# Filters :-
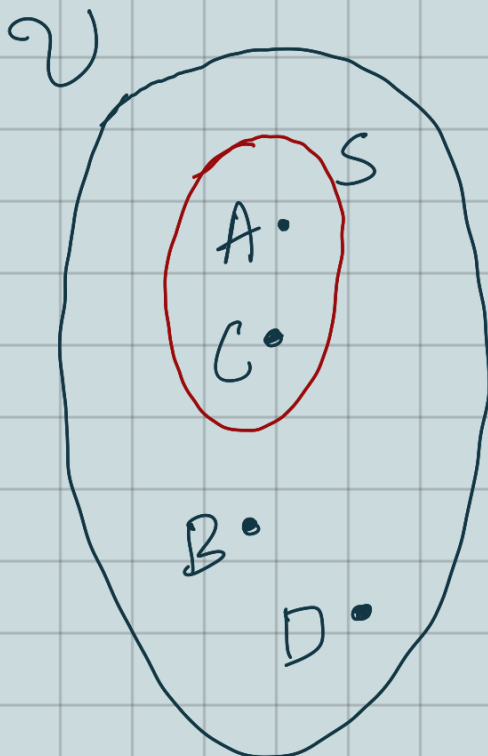
Filters represent a set approximately; trading of accuracy for space efficiency.

→ Bloom filter      } Dynamic
→ quotient filter         } deletes
→ cuckoo filter
→ XOR filter        } Static filter
→ Ribbon filter.

Static :- the set of items is known in advance
Dynamic :- the set of items is not known in advance.



A : ✓

B : X

C : ✓

D : ✓  false positive

A filter guarantees a false-positive rate $\epsilon$

if $q \in S$, return $\checkmark$    with prob. 1

if $q \notin S$, return $\begin{cases} \times \text{ with prob.} > 1 - \epsilon \\ \checkmark \text{ with prob.} \leq \epsilon \end{cases}$

Filter have one-sided errors.
(No false Negatives)

Space usage:

Filter

$\geq n \lg \dfrac{1}{\epsilon}$ bits

Dictionary
(hash table)

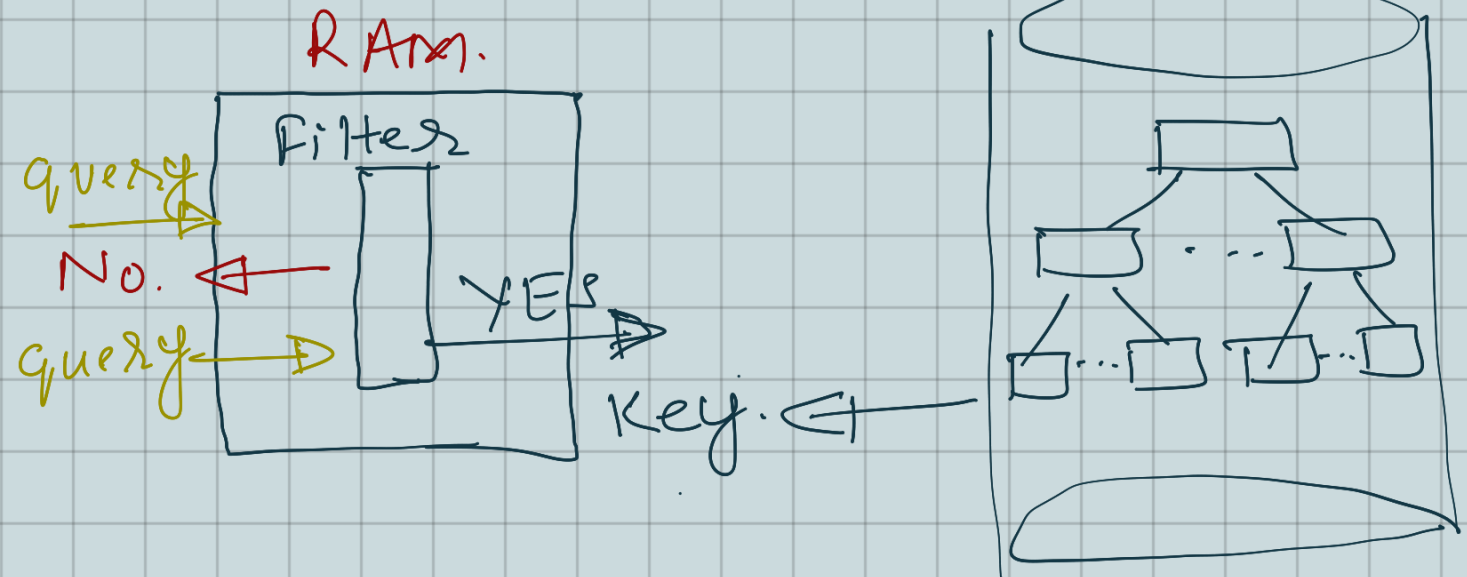$\Omega\left(n \cdot \lg |\mathcal{U}|\right)$ bits

for most practical purposes:

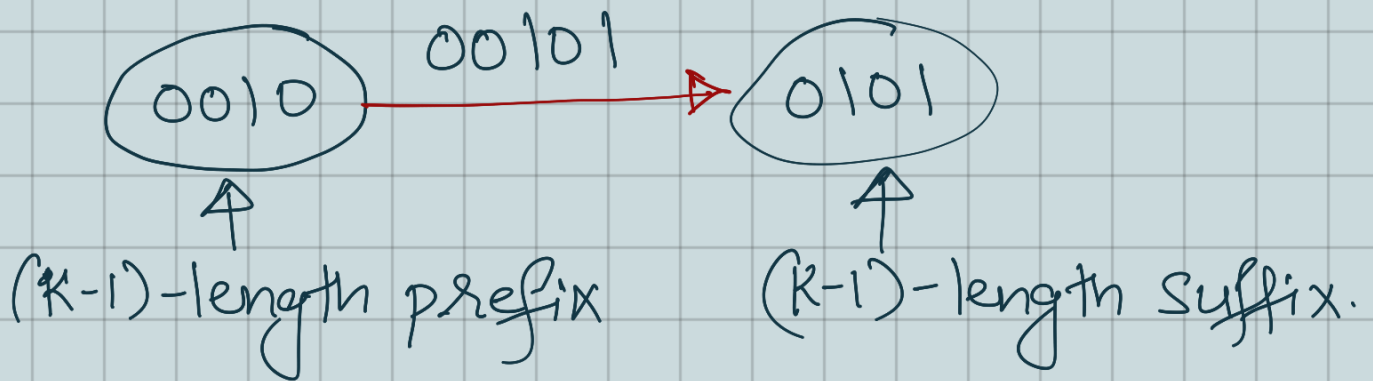$\epsilon = 2\%$, a filter requires $\sim 8$ bits/item

# Filter use case:-

## Databases:



→ In a database, disk accesses are costly and dominate the performance.

→ RAM is much smaller in size than Disk.
→ filter represent the key set approx. in RAM
→ filters can help avoid going to disk for negative queries.

→ For use in databases, filters are often required to support insert, queries, deletes.

# de Bruijn graph:

$$0010 \xrightarrow{\;00101\;} 0101$$

(K-1)-length prefix      (K-1)-length suffix.

An edge is a length-K string connecting its two (K-1) substrings.

## In genomics :-

C A C T G A A → Read.

K-mers (K=4)
$$\begin{cases} \text{C A C T} \\ \text{A C T G} \\ \text{C T G A} \\ \text{T G A A} \end{cases}$$

$$\boxed{CACT} \longrightarrow \boxed{ACTG} \longrightarrow \boxed{CTGA} \longrightarrow \boxed{TGAA}$$
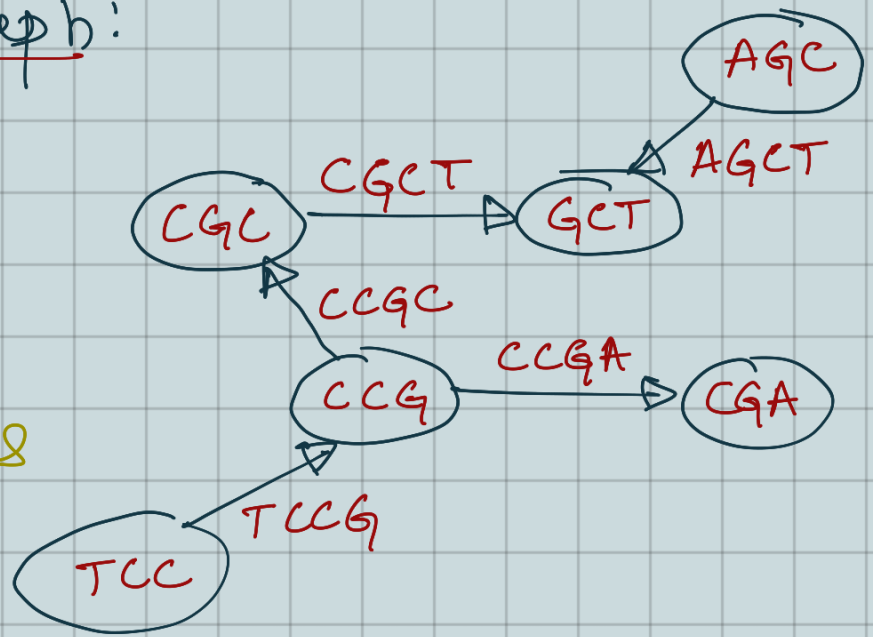
# de Bruijn graph:

## Set:

```
T  C  C  G
C  C  G  C
C  C  G  A
C  G  C  T
A  G  C  T
```

} Edges

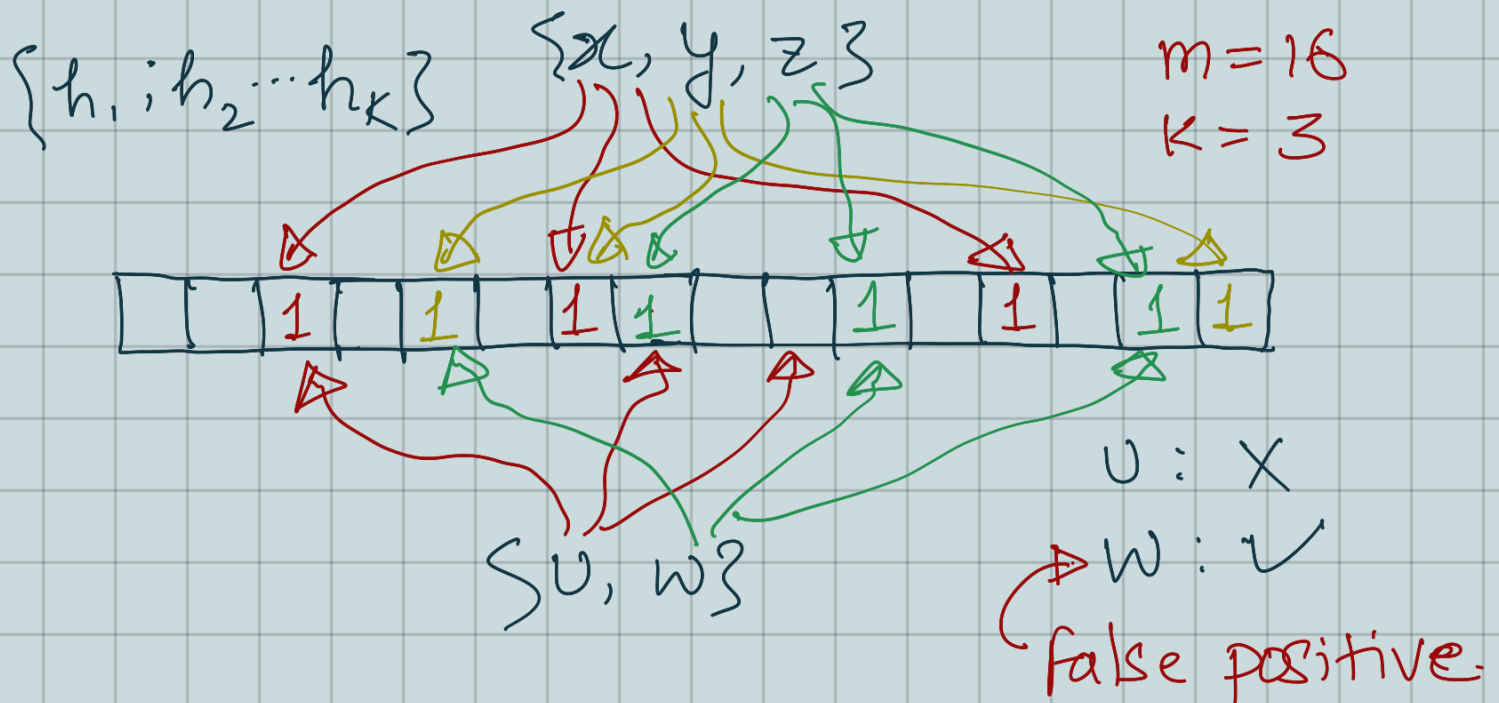

A filter can be used to represent a de Bruijn graph approximately.

→ Using a filter to traverse the deBruijn graph causes a small number of topological errors.

Question:- How can you remove the topological errors?

# Bloom filter (BF):

A BF consists of a bit vector of size $m$ and $K$ hash functions.

$\{h_1 ; h_2 \cdots h_K\}$  $\{x, y, z\}$  $m = 16$

$K = 3$

| | | 1 | | 1 | | 1 | 1 | | | 1 | | 1 | | 1 | 1 |

$\{U, W\}$

$U : X$

$W : \checkmark$

false positive.

Space : $\approx 1.44 \, n \, \lg\left(\frac{1}{\epsilon}\right)$ bits.

→ Bloom filters do not support deletes.

⇒ Q: how can you delete in Bloom filter?

→ Q: Can you merge two Bloom filters?

# False positive analysis:

Two parameter:- $m, K$

$$\Pr \left( \begin{array}{l} \text{a certain bit is not set 1} \\ \text{by a certain hash ftn.} \end{array} \right) = 1 - \frac{1}{m}$$

$$\Pr \left( \begin{array}{l} \text{a certain bit is not set 1} \\ \text{by any hash ftn.} \end{array} \right) = \left( 1 - \frac{1}{m} \right)^{K}$$

$$\Pr \left( \begin{array}{l} \text{a certain bit is not set 1} \\ \text{after } n \text{ insertions} \end{array} \right) = \left( 1 - \frac{1}{m} \right)^{Kn}$$

$$\boxed{\lim_{m \to \infty} \left( 1 - \frac{1}{m} \right)^{m} = \frac{1}{e}}$$

$$= e^{-Kn/m}$$

$$\Pr \left( \begin{array}{l} \text{a certain bit is 1} \\ \text{after } n \text{ insertions} \end{array} \right) = 1 - e^{-Kn/m}$$

$$\Pr \left( \begin{array}{l} \text{all } K \text{ bits are 1 during} \\ \text{a query of item not in} \\ \text{the set} \end{array} \right) = \left( 1 - e^{-\frac{Kn}{m}} \right)^{K}$$

$$\|$$

False positive rate

→ Using more space reduces false positive rate

→ Adding more items increases false positive rate

→ For a given value of m, n. there is a sweet spot for the number of hash functions.

$$K = \frac{m}{n} \ln 2$$

$$\varepsilon = 2^{-K}$$

Q:- we have 6 TB of database. of 512B keys ( ≈ 12.88 Billion Keys) we want a false positive rate ∼ 1.56%
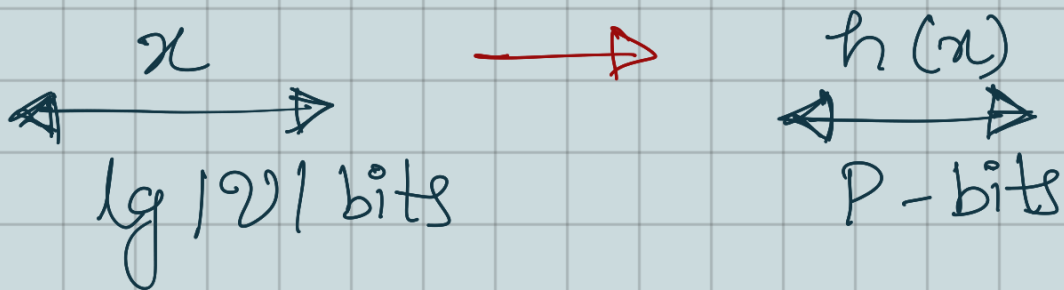
We need:-
6    37 bit hash functions

≈ 1 Byte / Key , Total: ∼ 12 GB

## Single-hashing filters:-

→ Use a hash function to hash items to a p-bit fingerprint

→ Store fingerprints compactly in a hash ~~table~~.

$$x \xleftrightarrow{\hspace{2cm}} \quad \longrightarrow \quad h(x) \xleftrightarrow{\hspace{2cm}}$$

$$\lg |\mathcal{U}| \text{ bits} \qquad\qquad P\text{-bits}$$

→ Only source of false positives:

→ Two distinct elements $x$ and $y$, where $h(x) = h(y)$

→ if $x$ is stored and $y$ isn't,

query $(y)$ gives a false positive

$$Pr(x \text{ & } y \text{ collide}) = \frac{1}{2^P}$$

# How to store fingerprints compactly.

## Quotienting :-

$h(x)$
$\leftarrow P \rightarrow$

$\longrightarrow$

$\boxed{b(x) \mid t(x)}$
$\leftarrow q \rightarrow \leftarrow r \rightarrow$

$b(x)$

$t(x)$
$t(y)$

$b(y)$

$b(x)$: location in the table
$t(x)$: $\bot$ stored in the table

Q: how handle collisions?

$\hookrightarrow$ Linear probing
$\hookrightarrow$ Robin hood hashing.

occupieds
$\hookrightarrow$

$\rightarrow$

runends

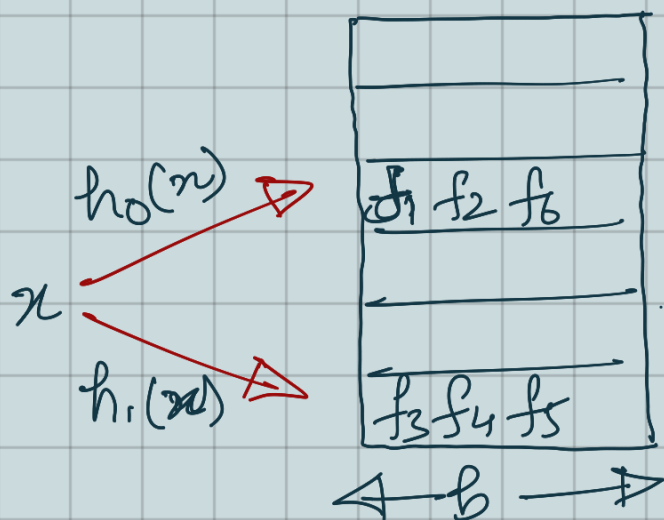|   |   | e | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|   |   | d | p | | | | | | | | | | |
|   |   | a | m | | | | | | | | | | |
|   |   | 1 | 1 | | | | | | | | | | |
|   |   | 1 | 1 | | | | | | | | | | |
|   | a | d | e | m | p | | | | | | | | |

$$\text{Space} : \approx n \lg\left(\frac{1}{\epsilon}\right) + 2.125\, n$$

$$\text{false positive rate} : \frac{1}{2^r}$$

# Cuckoo hashing:-

Typically,
$b = 4$.



① Compute $h_0(x)$ & $h_1(x)$
② Insert $f(x)$ into emptier block
③ Kick an item if needed.

Note:- $h_0(x)$ & $h_1(x)$ need to be dependent
      to support kicking

Space : $\approx n \lg\left(\dfrac{f}{\epsilon}\right) + 3n$

false positive rate: $\approx \dfrac{2b}{2^f}$