# Gems of TCS

## Streaming Algorithms

Sasha Golovnev

February 4, 2021

# 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
  $x_1, x_2, x_3, \ldots, x_n$

# STREAMING ALGORITHMS

- Massively long stream of data
  Instagram, search queries, network packets
  $x_1, x_2, x_3, \ldots, x_n$

- Data has grown: we can't afford even stoing it

# STREAMING ALGORITHMS

- Massively long stream of data
  Instagram, search queries, network packets
  $x_1, x_2, x_3, \ldots, x_n$

- Data has grown: we can't afford even stoing it

- $n$ inputs, space $\sqrt{n}$;  $\log^{10} n$;  $\log n$

# Streaming Algorithms

- Massively long stream of data
  Instagram, search queries, network packets
  $x_1, x_2, x_3, \ldots, x_n$

- Data has grown: we can't afford even stoing it

- $n$ inputs, space $\sqrt{n}$; $\quad \log^{10} n$; $\quad \log n$

- Efficient processing of stream

# STREAMING ALGORITHMS

- Massively long stream of data
  Instagram, search queries, network packets
  $x_1, x_2, x_3, \ldots, x_n$

- Data has grown: we can't afford even stoing it

- $n$ inputs, space $\sqrt{n}$;   $\log^{10} n$;   $\log n$

- Efficient processing of stream

- Mostly randomized algorithms

# Missing Number

- Stream contains $n$ distinct numbers in range $\{0, \ldots, n\}$     all but one

# MISSING NUMBER

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

- Return the only missing number

Could sort — linear space, go through stream many times

# MISSING NUMBER

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

- Return the only missing number

- Efficient algorithm?

- Compute sum of all elements in stream:

$$s = x_1 + \ldots x_n$$

# STREAMING ALGORITHM

- Compute sum of all elements in stream:

$$S = x_1 + \ldots x_n$$

- Sum of all numbers in range $\{0, \ldots, n\}$ is $S = \frac{n(n+1)}{2}$

# STREAMING ALGORITHM

- Compute sum of all elements in stream:

$$\overset{\text{0 21 7456}}{S = x_1 + \ldots x_n}$$

- Sum of all numbers in range $\{0, \ldots, n\}$ is
$S = \frac{n(n+1)}{2}$   $0+1+2'-- + 6 - (0+2+1+7+4+5+6)$

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

See element $\rightarrow$ process quickly - one addition

$O(\log n)$

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

- 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 $\{0, \ldots, n\}$

# Two Missing Elements

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

- Return both missing numbers

# Two Missing Elements

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

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

$$S = 0 + 1 + \dots + n$$

- Return both missing numbers

$$S - s = a + b$$

- Efficient algorithm?

Don't want to sort

# Streaming Algorithm

- Compute <span style="color:orange">sum and sum of squares</span> of all elements in stream:

$$s = x_1 + \ldots x_{n-1} \quad \text{← sum of inputs}$$

$$t = x_1^2 + \ldots x_{n-1}^2$$
$$\text{↖ sum of their squares}$$

# Streaming Algorithm

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

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

- Sum of all numbers in range $\{0, \ldots, n\}$ is
$S = \frac{n(n+1)}{2}$
Sum of squares of all numbers in range $\{0, \ldots, n\}$
is $T = \boxed{\frac{n(n+1)(2n+1)}{6}}$  $\quad T = \sum_{i=0}^{n} i^2$

# Streaming Algorithm

- If missing numbers are $a$ and $b$, then

I know

we know

we know

$$u = \underline{a + b} = \underline{S - s}$$

$$v = \underline{a^2 + b^2} = \boxed{T - t}$$

Want to find $a, b$

$$w = \frac{u^2 - v}{2} = \frac{(a+b)^2 - a^2 - b^2}{2} = \frac{2ab}{2} = ab$$

$u = a + b$    I know $u, w$

$w = a \cdot b$    Want to find $a, b$

$a = u - b$

$w = (u - b) b$

$b^2 - \underline{ub} + \underline{w} = 0$

$\sim$ I know

Find $b$

$\textcircled{D} = u^2 - 4w$

$b = \dfrac{u \pm \sqrt{u^2 - 4w}}{2}$

Two solutions are the two

missing els

# STREAMING ALGORITHM

- If missing numbers are *a* and *b*, then

$$S = x_1 + \cdots + x_n$$
$$T = x_1^2 + \cdots + x_n^2$$

$$a + b = S - s$$

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

This can be generalized

k els are missing

$$S_1 = x_1 + \cdots + x_n$$
$$S_2 = x_1^2 + \cdots + x_n^2$$
$$\overline{\phantom{xx}} \ \ \overline{\phantom{xx}}$$
$$S_k = x_1^k + \cdots + x_n^k$$

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

# Majority Element

# Majority Element

n – length of stream

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

# Majority Element

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

- Find it!

  Sort , Median

  We can't affand stowing input

candidate for Maj

- count $\leftarrow$ 0;  m $\leftarrow \perp$ **Null**

# STREAMING ALGORITHM

- count $\leftarrow 0$;   m $\leftarrow \perp$

- For each element $x_i$ of Stream:

- count $\leftarrow 0$; m $\leftarrow \perp$

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

new candidate

# STREAMING ALGORITHM

- count $\leftarrow 0$;   m $\leftarrow \perp$

- For each element $x_i$ of Stream:
    - If **count** $= 0$, then **m** $\leftarrow x_i$ and **count $\leftarrow 1$**
    - ElseIf $x_i = $ **m**, then **count $++$**

# STREAMING ALGORITHM

- count $\leftarrow 0$;  m $\leftarrow \perp$

- For each element $x_i$ of Stream:
  - If count $= 0$, then m $\leftarrow x_i$ and count $\leftarrow 1$
  - ElseIf $x_i =$ m, then count $++$
  - Else count $--$   $x_i \neq m$

# STREAMING ALGORITHM

two variables

1 pass

processes input els: efficiently

- count $\leftarrow 0$;  m $\leftarrow \perp$

- For each element $x_i$ of Stream:
  - If count $= 0$, then m $\leftarrow x_i$ and **count $\leftarrow$ 1**
  - ElseIf $x_i = $ m, then count $++$
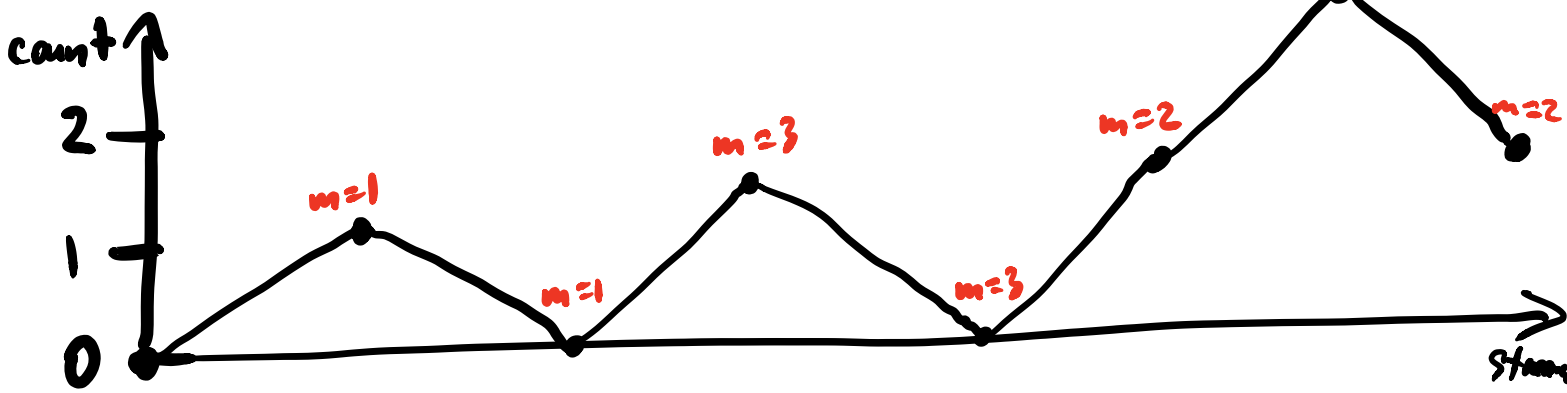  - Else count $--$   $x_i \neq m$
- Return m

$n = 7$     Max $= 2$

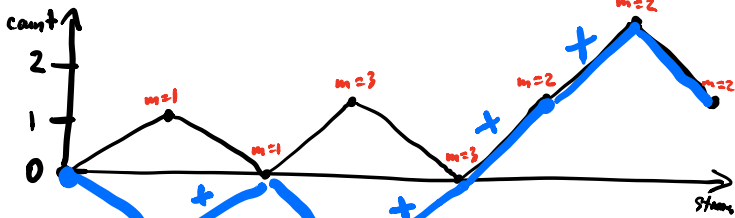| $x_i$ | 1 | 2 | 3 | 2 | 2 | 2 | 1 |
|---|---|---|---|---|---|---|---|
| $m$ | 1 | 1 | 3 | 3 | 2 | 2 | 2 |
| count | 0 | 1 | 0 | 1 | 0 | 1 | 2 | 1 |

Return 2



count ≥ 0

Maj = 2

$x_i$   1   2   3   2   2   2   1



variable count'
this is only for analysis
not for algorithm

$$count' = \begin{cases} \boxed{\text{Count}} & \text{if } m = Maj \quad \text{— time Maj} \in \mathbb{t} \\ -count & \text{if } m \neq Maj \end{cases}$$

— When I see **Maj**, increment **count'**

**Proof:** if $m = Maj \Rightarrow count++ \Rightarrow count'++$
if $m \neq Maj \Rightarrow count-- \Rightarrow count'++$ □

— See $Maj > \frac{n}{2}$ times $\Rightarrow$ count' is incremented $> \frac{n}{2}$
$\Rightarrow$ count' is decremented $< \frac{n}{2}$

In the end,
— $\sqrt{}$ count' $> 0 \Rightarrow$ count' $=$ count $\Rightarrow m = Maj \Rightarrow$
output Maj

$\cancel{1} \quad \cancel{2} \quad \cancel{3} \quad \cancel{2} \quad \cancel{2} \quad 2 \quad \cancel{1}$

- Pairs up distinct els
- Kills all these pairs
- Remaining els are Maj

$[ \times \times \times \cdots \boxed{m} | \times \times \times \times \cdots a_{n+1} ]$

$n/2+1 \qquad \frac{n}{2}-1$

Majority els remain

Assume there is Maj in stream
$$\left( > \frac{n}{2} \text{ occurrences} \right),$$

## Find it

without this assumption, we'll make two passes through input stn.

— candidate $\underline{\underline{m}}$

— count how many times $m$ appears in stream

---

For $k \in \mathbb{N}$, $k$-Heavy Hitters : Find all els that appear $\geqslant \frac{n}{k}$ in the stream

($\leq k$ such els) $\qquad$ Maj = case $k=2$

# MISRA-GRIES ALGORITHM

- $count_1, \ldots, count_k \leftarrow 0; \quad m_1, \ldots, m_k \leftarrow \perp$

- For each element $x_i$ of Stream:
    - If $x_i = m_j$ then $count_j$ ++
    - Else
        - Let $count_j$ be min in $count_1, \ldots count_k$
        - If $count_j = 0$, then $m_j = x_i; \quad count_j = 1$
        - Else $count_1$ --, $\ldots$, $count_k$ --
- Return $m_1, \ldots, m_k$ ⟵ contain all els of stream that appear $\geq n/k$ times

# 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"

~~1.1.1.1.~~   (1. 2. 3.4.)   ~~2.2.2.2.~~   (1.2.3.4)

EQ:   <u>1. 2. 3. 4</u>   <u>1. 2. 3. 4.</u>
      $\leq n$ of them in the stream
output length of the stream

- Router receives stream of network packages

- Want to count number of packages from IP "1.2.3.4"

- Efficient algorithm?

$count = 0$

See input: $count++$

In the end, $count = $ stream length

$\log_2(n+1)$ bits

Can we use fewer than $\log n$ bits?

Input length $\leq n$
outputs $\in \{0, 1, 2, \dots, n-1, n\}$

2 bits

| 0 | 1 |
|---|---|

00
01
10
11

$\leq 4$ distinct answers

IF $< \log_2 n$ bits $\Rightarrow$ different answers
$< 2^{\log_2 n} = n$

$\lceil \log_2 (n+1) \rceil$ bits is optimal

- Router receives stream of network packages

- Want to count number of packages from
  IP "1.2.3.4"

- Efficient algorithm?

- Efficient approximate algorithm?

$\log \log n \leftarrow$ exponentially better than previous sol

# OVERVIEW

n

Trivial alg   stores n

n=147 [1 0 0 1 0 0 1 1] ← log n

what if instead of storing **n in binary**

I'm storing the length of "n in binary"

Instead of storing 147, I'd store numbers

[1 0 0 0]

— Instead of log n bits to write $\{1, \dots, n\}$
store log log n bits to write $\{1, \dots, \log n\}$

If length = 4 Then

8 [1 0 0 0]

15 [1 1 1 1]

$8 \leq n \leq 15$

$2^{length-1} \leq n < 2^{length}$

$c$ — $\approx$ length of $n$ in binary

$n \approx 2^c$

when should increment $c$?

I want to increment $c$ after seeing $\boxed{2^c}$ new els

Now see new el

$\quad$ w.p. $\frac{1}{2^c}$ $\quad$ $c$++

$\quad$ w.p. $(1 - \frac{1}{2^c})$ don't update $c$

After seeing $2^c$ els, I expect to increment $c$ once

# MORRIS ALGORITHM

$$n = 0$$

$$2^c - 1 \approx n$$

- $c \leftarrow 0$

# 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

$$n \approx 2^c - 1$$

- $c \leftarrow 0$

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

- Thm

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

- Manhors:

$$E[\text{output} > 5n] < \frac{1}{5}$$

1. Manhor's ineq:   $\Pr[\text{output} \notin [n-O(n), n+O(n)] <$

$$< 0.9$$

2. Amplify prob. of success by repeating this alg several times:

$$\Pr[\text{output} \notin [\tfrac{n}{2}, 2n]] < 0.01$$

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

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

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

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

- Similar inequalities show that
  $\text{Pr}[\text{output} \in [n - O(n), n + O(n)] \geq 0.9$

# 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 ineqaulity: from Expectation to Probability

# SUMMARY

- One pass through stream may be sufficient

- Use Randomness and Approximation

- Markov's ineqaulity: from Expectation to Probability

  *Max-Cut*

- Amplify probability by Repetitions