

## Streaming and Sketching

What kind of useful computation can be done in small space

Model :

sequence of 'tokens'

$a_1, a_2, \dots, a_n$  (stream)

Each  $a_i \in \{0, 1\}^b$   $b$  bit string

Trivial space requirement :  $n \cdot b$

Space allowed for 'algorithm' :  $S \sim o(n)$ , linear dependence on  $b$   
e.g.,  $\sqrt{n}$ ,  $2^{\sqrt{\log n}}$ ,  $\text{poly}(\log n)$ ,  $\log n$

- Objectives:
- ① Most frequent element ✓
  - ② # distinct elements
  - ③ Frequency moments

→ Goal : Find frequent elements in data stream

### Micro-Gries Algorithm

- Fix parameter  $k$ . Output all elements that appear at least  $\frac{n}{k+1}$  times (possibly with

elements that are less frequent)

Algorithm:

(1)  $K$  tokens :  $b_1, b_2, \dots, b_K$ ,  
all initialized to  $\perp$ .  
 $K$  numbers :  $c_1, c_2, \dots, c_K$ ,  
all initialized to  $\perp$ .

(2) When  $a_i$  appears :

- if one of the  $b_j$ 's is equal to  $a_i$  : increase  $c_j$  by 1
- if some  $b_j$  is  $\perp$ , set  $b_j = a_i$  and  $c_j = 1$ .
- none of  $b_j$ 's are  $a_i$  (or  $\perp$ ) : decrease  $c_j$  by 1 for all  $j=1, \dots, K$ .

(if  $c_j = 0$ , then set  $b_j = \perp$ )

(3) Output :  $b_1, \dots, b_K$ .

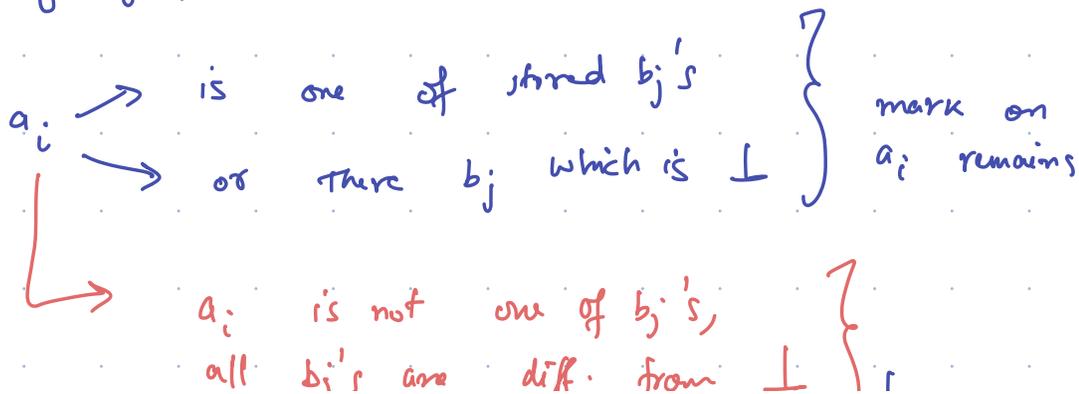
- Space requirement :  $S = \underbrace{K \cdot b}_{\text{tokens}} + \underbrace{K \cdot \lceil \log n \rceil}_{\text{counts}}$

Example  
 $K = 2$



Proof of correctness:

On seeing  $a_i$ , mark it as 'red'.



- j -                      0                      1                      2
- (1) remove mark from  $a_i$   
 (2) remove " " each  $b_j$  : from the first marked copy of  $b_j$ .

Invariant: at any point, there are exactly  $c_j$  marked copies of  $b_j$  ( $\forall j=1, \dots, k$ )

→ Proof is by induction on  $i$ .  
 $[i=0, \dots, n]$

$i=0$  ✓

Follows at step  $i$  directly from algorithm.

Now, note that at any 'removal step', exactly  $k+1$  marks are removed.

⇒ at most  $\frac{n}{k+1}$  'removal steps'.

$\Rightarrow$  if a token appears  $\geq \frac{n}{k+1}$  times,  
then some marked copy remains.

□

## Data Sketching

streaming + data structures

tokens:  $a_1, a_2, \dots, a_n$ , each  $a_i \in \{0, 1\}^b$

space:  $S$

↳ a data structure to answer queries.

→ Query: frequency of  $x$ ?  $f_x$

Trivial bounds:  $\rightarrow S = b \cdot n \rightarrow$  store stream

$\rightarrow S = 2^b \cdot \log n \rightarrow$  store frequency vector.

$\vec{f}$ :  $2^b$  coordinates indexed by  $x \in \{0, 1\}^b$ .

$x^{\text{th}}$  coordinate:  $f_x$ .

$S \sim \text{poly}(\log n, b)$ .

Algorithm: that w.p  $1 - \epsilon$  outputs  $\tilde{f}_x$   
 s.t  $|\tilde{f}_x - f_x| \leq \epsilon \cdot n$ .

Idea: 'store hash of frequency vectors'

→ Sample uniformly a 2-universal hash function  $h: \{0, 1\}^b \rightarrow [B]$ . ( $B$  to be chosen later)

→ Let  $C$  be an array of length  $B$ , all entries initialized to 0.

→ for  $i = 1$  to  $n$ :  
     on seeing  $a_i$ :  
     let  $k = h(a_i)$ .  
      $C[k] \leftarrow C[k] + 1$ ;  
 end for.

→ On query  $x$ , output  $C[h(x)] \equiv \tilde{f}_x$

Observation  $\tilde{f}_x \geq f_x$  [could be more due to collisions]

Space:  $B \cdot \log n$  + space for hash  $f_x$  ←  $O(b \cdot \log B)$