

- Given a sequence of strings S_1, \dots, S_n (maybe duplicate)
 \Rightarrow filter to keep just the distinct one (1 copy per string)
 when memory size M ; relation between n, M , total length of the strings

1) Bloom Filter: in this case $n = \Theta(M)$

$$m = c \cdot n < M$$

\rightarrow size of Bloom F

$$k = \frac{m}{n} \ln 2 = \frac{c \cdot n}{n} \ln 2 = c \ln 2$$

\Rightarrow - SCAN THE STRINGS S_i and \forall of them set

// take a string and check,
 if all pos are 1 strings
 \Rightarrow belongs to set
 \Rightarrow we don't do anything

check if $BI[h_j(s_i)] = 1 \quad \forall j=1, k$
 $\forall j = 1, \dots, k$
 $\forall i = 1, \dots, n$

else print S_i
 and

$\forall j = 1, \dots, k$ set

$BI[h_j(s_i)] = 1$

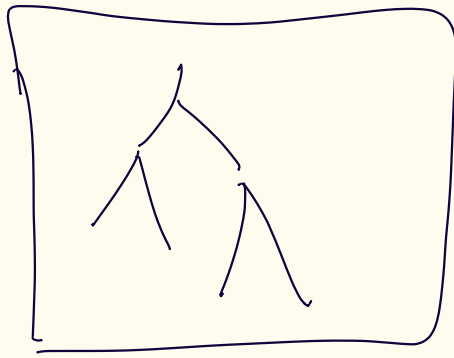
// for groups

Cost \rightarrow total length of the strings

$O\left(\frac{N}{B}\right)$ Total \Rightarrow It's a scan and all checks are done in M

time $O(n \cdot k + N)$
 \leftarrow insertion \leftarrow scan

TIME is the total length of the distinct strings
 A's in $M \Rightarrow < M$



~~~~~>>

scan on diff and check if  $i$  is in the  
 A's

BETTER  $\Rightarrow$  PATRICIA TREE: depends on the # of strings  
 & their total length

# string  $S_i$

$\Rightarrow$  do a Patricia <sup>tree</sup> search to know if  
 $S_i \in PT$  if present, do not

$\Rightarrow$  if is not present

$\Rightarrow$  print

insert  $S_i$  in  $PT$

read string

GO DOWN

LCP ( $S_i$ , reference)



BUT

$(n < M)$

$\Rightarrow$  Patricia tree  
 uses A's  
 in  $M$

BEST SOLUTION  $\Rightarrow$  hash

$\forall s_i \rightarrow h(s_i) \xrightarrow{\text{create pair}} \langle h(s_i), i \rangle$

and use an externally memory sorter

$\Rightarrow$  SORT

$\Rightarrow$  keep  $i$  of the first after sort

SPECIAL BP

7 cells  $m = 7$

insert (1 2 3 4)

| Key | $h_1$ | $h_2$ |
|-----|-------|-------|
| 1   | 2     | 3     |
| 2   | 4     | 6     |
| 3   | 6     | 2     |
| 4   | 1     | 5     |

$$h_1: x \bmod 7$$

$$h_2: 3x \bmod 7$$

| 0 | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 1 | 1 | 2 |

metto un 1 qui sotto che entra in una cella

Query (5)

$$h_1(5) = 10 \bmod 7 = 3$$

$$h_2(5) = 15 \bmod 7 = 1$$

$\Rightarrow$  Questo è errore  
perché 5 is not  
in set of keys

Il numero ottimale di funzioni di hash per il Bloom filter di size  $m$  and  $n$  keys

$$\Rightarrow k = \frac{m}{n} \ln 2$$

probabilità di avere una cella  $\Rightarrow$

$$e^{-\frac{km}{n}} = e^{-\frac{\frac{m}{n} \ln 2 \cdot n}{n}} = e^{-\ln 2} = \frac{1}{e^{\ln 2}} = \frac{1}{2}$$