

GEMS OF TCS

STREAMING ALGORITHMS

Sasha Golovnev

September 6, 2023

FRUIT GAME

Credit: Jelani Nelson

(<https://www.youtube.com/watch?v=CorP4I23wOo&t=2434s>)

STREAMING ALGORITHMS

- Massively long stream of data

STREAMING ALGORITHMS

- Massively long stream of data
Instagram, search queries, network packets

STREAMING ALGORITHMS

- Massively long stream of data
Instagram, search queries, network packets
 $X_1, X_2, X_3, \dots, X_n$

STREAMING ALGORITHMS

- Massively long stream of data
Instagram, search queries, network packets
 $X_1, X_2, X_3, \dots, X_n$
- Data has grown: we can't afford even storing it

STREAMING ALGORITHMS

- Massively long stream of data
Instagram, search queries, network packets
 $X_1, X_2, X_3, \dots, X_n$
- Data has grown: we can't afford even storing it
- n inputs, space \sqrt{n} ; $\log^{10} n$; $\log n$

STREAMING ALGORITHMS

- Massively long stream of data
Instagram, search queries, network packets
 $X_1, X_2, X_3, \dots, X_n$
- Data has grown: we can't afford even storing it
- n inputs, space \sqrt{n} ; $\log^{10} n$; $\log n$
- Efficient processing of stream

STREAMING ALGORITHMS

- Massively long stream of data
Instagram, search queries, network packets
 $X_1, X_2, X_3, \dots, X_n$
- Data has grown: we can't afford even storing it
- n inputs, space \sqrt{n} ; $\log^{10} n$; $\log n$
- Efficient processing of stream
- Mostly randomized algorithms

Missing Number

MISSING NUMBER

- Stream contains n **distinct** numbers in range $\{0, \dots, n\}$

MISSING NUMBER

- Stream contains n **distinct** numbers in range $\{0, \dots, n\}$
- Return the only **missing** number

MISSING NUMBER

- Stream contains n **distinct** numbers in range $\{0, \dots, n\}$
- Return the only **missing** number
- Efficient algorithm?

STREAMING ALGORITHM

- Compute sum of **all** elements in stream:

$$S = x_1 + \dots x_n$$

STREAMING ALGORITHM

- Compute sum of **all** elements in stream:

$$S = x_1 + \dots + x_n$$

- Sum of all numbers in range $\{0, \dots, n\}$ is

$$S = \frac{n(n+1)}{2}$$

STREAMING ALGORITHM

- Compute sum of **all** elements in stream:

$$S = x_1 + \dots + x_n$$

- Sum of all numbers in range $\{0, \dots, n\}$ is

$$S = \frac{n(n+1)}{2}$$

- Missing number is $S - s = \frac{n(n+1)}{2} - s$

STREAMING ALGORITHM

- Compute sum of **all** elements in stream:

$$S = x_1 + \dots + x_n$$

- Sum of all numbers in range $\{0, \dots, n\}$ is

$$S = \frac{n(n+1)}{2}$$

- Missing number is $S - s = \frac{n(n+1)}{2} - s$
- One pass through stream, efficient processing, $O(\log n)$ space

TWO MISSING ELEMENTS

- Stream contains $n - 1$ **distinct** numbers in range $\{0, \dots, n\}$

TWO MISSING ELEMENTS

- Stream contains $n - 1$ **distinct** numbers in range $\{0, \dots, n\}$
- Return **both missing** numbers

TWO MISSING ELEMENTS

- Stream contains $n - 1$ **distinct** numbers in range $\{0, \dots, n\}$
- Return **both missing** numbers
- Efficient algorithm?

STREAMING ALGORITHM

- Compute **sum and sum of squares** of all elements in stream:

$$S = x_1 + \dots x_{n-1}$$

$$t = x_1^2 + \dots x_{n-1}^2$$

STREAMING ALGORITHM

- Compute **sum and sum of squares** of all elements in stream:

$$S = x_1 + \dots + x_{n-1}$$

$$t = x_1^2 + \dots + x_{n-1}^2$$

- Sum of all numbers in range $\{0, \dots, n\}$ is

$$S = \frac{n(n+1)}{2}$$

Sum of squares of all numbers in range

$\{0, \dots, n\}$ is $T = \frac{n(n+1)(2n+1)}{6}$

STREAMING ALGORITHM

- If missing numbers are a and b , then

$$a + b = S - s$$

$$a^2 + b^2 = T - t$$

STREAMING ALGORITHM

- If missing numbers are a and b , then

$$a + b = S - s$$

$$a^2 + b^2 = T - t$$

- One pass through stream, efficient processing,
 $O(\log n)$ space

Majority Element

MAJORITY ELEMENT

- Stream has element occurring $> n/2$ times

MAJORITY ELEMENT

- Stream has element occurring $> n/2$ times
- Find it!

STREAMING ALGORITHM

- $\text{count} \leftarrow 0; m \leftarrow \perp$

STREAMING ALGORITHM

- $\text{count} \leftarrow 0; m \leftarrow \perp$
- For each element x_i of Stream:

STREAMING ALGORITHM

- $\text{count} \leftarrow 0; m \leftarrow \perp$
- For each element x_i of Stream:
 - If $\text{count} = 0$, then $m \leftarrow x_i$ and $\text{count} \leftarrow 1$

STREAMING ALGORITHM

- $\text{count} \leftarrow 0; m \leftarrow \perp$
- For each element x_i of Stream:
 - If $\text{count} = 0$, then $m \leftarrow x_i$ and $\text{count} \leftarrow 1$
 - Else if $x_i = m$, then $\text{count} ++$

STREAMING ALGORITHM

- $\text{count} \leftarrow 0; m \leftarrow \perp$
- For each element x_i of Stream:
 - If $\text{count} = 0$, then $m \leftarrow x_i$ and $\text{count} \leftarrow 1$
 - Else if $x_i = m$, then $\text{count} ++$
 - Else $\text{count} --$

STREAMING ALGORITHM

- $\text{count} \leftarrow 0$; $m \leftarrow \perp$
- For each element x_i of Stream:
 - If $\text{count} = 0$, then $m \leftarrow x_i$ and $\text{count} \leftarrow 1$
 - Else if $x_i = m$, then $\text{count} ++$
 - Else $\text{count} --$
- Return m

EXAMPLE

PROOF

ANOTHER VIEW

MISRA-GRIES ALGORITHM

MISRA-GRIES ALGORITHM

- $\text{count}_1, \dots, \text{count}_k \leftarrow 0$; $m_1, \dots, m_k \leftarrow \perp$
- For each element x_i of Stream:
 - If $x_i = m_j$, then $\text{count}_j ++$
 - Else
 - Let count_j be min in $\text{count}_1, \dots, \text{count}_k$
 - If $\text{count}_j = 0$, then $m_j = x_i$; $\text{count}_j = 1$
 - Else $\text{count}_1 --, \dots, \text{count}_k --$
- Return m_1, \dots, m_k

Approximate Counting

- Router receives stream of network packages

- Router receives stream of network packages
- Want to count number of packages from IP "1.2.3.4"

- Router receives stream of network packages
- Want to count number of packages from IP "1.2.3.4"
- Efficient algorithm?

- Router receives stream of network packages
- Want to count number of packages from IP "1.2.3.4"
- Efficient algorithm?
- Efficient **approximate** algorithm?

OVERVIEW

MORRIS ALGORITHM

MORRIS ALGORITHM

- $c \leftarrow 0$

MORRIS ALGORITHM

- $c \leftarrow 0$
- When see next element:
 - with probability $\frac{1}{2^c}$ increment c
 - with probability $1 - \frac{1}{2^c}$ do nothing

MORRIS ALGORITHM

- $c \leftarrow 0$
- When see next element:
 - with probability $\frac{1}{2^c}$ increment c
 - with probability $1 - \frac{1}{2^c}$ do nothing
- Return $2^c - 1$

ANALYSIS

PROBABILITY OF SUCCESS

- $\mathbb{E}[\text{output}] = n$

PROBABILITY OF SUCCESS

- $\mathbb{E}[\text{output}] = n$
- By Markov's, $\Pr[\text{output} \geq 2n] \leq 1/2$

PROBABILITY OF SUCCESS

- $\mathbb{E}[\text{output}] = n$
- By Markov's, $\Pr[\text{output} \geq 2n] \leq 1/2$
- Similar inequalities show that $\Pr[\text{output} \in [n - O(n), n + O(n)]] \geq 0.9$

PROBABILITY OF SUCCESS

- $\mathbb{E}[\text{output}] = n$
- By Markov's, $\Pr[\text{output} \geq 2n] \leq 1/2$
- Similar inequalities show that $\Pr[\text{output} \in [n - O(n), n + O(n)]] \geq 0.9$
- Again, repeating Algorithm several times significantly amplifies probability of success

SUMMARY

- One pass through stream may be sufficient

SUMMARY

- One pass through stream may be sufficient
- Use Randomness and Approximation

SUMMARY

- One pass through stream may be sufficient
- Use Randomness and Approximation
- Markov's inequality: from Expectation to Probability

SUMMARY

- One pass through stream may be sufficient
- Use Randomness and Approximation
- Markov's inequality: from Expectation to Probability
- Amplify probability by Repetitions