Theta^1, \Theta^2, \dots, \Theta^k$  be  $k$  independent instances of MORRIS

2: **Return**  $\bar{\Theta} := \frac{1}{k} \sum_{i=1}^k \Theta^i$

- Clearly,  $\mathbf{E}[\bar{\Theta}] = n$ .

**Idea:** Reduce Variance by Running Independent Instances and Taking Average.

IMPROVED MORRIS ALGORITHM( $G$ )

- 1: Let  $\Theta^1, \Theta^2, \dots, \Theta^k$  be  $k$  independent instances of MORRIS
- 2: **Return**  $\bar{\Theta} := \frac{1}{k} \sum_{i=1}^k \Theta^i$

- Clearly,  $\mathbf{E} [\bar{\Theta}] = n$ . For the variance,

**Idea:** Reduce Variance by Running Independent Instances and Taking Average.

IMPROVED MORRIS ALGORITHM(G)

- 1: Let  $\Theta^1, \Theta^2, \dots, \Theta^k$  be  $k$  independent instances of MORRIS
- 2: **Return**  $\bar{\Theta} := \frac{1}{k} \sum_{i=1}^k \Theta^i$

- Clearly,  $\mathbf{E} [\bar{\Theta}] = n$ . For the **variance**,

$$\mathbf{V} [\bar{\Theta}] = \frac{1}{k^2} \cdot \mathbf{V} \left[ \sum_{i=1}^k \Theta^i \right]$$



**Idea:** Reduce Variance by Running Independent Instances and Taking Average.

IMPROVED MORRIS ALGORITHM(G)

- 1: Let  $\Theta^1, \Theta^2, \dots, \Theta^k$  be  $k$  independent instances of MORRIS
- 2: **Return**  $\bar{\Theta} := \frac{1}{k} \sum_{i=1}^k \Theta^i$

- Clearly,  $\mathbf{E} [\bar{\Theta}] = n$ . For the **variance**,

$$\mathbf{V} [\bar{\Theta}] = \frac{1}{k^2} \cdot \mathbf{V} \left[ \sum_{i=1}^k \Theta^i \right] = \frac{1}{k} \cdot \mathbf{V} [\Theta^1] \leq \frac{1}{k} \cdot \frac{n^2}{2}$$

**Idea:** Reduce Variance by Running Independent Instances and Taking Average.

IMPROVED MORRIS ALGORITHM(G)

- 1: Let  $\Theta^1, \Theta^2, \dots, \Theta^k$  be  $k$  independent instances of MORRIS
- 2: **Return**  $\bar{\Theta} := \frac{1}{k} \sum_{i=1}^k \Theta^i$

- Clearly,  $\mathbf{E} [\bar{\Theta}] = n$ . For the **variance**,

$$\mathbf{V} [\bar{\Theta}] = \frac{1}{k^2} \cdot \mathbf{V} \left[ \sum_{i=1}^k \Theta^i \right] = \frac{1}{k} \cdot \mathbf{V} [\Theta^1] \leq \frac{1}{k} \cdot \frac{n^2}{2}$$

- Hence using **Chebyshev**,

$$\mathbf{P} [ |\bar{\Theta} - n| \geq \epsilon \cdot n ] \leq \frac{1}{2k\epsilon^2}.$$

## Analysis (3/3)

**Idea:** Reduce Variance by Running Independent Instances and Taking Average.

IMPROVED MORRIS ALGORITHM(G)

- 1: Let  $\Theta^1, \Theta^2, \dots, \Theta^k$  be  $k$  independent instances of MORRIS
- 2: **Return**  $\bar{\Theta} := \frac{1}{k} \sum_{i=1}^k \Theta^i$

- Clearly,  $\mathbf{E} [\bar{\Theta}] = n$ . For the **variance**,

$$\mathbf{V} [\bar{\Theta}] = \frac{1}{k^2} \cdot \mathbf{V} \left[ \sum_{i=1}^k \Theta^i \right] = \frac{1}{k} \cdot \mathbf{V} [\Theta^1] \leq \frac{1}{k} \cdot \frac{n^2}{2}$$

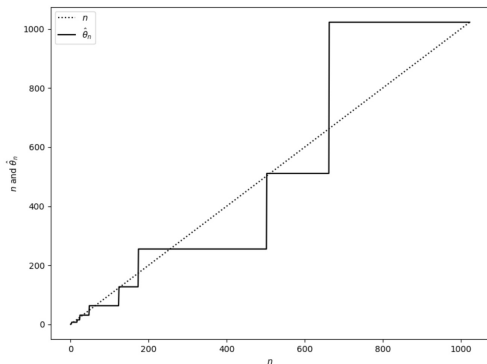
- Hence using **Chebyshev**,

$$\mathbf{P} [ |\bar{\Theta} - n| \geq \epsilon \cdot n ] \leq \frac{1}{2k\epsilon^2}.$$

### Conclusion

For any  $\epsilon, \delta < 1$ , the IMPROVED MORRIS ALG. with  $k \geq \frac{1}{2\epsilon^2\delta}$  satisfies:

$$\mathbf{P} [ |\bar{\Theta} - n| \leq \epsilon \cdot n ] \geq 1 - \delta.$$



A run of Morris's algorithm on  $n = 1024$  data points

(source: <http://gregorygundersen.com/blog/2019/11/11/morris-algorithm/>)

# Outline

---

Introduction

Approximate Counting

**Distinct Elements and Frequency Moments**

Extra Material (non-examinable): An Algorithm for  $F_0$  in the Turnstile Model

## Norm Estimation: the Alon-Matias-Szegedy algorithm

---

$F_p$ -norm (Frequency Moments)

Let  $U$  with  $|U| = n$ . For  $i \in U$ , let  $f_i$  be the number of occurrences of  $i \in U$  in the stream  $S$ .

## Norm Estimation: the Alon-Matias-Szegedy algorithm

$F_p$ -norm (Frequency Moments)

Let  $U$  with  $|U| = n$ . For  $i \in U$ , let  $f_i$  be the number of occurrences of  $i \in U$  in the stream  $\mathcal{S}$ . Then for any  $p > 0$ , the  $F_p$ -norm is defined by

$$F_p := \sum_{i \in U} f_i^p.$$

## Norm Estimation: the Alon-Matias-Szegedy algorithm

$F_p$ -norm (Frequency Moments)

Let  $U$  with  $|U| = n$ . For  $i \in U$ , let  $f_i$  be the number of occurrences of  $i \in U$  in the stream  $\mathcal{S}$ . Then for any  $p > 0$ , the  $F_p$ -norm is defined by

$$F_p := \sum_{i \in U} f_i^p.$$

- $F_1$  = total number of items in stream  $\mathcal{S}$ .
- $F_0$  = total number of distinct items in stream  $\mathcal{S}$ .



## Norm Estimation: the Alon-Matias-Szegedy algorithm

$F_p$ -norm (Frequency Moments)

Let  $U$  with  $|U| = n$ . For  $i \in U$ , let  $f_i$  be the number of occurrences of  $i \in U$  in the stream  $\mathcal{S}$ . Then for any  $p > 0$ , the  $F_p$ -norm is defined by

$$F_p := \sum_{i \in U} f_i^p.$$

- $F_1$  = total number of items in stream  $\mathcal{S}$ .
- $F_0$  = total number of distinct items in stream  $\mathcal{S}$ .

Alon, Matias, and Szegedy (1996) presented a systematical study for approximating frequency moments.

## Norm Estimation: the Alon-Matias-Szegedy algorithm

$F_p$ -norm (Frequency Moments)

Let  $U$  with  $|U| = n$ . For  $i \in U$ , let  $f_i$  be the number of occurrences of  $i \in U$  in the stream  $\mathcal{S}$ . Then for any  $p > 0$ , the  $F_p$ -norm is defined by

$$F_p := \sum_{i \in U} f_i^p.$$

- $F_1$  = total number of items in stream  $\mathcal{S}$ .
- $F_0$  = total number of distinct items in stream  $\mathcal{S}$ .

Alon, Matias, and Szegedy (1996) presented a systematical study for approximating frequency moments.

- $F_0, F_1, F_2$  can be approximated in space logarithmic in  $n$  and  $|\mathcal{S}|$ .

## Norm Estimation: the Alon-Matias-Szegedy algorithm

$F_p$ -norm (Frequency Moments)

Let  $U$  with  $|U| = n$ . For  $i \in U$ , let  $f_i$  be the number of occurrences of  $i \in U$  in the stream  $\mathcal{S}$ . Then for any  $p > 0$ , the  $F_p$ -norm is defined by

$$F_p := \sum_{i \in U} f_i^p.$$

- $F_1$  = total number of items in stream  $\mathcal{S}$ .
- $F_0$  = total number of distinct items in stream  $\mathcal{S}$ .

Alon, Matias, and Szegedy (1996) presented a systematical study for approximating frequency moments.

- $F_0, F_1, F_2$  can be approximated in space logarithmic in  $n$  and  $|\mathcal{S}|$ .
- Approximating  $F_p$  for  $p \geq 6$  requires  $n^{\Omega(1)}$  space.

## Norm Estimation: the Alon-Matias-Szegedy algorithm

$F_p$ -norm (Frequency Moments)

Let  $U$  with  $|U| = n$ . For  $i \in U$ , let  $f_i$  be the number of occurrences of  $i \in U$  in the stream  $\mathcal{S}$ . Then for any  $p > 0$ , the  $F_p$ -norm is defined by

$$F_p := \sum_{i \in U} f_i^p.$$

- $F_1$  = total number of items in stream  $\mathcal{S}$ .
- $F_0$  = total number of distinct items in stream  $\mathcal{S}$ .

Alon, Matias, and Szegedy (1996) presented a systematical study for approximating frequency moments.

- $F_0, F_1, F_2$  can be approximated in space logarithmic in  $n$  and  $|\mathcal{S}|$ .
- Approximating  $F_p$  for  $p \geq 6$  requires  $n^{\Omega(1)}$  space.
- The paper won 2005 Gödel Award for “their foundational contribution to streaming algorithms”.

## Important Tool: Pairwise independent Hash Functions

---

We will focus on the simpler case of  $F_0$ , the number of distinct elements.

## Important Tool: Pairwise independent Hash Functions

We will focus on the simpler case of  $F_0$ , the number of distinct elements.

### Pairwise Independence

A family of functions  $H = \{h \mid h : U \mapsto [n]\}$  is **pairwise independent** if, for any  $h$  chosen uniformly at random from  $H$ , the following holds:

1.  $h(x)$  is **uniformly distributed** in  $[n] = \{1, 2, \dots, n\}$  for any  $x \in U$ ;
2. For any  $x_1 \neq x_2 \in U$ ,  $h(x_1)$  and  $h(x_2)$  are **independent**.

## Important Tool: Pairwise independent Hash Functions

We will focus on the simpler case of  $F_0$ , the number of distinct elements.

### Pairwise Independence

A family of functions  $H = \{h \mid h : U \mapsto [n]\}$  is **pairwise independent** if, for any  $h$  chosen uniformly at random from  $H$ , the following holds:

1.  $h(x)$  is **uniformly distributed** in  $[n] = \{1, 2, \dots, n\}$  for any  $x \in U$ ;
2. For any  $x_1 \neq x_2 \in U$ ,  $h(x_1)$  and  $h(x_2)$  are **independent**.

### Theorem (Fact)

Let  $n$  be a **prime number**, and let  $h_{a,b}(x) = (ax + b) \bmod n$ . Define

$$H = \{h_{a,b} \mid 0 \leq a, b \leq n - 1\}.$$

Then  $H$  is a family of **pairwise independent hash functions**.

## Intuition behind the AMS algorithm

---

Assume that we have a random hash function  $h$ .



## Intuition behind the AMS algorithm

---

Assume that we have a random hash function  $h$ . Define

$$\rho(x) := \max_{i \geq 0} \left\{ i : x \bmod 2^i = 0 \right\},$$

which is the number of consecutive 0's among the lowest bits of  $x$ .

## Intuition behind the AMS algorithm

---

Assume that we have a random hash function  $h$ . Define

$$\rho(x) := \max_{i \geq 0} \left\{ i : x \bmod 2^i = 0 \right\},$$

which is the number of consecutive 0's among the lowest bits of  $x$ .

**Example:**  $\rho(2) = 1, \rho(3) = 0, \rho(4) = 2, \rho(8) = 3, \rho(16) = 4, \rho(17) = 0.$

## Intuition behind the AMS algorithm

---

Assume that we have a random hash function  $h$ . Define

$$\rho(x) := \max_{i \geq 0} \left\{ i : x \bmod 2^i = 0 \right\},$$

which is the number of consecutive 0's among the lowest bits of  $x$ .

**Example:**  $\rho(2) = 1, \rho(3) = 0, \rho(4) = 2, \rho(8) = 3, \rho(16) = 4, \rho(17) = 0$ .

**Observation.** Since  $h(x)$  is uniformly distributed over  $[n]$ , the following holds:

- with probability  $1/2$ , we have  $\rho(h(x)) \geq 1$

## Intuition behind the AMS algorithm

Assume that we have a **random hash function**  $h$ . Define

$$\rho(x) := \max_{i \geq 0} \left\{ i : x \bmod 2^i = 0 \right\},$$

which is the **number of consecutive 0's** among the lowest bits of  $x$ .

**Example:**  $\rho(2) = 1, \rho(3) = 0, \rho(4) = 2, \rho(8) = 3, \rho(16) = 4, \rho(17) = 0$ .

**Observation.** Since  $h(x)$  is uniformly distributed over  $[n]$ , the following holds:

- with probability  $1/2$ , we have  $\rho(h(x)) \geq 1$
- with probability  $1/4$ , we have  $\rho(h(x)) \geq 2$
- with probability  $1/8$ , we have  $\rho(h(x)) \geq 3$

## Intuition behind the AMS algorithm

Assume that we have a random hash function  $h$ . Define

$$\rho(x) := \max_{i \geq 0} \left\{ i : x \bmod 2^i = 0 \right\},$$

which is the number of consecutive 0's among the lowest bits of  $x$ .

**Example:**  $\rho(2) = 1, \rho(3) = 0, \rho(4) = 2, \rho(8) = 3, \rho(16) = 4, \rho(17) = 0$ .

**Observation.** Since  $h(x)$  is uniformly distributed over  $[n]$ , the following holds:

- with probability  $1/2$ , we have  $\rho(h(x)) \geq 1$
- with probability  $1/4$ , we have  $\rho(h(x)) \geq 2$
- with probability  $1/8$ , we have  $\rho(h(x)) \geq 3$
- $\vdots$
- with probability  $1/2^r$ , we have  $\rho(h(x)) \geq r$

## Intuition behind the AMS algorithm

Assume that we have a random hash function  $h$ . Define

$$\rho(x) := \max_{i \geq 0} \left\{ i : x \bmod 2^i = 0 \right\},$$

which is the number of consecutive 0's among the lowest bits of  $x$ .

**Example:**  $\rho(2) = 1, \rho(3) = 0, \rho(4) = 2, \rho(8) = 3, \rho(16) = 4, \rho(17) = 0$ .

**Observation.** Since  $h(x)$  is uniformly distributed over  $[n]$ , the following holds:

- with probability  $1/2$ , we have  $\rho(h(x)) \geq 1$
- with probability  $1/4$ , we have  $\rho(h(x)) \geq 2$
- with probability  $1/8$ , we have  $\rho(h(x)) \geq 3$
- $\vdots$
- with probability  $1/2^r$ , we have  $\rho(h(x)) \geq r$

Since  $n$  is not a power of 2, this probability is in fact equal to  $\frac{\lfloor n/2^r \rfloor}{n} \approx 1/2^r - o(1)$ .

## The AMS Algorithm

---

### AMS ALGORITHM

- 1: Choose a random hash function  $h : [n] \rightarrow [n]$
- 2:  $Z \leftarrow 0$
- 3: **while** item  $x$  from stream  $\mathcal{S}$  arrives
- 4:     **if**  $\rho(h(x)) > Z$  **then**  $Z \leftarrow \rho(h(x))$
- 5: **return**  $2^{Z+1/2}$

## The AMS Algorithm

---

### AMS ALGORITHM

- 1: Choose a random hash function  $h : [n] \rightarrow [n]$
- 2:  $Z \leftarrow 0$
- 3: **while** item  $x$  from stream  $S$  arrives
- 4:     **if**  $\rho(h(x)) > Z$  **then**  $Z \leftarrow \rho(h(x))$
- 5: **return**  $2^{Z+1/2}$

$Z \leftarrow \max\{Z, \rho(h(x))\}$



# The AMS Algorithm

## AMS ALGORITHM

- 1: Choose a random hash function  $h : [n] \rightarrow [n]$
- 2:  $Z \leftarrow 0$
- 3: **while** item  $x$  from stream  $S$  arrives
- 4:     **if**  $\rho(h(x)) > Z$  **then**  $Z \leftarrow \rho(h(x))$
- 5: **return**  $2^{Z+1/2}$

$Z \leftarrow \max\{Z, \rho(h(x))\}$

### Analysis of AMS Algorithm

With constant probability  $> 0$ , the algorithm's output satisfies

$$2^{Z+1/2} \in [F_0/3, 3 \cdot F_0].$$

## The AMS Algorithm

### AMS ALGORITHM

- 1: Choose a random hash function  $h : [n] \rightarrow [n]$
- 2:  $Z \leftarrow 0$
- 3: **while** item  $x$  from stream  $S$  arrives
- 4:     **if**  $\rho(h(x)) > Z$  **then**  $Z \leftarrow \rho(h(x))$
- 5: **return**  $2^{Z+1/2}$

$$Z \leftarrow \max\{Z, \rho(h(x))\}$$

### Analysis of AMS Algorithm

With constant probability  $> 0$ , the algorithm's output satisfies

$$2^{Z+1/2} \in [F_0/3, 3 \cdot F_0].$$

We get an  $(O(1), \delta)$ -approximation of  $F_0$  by running  $\Theta(\log(1/\delta))$  independent copies of the algorithm and returning the **median**.

## The AMS Algorithm

### AMS ALGORITHM

- 1: Choose a random hash function  $h : [n] \rightarrow [n]$
- 2:  $Z \leftarrow 0$
- 3: **while** item  $x$  from stream  $S$  arrives
- 4:     **if**  $\rho(h(x)) > Z$  **then**  $Z \leftarrow \rho(h(x))$
- 5: **return**  $2^{Z+1/2}$

$$Z \leftarrow \max\{Z, \rho(h(x))\}$$

### Analysis of AMS Algorithm

With constant probability  $> 0$ , the algorithm's output satisfies

$$2^{Z+1/2} \in [F_0/3, 3 \cdot F_0].$$

We get an  $(O(1), \delta)$ -approximation of  $F_0$  by running  $\Theta(\log(1/\delta))$  independent copies of the algorithm and returning the **median**.

Recall  $(\epsilon, \delta)$ -approximation:

$$\mathbf{P}[\text{Output} \in (1 - \epsilon, 1 + \epsilon) \cdot \text{Exact}] \geq 1 - \delta$$

## Example of the AMS Algorithm

---

- Assume  $n = 101$  (which is prime)
- The hash function is  $h(x) = (ax + b) \bmod n$  with  $a = 28$ ,  $b = 16$

## Example of the AMS Algorithm

---

- Assume  $n = 101$  (which is prime)
- The hash function is  $h(x) = (ax + b) \bmod n$  with  $a = 28$ ,  $b = 16$
- The data stream is:

$$S = (25, 76, 14, 51, 25, 14, 76, 76, 3, 51, 96, 14, 67, 3, 15, 25, 2, 76, 14, 71)$$

- $F_0 = 10$ , as the following numbers appeared:  $\{2, 3, 14, 15, 25, 51, 67, 71, 76, 96\}$

## Example of the AMS Algorithm

- Assume  $n = 101$  (which is prime)
- The hash function is  $h(x) = (ax + b) \bmod n$  with  $a = 28$ ,  $b = 16$
- The data stream is:

$$S = (25, 76, 14, 51, 25, 14, 76, 76, 3, 51, 96, 14, 67, 3, 15, 25, 2, 76, 14, 71)$$

- $F_0 = 10$ , as the following numbers appeared:  $\{2, 3, 14, 15, 25, 51, 67, 71, 76, 96\}$

$x$	$h(x)$	Binary Representation							$\rho(h(x))$
2	72	1	0	0	1	0	0	0	3
3	100	1	1	0	0	1	0	0	2
14	4	0	0	0	0	1	0	0	2
15	32	0	1	0	0	0	0	0	5
25	9	0	0	0	1	0	0	1	0
51	30	0	0	1	1	1	1	0	1
67	74	1	0	0	1	0	1	0	1
71	85	1	0	1	0	1	0	1	0
76	23	0	0	1	0	1	1	1	0
96	78	1	0	0	1	1	1	0	1

## Example of the AMS Algorithm

- Assume  $n = 101$  (which is prime)
- The hash function is  $h(x) = (ax + b) \bmod n$  with  $a = 28$ ,  $b = 16$
- The data stream is:

$$S = (25, 76, 14, 51, 25, 14, 76, 76, 3, 51, 96, 14, 67, 3, 15, 25, 2, 76, 14, 71)$$

- $F_0 = 10$ , as the following numbers appeared:  $\{2, 3, 14, 15, 25, 51, 67, 71, 76, 96\}$

$x$	$h(x)$	Binary Representation							$\rho(h(x))$
2	72	1	0	0	1	0	0	0	3
3	100	1	1	0	0	1	0	0	2
14	4	0	0	0	0	1	0	0	2
15	32	0	1	0	0	0	0	0	5
25	9	0	0	0	1	0	0	1	0
51	30	0	0	1	1	1	1	0	1
67	74	1	0	0	1	0	1	0	1
71	85	1	0	1	0	1	0	1	0
76	23	0	0	1	0	1	1	1	0
96	78	1	0	0	1	1	1	0	1

returned estimate:  
 $2^{5+1/2} \approx 45.25$

## Analysis (1/2)

---

Let  $X_{r,j}$  be a 0/1 indicator random variable such that

$$X_{r,j} = 1 \Leftrightarrow \rho(h(j)) \geq r.$$

We say item  $j$  reaches level  $r$  if  $X_{r,j} = 1$ .



## Analysis (1/2)

---

Let  $X_{r,j}$  be a 0/1 indicator random variable such that

$$X_{r,j} = 1 \Leftrightarrow \rho(h(j)) \geq r.$$

We say item  $j$  reaches level  $r$  if  $X_{r,j} = 1$ .

Let  $Y_r = \sum_{j \in \mathcal{S}} X_{r,j}$  be the number of items  $j$  reaching level  $r$ .

## Analysis (1/2)

---

Let  $X_{r,j}$  be a 0/1 indicator random variable such that

$$X_{r,j} = 1 \Leftrightarrow \rho(h(j)) \geq r.$$

We say item  $j$  reaches level  $r$  if  $X_{r,j} = 1$ .

Let  $Y_r = \sum_{j \in \mathcal{S}} X_{r,j}$  be the number of items  $j$  reaching level  $r$ .

Using that  $h(j)$  is uniformly distributed, we conclude

$$\mathbf{E}[X_{r,j}] = \mathbf{P}[\rho(h(j)) \geq r] = \mathbf{P}[h(j) \bmod 2^r = 0] = 2^{-r}.$$

definition of function  $\rho$

## Analysis (1/2)

---

Let  $X_{r,j}$  be a 0/1 indicator random variable such that

$$X_{r,j} = 1 \Leftrightarrow \rho(h(j)) \geq r.$$

We say item  $j$  reaches level  $r$  if  $X_{r,j} = 1$ .

Let  $Y_r = \sum_{j \in \mathcal{S}} X_{r,j}$  be the number of items  $j$  reaching level  $r$ .

Using that  $h(j)$  is **uniformly distributed**, we conclude

$$\mathbf{E}[X_{r,j}] = \mathbf{P}[\rho(h(j)) \geq r] = \mathbf{P}[h(j) \bmod 2^r = 0] = 2^{-r}.$$

definition of function  $\rho$

By **linearity of expectation**, we have

$$\mathbf{E}[Y_r] = \sum_{j \in \mathcal{S}} \mathbf{E}[X_{r,j}] = \frac{F_0}{2^r},$$

## Analysis (1/2)

---

Let  $X_{r,j}$  be a 0/1 indicator random variable such that

$$X_{r,j} = 1 \Leftrightarrow \rho(h(j)) \geq r.$$

We say item  $j$  reaches level  $r$  if  $X_{r,j} = 1$ .

Let  $Y_r = \sum_{j \in \mathcal{S}} X_{r,j}$  be the number of items  $j$  reaching level  $r$ .

Using that  $h(j)$  is **uniformly distributed**, we conclude

$$\mathbf{E}[X_{r,j}] = \mathbf{P}[\rho(h(j)) \geq r] = \mathbf{P}[h(j) \bmod 2^r = 0] = 2^{-r}.$$

definition of function  $\rho$

By **linearity of expectation**, we have

$$\mathbf{E}[Y_r] = \sum_{j \in \mathcal{S}} \mathbf{E}[X_{r,j}] = \frac{F_0}{2^r},$$

$$\mathbf{V}[Y_r] = \sum_{j \in \mathcal{S}} \mathbf{V}[X_{r,j}]$$

## Analysis (1/2)

Let  $X_{r,j}$  be a 0/1 indicator random variable such that

$$X_{r,j} = 1 \Leftrightarrow \rho(h(j)) \geq r.$$

We say item  $j$  reaches level  $r$  if  $X_{r,j} = 1$ .

Let  $Y_r = \sum_{j \in \mathcal{S}} X_{r,j}$  be the number of items  $j$  reaching level  $r$ .

Using that  $h(j)$  is **uniformly distributed**, we conclude

$$\mathbf{E}[X_{r,j}] = \mathbf{P}[\rho(h(j)) \geq r] = \mathbf{P}[h(j) \bmod 2^r = 0] = 2^{-r}.$$

definition of function  $\rho$

By **linearity of expectation**, we have

$$\mathbf{E}[Y_r] = \sum_{j \in \mathcal{S}} \mathbf{E}[X_{r,j}] = \frac{F_0}{2^r},$$

$$\mathbf{V}[Y_r] = \sum_{j \in \mathcal{S}} \mathbf{V}[X_{r,j}]$$

using **pairwise independence of  $h$** !

## Analysis (1/2)

Let  $X_{r,j}$  be a 0/1 indicator random variable such that

$$X_{r,j} = 1 \Leftrightarrow \rho(h(j)) \geq r.$$

We say item  $j$  reaches level  $r$  if  $X_{r,j} = 1$ .

Let  $Y_r = \sum_{j \in \mathcal{S}} X_{r,j}$  be the number of items  $j$  reaching level  $r$ .

Using that  $h(j)$  is **uniformly distributed**, we conclude

$$\mathbf{E}[X_{r,j}] = \mathbf{P}[\rho(h(j)) \geq r] = \mathbf{P}[h(j) \bmod 2^r = 0] = 2^{-r}.$$

definition of function  $\rho$

By **linearity of expectation**, we have

$$\mathbf{E}[Y_r] = \sum_{j \in \mathcal{S}} \mathbf{E}[X_{r,j}] = \frac{F_0}{2^r},$$

$$\mathbf{V}[Y_r] = \sum_{j \in \mathcal{S}} \mathbf{V}[X_{r,j}] \leq \sum_{j \in \mathcal{S}} \mathbf{E}[X_{r,j}^2]$$

using **pairwise independence of  $h$** !

## Analysis (1/2)

Let  $X_{r,j}$  be a 0/1 indicator random variable such that

$$X_{r,j} = 1 \Leftrightarrow \rho(h(j)) \geq r.$$

We say item  $j$  reaches level  $r$  if  $X_{r,j} = 1$ .

Let  $Y_r = \sum_{j \in \mathcal{S}} X_{r,j}$  be the number of items  $j$  reaching level  $r$ .

Using that  $h(j)$  is **uniformly distributed**, we conclude

$$\mathbf{E}[X_{r,j}] = \mathbf{P}[\rho(h(j)) \geq r] = \mathbf{P}[h(j) \bmod 2^r = 0] = 2^{-r}.$$

definition of function  $\rho$

By **linearity of expectation**, we have

$$\mathbf{E}[Y_r] = \sum_{j \in \mathcal{S}} \mathbf{E}[X_{r,j}] = \frac{F_0}{2^r},$$

$$\mathbf{V}[Y_r] = \sum_{j \in \mathcal{S}} \mathbf{V}[X_{r,j}] \leq \sum_{j \in \mathcal{S}} \mathbf{E}[X_{r,j}^2] = \sum_{j \in \mathcal{S}} \mathbf{E}[X_{r,j}] = \frac{F_0}{2^r}$$

using **pairwise independence of  $h$** !

## Analysis (2/2)

---

We have proved  $\mathbf{E}[Y_r] = \frac{F_0}{2r}$  and  $\mathbf{V}[Y_r] \leq \frac{F_0}{2r}$ .



## Analysis (2/2)

---

We have proved  $\mathbf{E}[Y_r] = \frac{F_0}{2^r}$  and  $\mathbf{V}[Y_r] \leq \frac{F_0}{2^r}$ .

By [Markov's inequality](#), we have

$$\mathbf{P}[Y_r > 0] = \mathbf{P}[Y_r \geq 1] \leq \frac{\mathbf{E}[Y_r]}{1} = \frac{F_0}{2^r}.$$

## Analysis (2/2)

---

We have proved  $\mathbf{E}[Y_r] = \frac{F_0}{2^r}$  and  $\mathbf{V}[Y_r] \leq \frac{F_0}{2^r}$ .

By **Markov's inequality**, we have

$$\mathbf{P}[Y_r > 0] = \mathbf{P}[Y_r \geq 1] \leq \frac{\mathbf{E}[Y_r]}{1} = \frac{F_0}{2^r}.$$

By **Chebyshev's inequality**, we have

$$\mathbf{P}[Y_r = 0] \leq \mathbf{P}[|Y_r - \mathbf{E}[Y_r]| \geq F_0/2^r] \leq \frac{\mathbf{V}[Y_r]}{(F_0/2^r)^2} \leq \frac{2^r}{F_0}.$$

## Analysis (2/2)

---

We have proved  $\mathbf{E}[Y_r] = \frac{F_0}{2^r}$  and  $\mathbf{V}[Y_r] \leq \frac{F_0}{2^r}$ .

By **Markov's inequality**, we have

$$\mathbf{P}[Y_r > 0] = \mathbf{P}[Y_r \geq 1] \leq \frac{\mathbf{E}[Y_r]}{1} = \frac{F_0}{2^r}.$$

By **Chebyshev's inequality**, we have

$$\mathbf{P}[Y_r = 0] \leq \mathbf{P}[|Y_r - \mathbf{E}[Y_r]| \geq F_0/2^r] \leq \frac{\mathbf{V}[Y_r]}{(F_0/2^r)^2} \leq \frac{2^r}{F_0}.$$

Let  $Z$  be the **final integer** the algo. keeps. So the algo. returns  $2^{Z+1/2}$ .

## Analysis (2/2)

---

We have proved  $\mathbf{E}[Y_r] = \frac{F_0}{2^r}$  and  $\mathbf{V}[Y_r] \leq \frac{F_0}{2^r}$ .

By **Markov's inequality**, we have

$$\mathbf{P}[Y_r > 0] = \mathbf{P}[Y_r \geq 1] \leq \frac{\mathbf{E}[Y_r]}{1} = \frac{F_0}{2^r}.$$

By **Chebyshev's inequality**, we have

$$\mathbf{P}[Y_r = 0] \leq \mathbf{P}[|Y_r - \mathbf{E}[Y_r]| \geq F_0/2^r] \leq \frac{\mathbf{V}[Y_r]}{(F_0/2^r)^2} \leq \frac{2^r}{F_0}.$$

Let  $Z$  be the **final integer** the algo. keeps. So the algo. returns  $2^{Z+1/2}$ .

Let  $p$  be the **smallest integer** such that  $2^{p+1/2} \geq 3F_0$ :

## Analysis (2/2)

---

We have proved  $\mathbf{E}[Y_r] = \frac{F_0}{2^r}$  and  $\mathbf{V}[Y_r] \leq \frac{F_0}{2^r}$ .

By **Markov's inequality**, we have

$$\mathbf{P}[Y_r > 0] = \mathbf{P}[Y_r \geq 1] \leq \frac{\mathbf{E}[Y_r]}{1} = \frac{F_0}{2^r}.$$

By **Chebyshev's inequality**, we have

$$\mathbf{P}[Y_r = 0] \leq \mathbf{P}[|Y_r - \mathbf{E}[Y_r]| \geq F_0/2^r] \leq \frac{\mathbf{V}[Y_r]}{(F_0/2^r)^2} \leq \frac{2^r}{F_0}.$$

Let  $Z$  be the **final integer** the algo. keeps. So the algo. returns  $2^{Z+1/2}$ .

Let  $p$  be the **smallest integer** such that  $2^{p+1/2} \geq 3F_0$ :

$$\mathbf{P}\left[2^{Z+1/2} \geq 3F_0\right]$$

## Analysis (2/2)

---

We have proved  $\mathbf{E}[Y_r] = \frac{F_0}{2^r}$  and  $\mathbf{V}[Y_r] \leq \frac{F_0}{2^r}$ .

By **Markov's inequality**, we have

$$\mathbf{P}[Y_r > 0] = \mathbf{P}[Y_r \geq 1] \leq \frac{\mathbf{E}[Y_r]}{1} = \frac{F_0}{2^r}.$$

By **Chebyshev's inequality**, we have

$$\mathbf{P}[Y_r = 0] \leq \mathbf{P}[|Y_r - \mathbf{E}[Y_r]| \geq F_0/2^r] \leq \frac{\mathbf{V}[Y_r]}{(F_0/2^r)^2} \leq \frac{2^r}{F_0}.$$

Let  $Z$  be the **final integer** the algo. keeps. So the algo. returns  $2^{Z+1/2}$ .

Let  $p$  be the **smallest integer** such that  $2^{p+1/2} \geq 3F_0$ :

$$\mathbf{P}\left[2^{Z+1/2} \geq 3F_0\right] = \mathbf{P}[Z \geq p]$$

## Analysis (2/2)

---

We have proved  $\mathbf{E}[Y_r] = \frac{F_0}{2^r}$  and  $\mathbf{V}[Y_r] \leq \frac{F_0}{2^r}$ .

By **Markov's inequality**, we have

$$\mathbf{P}[Y_r > 0] = \mathbf{P}[Y_r \geq 1] \leq \frac{\mathbf{E}[Y_r]}{1} = \frac{F_0}{2^r}.$$

By **Chebyshev's inequality**, we have

$$\mathbf{P}[Y_r = 0] \leq \mathbf{P}[|Y_r - \mathbf{E}[Y_r]| \geq F_0/2^r] \leq \frac{\mathbf{V}[Y_r]}{(F_0/2^r)^2} \leq \frac{2^r}{F_0}.$$

Let  $Z$  be the **final integer** the algo. keeps. So the algo. returns  $2^{Z+1/2}$ .

Let  $p$  be the **smallest integer** such that  $2^{p+1/2} \geq 3F_0$ :

$$\mathbf{P}\left[2^{Z+1/2} \geq 3F_0\right] = \mathbf{P}[Z \geq p] = \mathbf{P}[Y_p > 0]$$

## Analysis (2/2)

---

We have proved  $\mathbf{E}[Y_r] = \frac{F_0}{2^r}$  and  $\mathbf{V}[Y_r] \leq \frac{F_0}{2^r}$ .

By **Markov's inequality**, we have

$$\mathbf{P}[Y_r > 0] = \mathbf{P}[Y_r \geq 1] \leq \frac{\mathbf{E}[Y_r]}{1} = \frac{F_0}{2^r}.$$

By **Chebyshev's inequality**, we have

$$\mathbf{P}[Y_r = 0] \leq \mathbf{P}[|Y_r - \mathbf{E}[Y_r]| \geq F_0/2^r] \leq \frac{\mathbf{V}[Y_r]}{(F_0/2^r)^2} \leq \frac{2^r}{F_0}.$$

Let  $Z$  be the **final integer** the algo. keeps. So the algo. returns  $2^{Z+1/2}$ .

Let  $p$  be the **smallest integer** such that  $2^{p+1/2} \geq 3F_0$ :

$$\mathbf{P}\left[2^{Z+1/2} \geq 3F_0\right] = \mathbf{P}[Z \geq p] = \mathbf{P}[Y_p > 0] \leq \frac{F_0}{2^p}$$



## Analysis (2/2)

---

We have proved  $\mathbf{E}[Y_r] = \frac{F_0}{2^r}$  and  $\mathbf{V}[Y_r] \leq \frac{F_0}{2^r}$ .

By **Markov's inequality**, we have

$$\mathbf{P}[Y_r > 0] = \mathbf{P}[Y_r \geq 1] \leq \frac{\mathbf{E}[Y_r]}{1} = \frac{F_0}{2^r}.$$

By **Chebyshev's inequality**, we have

$$\mathbf{P}[Y_r = 0] \leq \mathbf{P}[|Y_r - \mathbf{E}[Y_r]| \geq F_0/2^r] \leq \frac{\mathbf{V}[Y_r]}{(F_0/2^r)^2} \leq \frac{2^r}{F_0}.$$

Let  $Z$  be the **final integer** the algo. keeps. So the algo. returns  $2^{Z+1/2}$ .

Let  $p$  be the **smallest integer** such that  $2^{p+1/2} \geq 3F_0$ :

$$\mathbf{P}\left[2^{Z+1/2} \geq 3F_0\right] = \mathbf{P}[Z \geq p] = \mathbf{P}[Y_p > 0] \leq \frac{F_0}{2^p} \leq \frac{\sqrt{2}}{3}.$$

## Analysis (2/2)

---

We have proved  $\mathbf{E}[Y_r] = \frac{F_0}{2^r}$  and  $\mathbf{V}[Y_r] \leq \frac{F_0}{2^r}$ .

By **Markov's inequality**, we have

$$\mathbf{P}[Y_r > 0] = \mathbf{P}[Y_r \geq 1] \leq \frac{\mathbf{E}[Y_r]}{1} = \frac{F_0}{2^r}.$$

By **Chebyshev's inequality**, we have

$$\mathbf{P}[Y_r = 0] \leq \mathbf{P}[|Y_r - \mathbf{E}[Y_r]| \geq F_0/2^r] \leq \frac{\mathbf{V}[Y_r]}{(F_0/2^r)^2} \leq \frac{2^r}{F_0}.$$

Let  $Z$  be the **final integer** the algo. keeps. So the algo. returns  $2^{Z+1/2}$ .

Let  $p$  be the **smallest integer** such that  $2^{p+1/2} \geq 3F_0$ :

$$\mathbf{P}\left[2^{Z+1/2} \geq 3F_0\right] = \mathbf{P}[Z \geq p] = \mathbf{P}[Y_p > 0] \leq \frac{F_0}{2^p} \leq \frac{\sqrt{2}}{3}.$$

Let  $q$  be the **largest integer** such that  $2^{q+1/2} \leq F_0/3$ :

## Analysis (2/2)

---

We have proved  $\mathbf{E}[Y_r] = \frac{F_0}{2^r}$  and  $\mathbf{V}[Y_r] \leq \frac{F_0}{2^r}$ .

By **Markov's inequality**, we have

$$\mathbf{P}[Y_r > 0] = \mathbf{P}[Y_r \geq 1] \leq \frac{\mathbf{E}[Y_r]}{1} = \frac{F_0}{2^r}.$$

By **Chebyshev's inequality**, we have

$$\mathbf{P}[Y_r = 0] \leq \mathbf{P}[|Y_r - \mathbf{E}[Y_r]| \geq F_0/2^r] \leq \frac{\mathbf{V}[Y_r]}{(F_0/2^r)^2} \leq \frac{2^r}{F_0}.$$

Let  $Z$  be the **final integer** the algo. keeps. So the algo. returns  $2^{Z+1/2}$ .

Let  $p$  be the **smallest integer** such that  $2^{p+1/2} \geq 3F_0$ :

$$\mathbf{P}\left[2^{Z+1/2} \geq 3F_0\right] = \mathbf{P}[Z \geq p] = \mathbf{P}[Y_p > 0] \leq \frac{F_0}{2^p} \leq \frac{\sqrt{2}}{3}.$$

Let  $q$  be the **largest integer** such that  $2^{q+1/2} \leq F_0/3$ :

$$\mathbf{P}\left[2^{Z+1/2} \leq F_0/3\right]$$

## Analysis (2/2)

---

We have proved  $\mathbf{E}[Y_r] = \frac{F_0}{2^r}$  and  $\mathbf{V}[Y_r] \leq \frac{F_0}{2^r}$ .

By **Markov's inequality**, we have

$$\mathbf{P}[Y_r > 0] = \mathbf{P}[Y_r \geq 1] \leq \frac{\mathbf{E}[Y_r]}{1} = \frac{F_0}{2^r}.$$

By **Chebyshev's inequality**, we have

$$\mathbf{P}[Y_r = 0] \leq \mathbf{P}[|Y_r - \mathbf{E}[Y_r]| \geq F_0/2^r] \leq \frac{\mathbf{V}[Y_r]}{(F_0/2^r)^2} \leq \frac{2^r}{F_0}.$$

Let  $Z$  be the **final integer** the algo. keeps. So the algo. returns  $2^{Z+1/2}$ .

Let  $p$  be the **smallest integer** such that  $2^{p+1/2} \geq 3F_0$ :

$$\mathbf{P}\left[2^{Z+1/2} \geq 3F_0\right] = \mathbf{P}[Z \geq p] = \mathbf{P}[Y_p > 0] \leq \frac{F_0}{2^p} \leq \frac{\sqrt{2}}{3}.$$

Let  $q$  be the **largest integer** such that  $2^{q+1/2} \leq F_0/3$ :

$$\mathbf{P}\left[2^{Z+1/2} \leq F_0/3\right] = \mathbf{P}[Z \leq q]$$

## Analysis (2/2)

---

We have proved  $\mathbf{E}[Y_r] = \frac{F_0}{2^r}$  and  $\mathbf{V}[Y_r] \leq \frac{F_0}{2^r}$ .

By **Markov's inequality**, we have

$$\mathbf{P}[Y_r > 0] = \mathbf{P}[Y_r \geq 1] \leq \frac{\mathbf{E}[Y_r]}{1} = \frac{F_0}{2^r}.$$

By **Chebyshev's inequality**, we have

$$\mathbf{P}[Y_r = 0] \leq \mathbf{P}[|Y_r - \mathbf{E}[Y_r]| \geq F_0/2^r] \leq \frac{\mathbf{V}[Y_r]}{(F_0/2^r)^2} \leq \frac{2^r}{F_0}.$$

Let  $Z$  be the **final integer** the algo. keeps. So the algo. returns  $2^{Z+1/2}$ .

Let  $p$  be the **smallest integer** such that  $2^{p+1/2} \geq 3F_0$ :

$$\mathbf{P}\left[2^{Z+1/2} \geq 3F_0\right] = \mathbf{P}[Z \geq p] = \mathbf{P}[Y_p > 0] \leq \frac{F_0}{2^p} \leq \frac{\sqrt{2}}{3}.$$

Let  $q$  be the **largest integer** such that  $2^{q+1/2} \leq F_0/3$ :

$$\mathbf{P}\left[2^{Z+1/2} \leq F_0/3\right] = \mathbf{P}[Z \leq q] \leq \mathbf{P}[Y_{q+1} = 0]$$

## Analysis (2/2)

---

We have proved  $\mathbf{E}[Y_r] = \frac{F_0}{2^r}$  and  $\mathbf{V}[Y_r] \leq \frac{F_0}{2^r}$ .

By **Markov's inequality**, we have

$$\mathbf{P}[Y_r > 0] = \mathbf{P}[Y_r \geq 1] \leq \frac{\mathbf{E}[Y_r]}{1} = \frac{F_0}{2^r}.$$

By **Chebyshev's inequality**, we have

$$\mathbf{P}[Y_r = 0] \leq \mathbf{P}[|Y_r - \mathbf{E}[Y_r]| \geq F_0/2^r] \leq \frac{\mathbf{V}[Y_r]}{(F_0/2^r)^2} \leq \frac{2^r}{F_0}.$$

Let  $Z$  be the **final integer** the algo. keeps. So the algo. returns  $2^{Z+1/2}$ .

Let  $p$  be the **smallest integer** such that  $2^{p+1/2} \geq 3F_0$ :

$$\mathbf{P}[2^{Z+1/2} \geq 3F_0] = \mathbf{P}[Z \geq p] = \mathbf{P}[Y_p > 0] \leq \frac{F_0}{2^p} \leq \frac{\sqrt{2}}{3}.$$

Let  $q$  be the **largest integer** such that  $2^{q+1/2} \leq F_0/3$ :

$$\mathbf{P}[2^{Z+1/2} \leq F_0/3] = \mathbf{P}[Z \leq q] \leq \mathbf{P}[Y_{q+1} = 0] \leq \frac{2^{q+1}}{F_0}$$

## Analysis (2/2)

---

We have proved  $\mathbf{E}[Y_r] = \frac{F_0}{2^r}$  and  $\mathbf{V}[Y_r] \leq \frac{F_0}{2^r}$ .

By **Markov's inequality**, we have

$$\mathbf{P}[Y_r > 0] = \mathbf{P}[Y_r \geq 1] \leq \frac{\mathbf{E}[Y_r]}{1} = \frac{F_0}{2^r}.$$

By **Chebyshev's inequality**, we have

$$\mathbf{P}[Y_r = 0] \leq \mathbf{P}[|Y_r - \mathbf{E}[Y_r]| \geq F_0/2^r] \leq \frac{\mathbf{V}[Y_r]}{(F_0/2^r)^2} \leq \frac{2^r}{F_0}.$$

Let  $Z$  be the **final integer** the algo. keeps. So the algo. returns  $2^{Z+1/2}$ .

Let  $p$  be the **smallest integer** such that  $2^{p+1/2} \geq 3F_0$ :

$$\mathbf{P}[2^{Z+1/2} \geq 3F_0] = \mathbf{P}[Z \geq p] = \mathbf{P}[Y_p > 0] \leq \frac{F_0}{2^p} \leq \frac{\sqrt{2}}{3}.$$

Let  $q$  be the **largest integer** such that  $2^{q+1/2} \leq F_0/3$ :

$$\mathbf{P}[2^{Z+1/2} \leq F_0/3] = \mathbf{P}[Z \leq q] \leq \mathbf{P}[Y_{q+1} = 0] \leq \frac{2^{q+1}}{F_0} \leq \frac{\sqrt{2}}{3}.$$

## Analysis (2/2)

We have proved  $\mathbf{E}[Y_r] = \frac{F_0}{2^r}$  and  $\mathbf{V}[Y_r] \leq \frac{F_0}{2^r}$ .

By **Markov's inequality**, we have

$$\mathbf{P}[Y_r > 0] = \mathbf{P}[Y_r \geq 1] \leq \frac{\mathbf{E}[Y_r]}{1} = \frac{F_0}{2^r}.$$

By **Chebyshev's inequality**, we have

$$\mathbf{P}[Y_r = 0] \leq \mathbf{P}[|Y_r - \mathbf{E}[Y_r]| \geq F_0/2^r] \leq \frac{\mathbf{V}[Y_r]}{(F_0/2^r)^2} \leq \frac{2^r}{F_0}.$$

Let  $Z$  be the **final integer** the algo. keeps. So the algo. returns  $2^{Z+1/2}$ .

Let  $p$  be the **smallest integer** such that  $2^{p+1/2} \geq 3F_0$ :

$$\mathbf{P}[2^{Z+1/2} \geq 3F_0] = \mathbf{P}[Z \geq p] = \mathbf{P}[Y_p > 0] \leq \frac{F_0}{2^p} \leq \frac{\sqrt{2}}{3}.$$

Let  $q$  be the **largest integer** such that  $2^{q+1/2} \leq F_0/3$ :

Union Bound: Error  $\leq 2 \cdot \frac{\sqrt{2}}{3} < 1$

$$\mathbf{P}[2^{Z+1/2} \leq F_0/3] = \mathbf{P}[Z \leq q] \leq \mathbf{P}[Y_{q+1} = 0] \leq \frac{2^{q+1}}{F_0} \leq \frac{\sqrt{2}}{3}.$$



## Analysis (2/2)

We have proved  $\mathbf{E}[Y_r] = \frac{F_0}{2^r}$  and  $\mathbf{V}[Y_r] \leq \frac{F_0}{2^r}$ .

By **Markov's inequality**, we have

$$\mathbf{P}[Y_r > 0] = \mathbf{P}[Y_r \geq 1] \leq \frac{\mathbf{E}[Y_r]}{1} = \frac{F_0}{2^r}.$$

By **Chebyshev's inequality**, we have

$$\mathbf{P}[Y_r = 0] \leq \mathbf{P}[|Y_r - \mathbf{E}[Y_r]| \geq F_0/2^r] \leq \frac{\mathbf{V}[Y_r]}{(F_0/2^r)^2} \leq \frac{2^r}{F_0}.$$

Let  $Z$  be the **final integer** the algo. keeps. So the algo. returns  $2^{Z+1/2}$ .

Let  $p$  be the **smallest integer** such that  $2^{p+1/2} \geq 3F_0$ :

$$\mathbf{P}[2^{Z+1/2} \geq 3F_0] = \mathbf{P}[Z \geq p] = \mathbf{P}[Y_p > 0] \leq \frac{F_0}{2^p} \leq \frac{\sqrt{2}}{3}.$$

Let  $q$  be the **largest integer** such that  $2^{q+1/2} \leq F_0/3$ :

Union Bound: Error  $\leq 2 \cdot \frac{\sqrt{2}}{3} < 1$

$$\mathbf{P}[2^{Z+1/2} \leq F_0/3] = \mathbf{P}[Z \leq q] \leq \mathbf{P}[Y_{q+1} = 0] \leq \frac{2^{q+1}}{F_0} \leq \frac{\sqrt{2}}{3}. \quad \square$$

## Final Remarks

---

- Durand and Flajolet (2003) proposed the LOGLOG algorithm for estimating  $F_0$
- Their algorithm condenses the [whole of Shakespeare's works](#) to a table of 256 “small bytes” of 4 bits each
- The estimate of the number of distinct words is  $\widetilde{F}_0 = 30897$ , while the true answer is  $F_0 = 28239$ , which represents a relative error +9.4%.

```
ghfffghfghggghggggghghheehfhfhhgghghghhfgffffhhhiigfhhffgfiihfhhh  
igigighfgihffffghigihghigfhhgeegeghgghhhgghhhfhidiigihighihehhhfgg  
hfgighigffghdieghhhggghhfhghhfiieffghghihifgggffihgihfggighgiiif  
fjgfgjhhjiiifhjgehgghfhfhjhiggghghihiggghihihgiighgfhlgjfgjjjml
```

# Outline

---

Introduction

Approximate Counting

Distinct Elements and Frequency Moments

Extra Material (non-examinable): An Algorithm for  $F_0$  in the Turnstile Model

## The AMS algorithm (cash register model)

---

Common approach for designing algorithms in the cash register model:

1. Sample the data items based on hashed values;

## The AMS algorithm (cash register model)

---

Common approach for designing algorithms in the cash register model:

1. Sample the data items based on hashed values;
2. Store the statistical information of the sampled items, or store the sampled items directly.

## The AMS algorithm (cash register model)

---

Common approach for designing algorithms in the cash register model:

1. Sample the data items based on hashed values;
2. Store the statistical information of the sampled items, or store the sampled items directly.

Downside of this framework:

- Sampling probability for the current item usually depends on the whole data stream that algorithm has seen so far.
- Deleting an item appeared before could potentially make the current statistical information useless! :(

## The AMS algorithm (cash register model)

---

Common approach for designing algorithms in the cash register model:

1. Sample the data items based on hashed values;
2. Store the statistical information of the sampled items, or store the sampled items directly.

Downside of this framework:

- Sampling probability for the current item usually depends on the whole data stream that algorithm has seen so far.
- Deleting an item appeared before could potentially make the current statistical information useless! :(

Sampling techniques are usually non-applicable in the turnstile model.

## Algorithm to approximate $F_2$ in the turnstile model

---

Algorithm to approximate  $F_2$  (simplified description)

- 1: Choose a **4-wise independent** hash function  $h : [n] \rightarrow \{-1, 1\}$
- 2:  $y = 0$
- 3: **while** item  $(x, \pm)$  from stream  $S$  arrives
- 4:     **if**  $x$  is inserted **then**  $y \leftarrow y + h(x)$
- 5:                     **else**  $y \leftarrow y - h(x)$
- 6: **return**  $Z := y^2$



## Algorithm to approximate $F_2$ in the turnstile model

Algorithm to approximate  $F_2$  (simplified description)

- 1: Choose a **4-wise independent** hash function  $h : [n] \rightarrow \{-1, 1\}$
- 2:  $y = 0$
- 3: **while** item  $(x, \pm)$  from stream  $S$  arrives
- 4:     **if**  $x$  is inserted **then**  $y \leftarrow y + h(x)$
- 5:                     **else**  $y \leftarrow y - h(x)$
- 6: **return**  $Z := y^2$

The algorithm runs in the turnstile model!

## Algorithm to approximate $F_2$ in the turnstile model

Algorithm to approximate  $F_2$  (simplified description)

- 1: Choose a **4-wise independent** hash function  $h : [n] \rightarrow \{-1, 1\}$
- 2:  $y = 0$
- 3: **while** item  $(x, \pm)$  from stream  $S$  arrives
- 4:     **if**  $x$  is inserted **then**  $y \leftarrow y + h(x)$
- 5:                             **else**  $y \leftarrow y - h(x)$
- 6: **return**  $Z := y^2$

The algorithm runs in the turnstile model!

**Key Lemma**

It holds that  $\mathbf{E}[Z] = F_2$  and  $\mathbf{V}[Z] \leq 2 \cdot (\sum_{i \in S} m_i^2)^2 = 2F_2^2$ .

## Algorithm to approximate $F_2$ in the turnstile model

Algorithm to approximate  $F_2$  (simplified description)

- 1: Choose a **4-wise independent** hash function  $h : [n] \rightarrow \{-1, 1\}$
- 2:  $y = 0$
- 3: **while** item  $(x, \pm)$  from stream  $S$  arrives
- 4:     **if**  $x$  is inserted **then**  $y \leftarrow y + h(x)$
- 5:                     **else**  $y \leftarrow y - h(x)$
- 6: **return**  $Z := y^2$

The algorithm runs in the turnstile model!

**Key Lemma**

It holds that  $\mathbf{E}[Z] = F_2$  and  $\mathbf{V}[Z] \leq 2 \cdot (\sum_{i \in S} m_i^2)^2 = 2F_2^2$ .

Hence, we can  $(\epsilon, \delta)$ -approximate  $F_2$ , by **running multiple copies of the algorithm in parallel and return the average value.**

## Algorithm to approximate $F_2$ in the turnstile model

Algorithm to approximate  $F_2$  (details)

- 1:  $t = \lceil 6/\varepsilon^2 \rceil$
- 2: Choose  $t$  **4-wise independent hash function**  $h_1, \dots, h_t$ , where

$$h_i : [n] \rightarrow \{-1, 1\}$$

- 3:  $y_i = 0$  for each  $i = 1, \dots, t$
- 4: **while** item  $(x, \pm)$  from stream  $\mathcal{S}$  arrives
- 5:     **if**  $x$  is inserted **then**  $y_i = y_i + h_i(x)$  for every  $1 \leq i \leq t$
- 6:     **else**  $y_i = y_i - h_i(x)$  for every  $1 \leq i \leq t$
- 7: **return**  $\frac{1}{t} \cdot \sum_{i=1}^t Z_i$ , where  $Z_i = y_i^2$

## Algorithm to approximate $F_2$ in the turnstile model

Algorithm to approximate  $F_2$  (details)

- 1:  $t = \lceil 6/\varepsilon^2 \rceil$
- 2: Choose  $t$  **4-wise independent hash function**  $h_1, \dots, h_t$ , where

$$h_i : [n] \rightarrow \{-1, 1\}$$

- 3:  $y_i = 0$  for each  $i = 1, \dots, t$
- 4: **while** item  $(x, \pm)$  from stream  $\mathcal{S}$  arrives
- 5:     **if**  $x$  is inserted **then**  $y_i = y_i + h_i(x)$  for every  $1 \leq i \leq t$
- 6:     **else**  $y_i = y_i - h_i(x)$  for every  $1 \leq i \leq t$
- 7: **return**  $\frac{1}{t} \cdot \sum_{i=1}^t Z_i$ , where  $Z_i = y_i^2$

### Analysis

With constant probability, the returned value of the algorithm lies in  $(1 - \varepsilon, 1 + \varepsilon) \cdot F_2$ . Moreover, the space complexity is  $O((1/\varepsilon^2) \log n)$  bits.