

# Data Structures and Algorithms

## Bloom Filters 3 & Cardinality Intro

CS 225

November 22, 2023

Brad Solomon



UNIVERSITY OF  
**ILLINOIS**  
URBANA - CHAMPAIGN

Department of Computer Science

Happy fall  
break!  
^ ^  
-

# Learning Objectives

Finish discussing bloom filters (and review bit vectors)

Applications  
↓

Introduce the concept of cardinality and cardinality estimation

# Bloom Filter: Error Rate

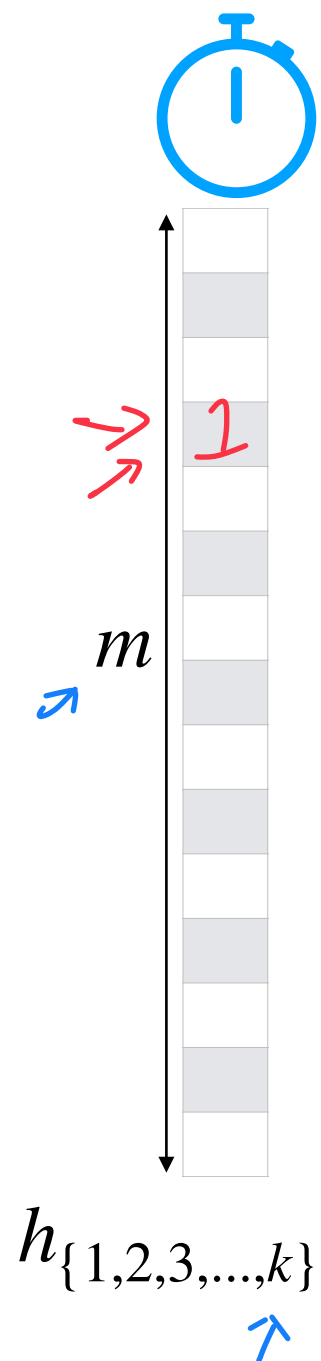
Given bit vector of size  $m$  and  $k$  SUHA hash function

What is our expected FPR after  $n$  objects are inserted?

The probability my bit is 1 after  $n$  objects inserted

$$\left( 1 - \left( 1 - \frac{1}{m} \right)^{nk} \right)^k$$

The number of [assumed independent] trials



# Bloom Filter: Error Rate

Vector of size  $m$ ,  $k$  SUHA hash function, and  $n$  objects

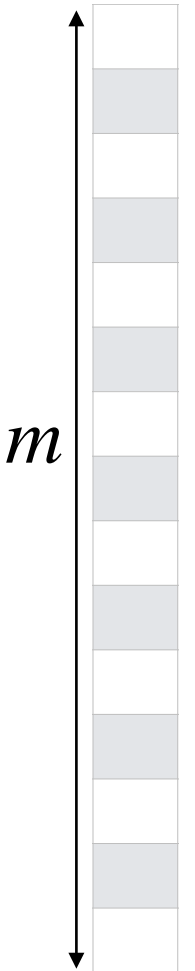
**To minimize the FPR, do we prefer...**

**(A) large  $k$**

**(B) small  $k$**

$$\left( 1 - \left( 1 - \frac{1}{m} \right)^{nk} \right)^k$$

$h_{\{1,2,3,\dots,k\}}$



# Bloom Filter: Error Rate

Vector of size  $m$ ,  $k$  SUHA hash function, and  $n$  objects

**(A) large  $k$**

$$\left(1 - \left(1 - \frac{1}{m}\right)^{nk}\right)^k$$

As  $k$  increases, this gets smaller!

↓ FPR approaches 1



$\equiv 100\%$  FPR

**(B) small  $k$**

$$\left(1 - \left(1 - \frac{1}{m}\right)^{nk}\right)^k$$

As  $k$  decreases, this gets smaller!

↳ Repeated random trials



# Bloom Filter: Optimal Error Rate

$k$  = hash  
 $M$  = BF size  
 $n$  = items

To build the optimal hash function, fix  $m$  and  $n$ !

**Claim:** The optimal hash function is when  $k^* = \ln 2 \cdot \frac{m}{n}$

$\alpha = \frac{n}{m}$

(1)  $\left(1 - \left(1 - \frac{1}{m}\right)^{nk}\right)^k \approx \left(1 - e^{\frac{-nk}{m}}\right)^k$

Shape of eq more clear

(2)  $\frac{d}{dk} \left(1 - e^{\frac{-nk}{m}}\right)^k \approx \frac{d}{dk} \left(k \ln\left(1 - e^{\frac{-nk}{m}}\right)\right)$

# Bloom Filter: Optimal Error Rate


**Claim 1:**  $\left(1 - \left(1 - \frac{1}{m}\right)^{nk}\right)^k \approx \left(1 - e^{-\frac{nk}{m}}\right)^k$

$$\left(1 - \frac{1}{m}\right)^{nk} = e^{\ln\left[\left(1 - \frac{1}{m}\right)^{nk}\right]}$$

$$x = e^{\ln(x)^y} = e^{y \ln(x)}$$

# Bloom Filter: Optimal Error Rate

**Claim 1:**  $\left(1 - \left(1 - \frac{1}{m}\right)^{nk}\right)^k \approx \left(1 - e^{-\frac{nk}{m}}\right)^k$

$$\begin{aligned} \left(1 - \frac{1}{m}\right)^{nk} &= e^{\ln\left[\left(1 - \frac{1}{m}\right)^{nk}\right]} \\ &= e^{\ln\left[\left(1 - \frac{1}{m}\right)\right]nk} \end{aligned}$$




# Bloom Filter: Optimal Error Rate

Taylor's expansion of  $\ln(1 + x)$ :

"Mercator Series"

$$x - \frac{x^2}{2} + \frac{x^3}{3} - \frac{x^4}{4} + \dots$$

$$e^{\ln\left(1 - \frac{1}{m}\right)^{nk}}$$

$$x = -\frac{1}{m}$$

↑  
This is small



$$\left(1 - \frac{1}{m}\right)^{nk} \approx e^{\frac{-nk}{m}}$$

# Bloom Filter: Optimal Error Rate

**Claim 1:**  $\left(1 - \left(1 - \frac{1}{m}\right)^{nk}\right)^k \approx \left(1 - e^{\frac{-nk}{m}}\right)^k$

$$\left(1 - \frac{1}{m}\right)^{nk} = e^{\ln\left[\left(1 - \frac{1}{m}\right)^{nk}\right]}$$

$$= e^{\ln\left[1 - \frac{1}{m}\right]nk}$$

$$\approx e^{\frac{-nk}{m}}$$



# Bloom Filter: Optimal Error Rate

**Claim 2:**  $\frac{d}{dk} \left( 1 - e^{-\frac{nk}{m}} \right)^k \approx \frac{d}{dk} \left( k \ln(1 - e^{-\frac{nk}{m}}) \right)$

**Fact:**  $\frac{d}{dx} \ln f(x) = \frac{1}{f(x)} \frac{df(x)}{dx}$

**TL;DR:**  $\min [f(x)] = \min [\ln f(x)]$

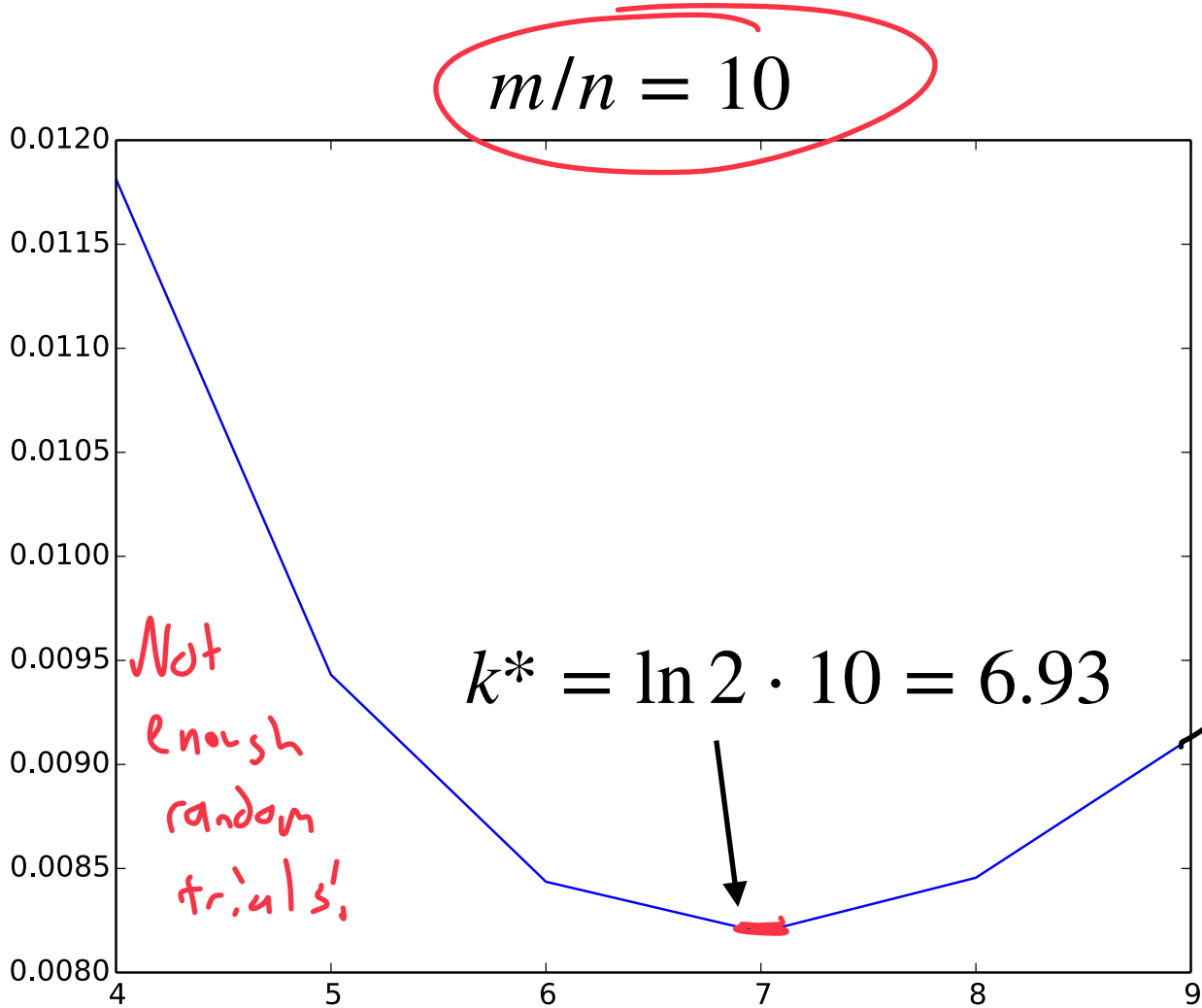
Derivative is zero when  $k^* = \ln 2 \cdot \frac{m}{n}$

Take  
deriv  
of  
 $\ln$

# Bloom Filter: Error Rate



FPR  $\left(1 - e^{-\frac{nk}{m}}\right)^k$



Not enough random trials!

$k^* = \ln 2 \cdot 10 = 6.93$

BF becomes too saturated w/ 1s

k small ☹️ ← # k hashes →

Figure by Ben Langmead

# Bloom Filter: Optimal Parameters

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

**Given any two values, we can optimize the third**

$n = 100$  items

$k = 3$  hashes

$3 = \ln 2 \left( \frac{m}{100} \right)$

$m =$

$\frac{300}{\ln 2} = \sim 433$  bits

$m = 100$  bits

$n = 20$  items

$k =$

$\ln 2 \left( \frac{100}{20} \right) \sim 3.5$  hashes  
3 ↙ ↘ 4

$m = 100$  bits

$k = 2$  items

$n =$

$\ln 2 \left( \frac{100}{2} \right) \sim 35$  items

# Bloom Filter: Optimal Parameters

$$m = \frac{nk}{\ln 2} \approx 1.44 \cdot nk$$

Optimal hash function is still  $O(n)$ !

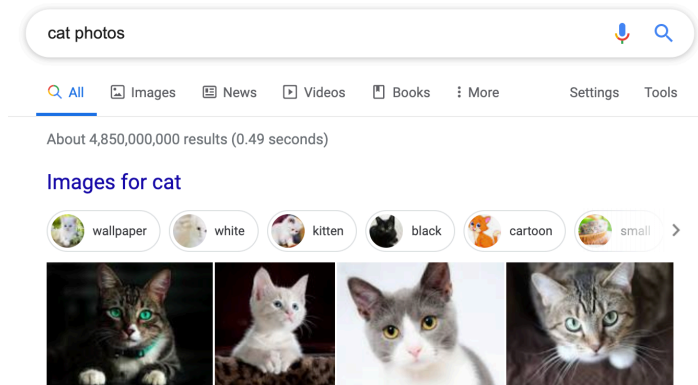
↳ Don't need to be optimal!



$n = 250,000$  files vs  $\sim 10^{15}$  nucleotides vs 260 TB

(character  
T, C, G, A)

↳  $1.44 \text{ bits} \times 250,000 = 45 \text{ KB} \sim 1 \text{ hash}$   
90 KB  $\approx$  hashes



$n = 60$  billion — 130 trillion

# Bloom Filters



A probabilistic data structure storing a set of values

$$h_{\{1,2,3,\dots,k\}}$$

Has three key properties:

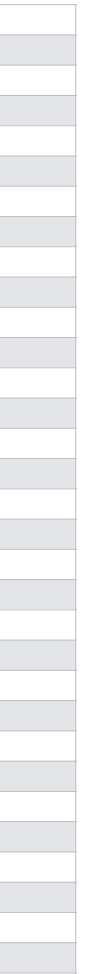
$k$ , number of hash functions

$n$ , expected number of insertions

$m$ , filter size in bits

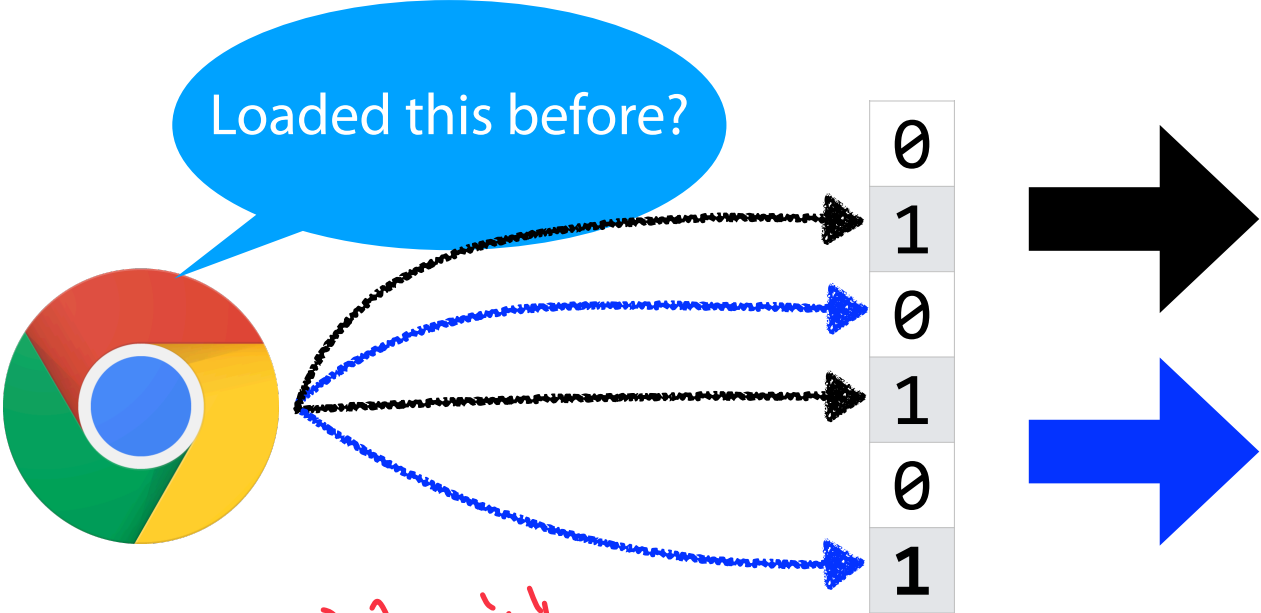
Expected false positive rate:  $\left(1 - \left(1 - \frac{1}{m}\right)^{nk}\right)^k \approx \left(1 - e^{\frac{-nk}{m}}\right)^k$

Optimal accuracy when:  $k^* = \ln 2 \cdot \frac{m}{n}$



# Bloom Filter: Website Caching

False Positive is not a problem!

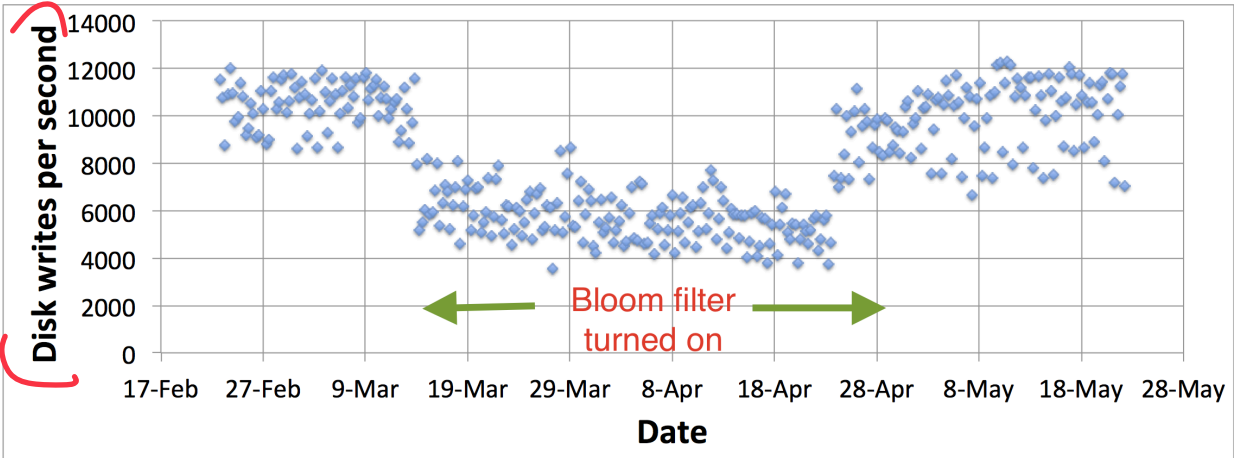


Cache this page!

↳ Caching a few extra pages isn't bad

Add to filter (but don't cache!)

↳ cache on 2nd visit  
 ↳ sometimes on first





# Bitwise Operators in C++

Array of  $n$  bits

How can we encode a bit vector in C++?

Vector<bool>

uint\_8 - 8 bits

||  
char - 8 bits

None are particularly good

work w/ extra steps!  
↳ Bit operations fix

Trivial - Vector<bool> doesn't make vector

↳ can't be multi-threaded    ↳ Doesn't implement parts of vector

..... | 0 1 | 2 | 3 | 4 5 6 7 | ← look up entire char  
..... | 1 0 | 0 | 1 0 1 0 | .....

# Bitwise Operators in C++

Traditionally, bit vectors are read from RIGHT to LEFT

**Warning: Lab\_Bloom won't do this**

64	32	16	8	4	2	2
↓	↓	↓	↓	↓	↓	↓
0	0	0	0	1	1	1

$= 1 + 2 + 4 = 7$

1	0	0	1	0	1	0
---	---	---	---	---	---	---

$= 64 + 8 + 2 = 74$

# Bitwise Operators in C++

↪ 3 is 1 or not

Let **A = 10110**  
Most      Least

Let **B = 01110**

~B: 10001

(invert all 1s and 0s)

A & B: 00110

(Both have to be 1 for 1  
else 0)

B & 00100

↪ 1 if 3 is 1  
0 otherwise

A | B: 11110

(if either is 1, set 1  
else 0)

A >> 2: 00101

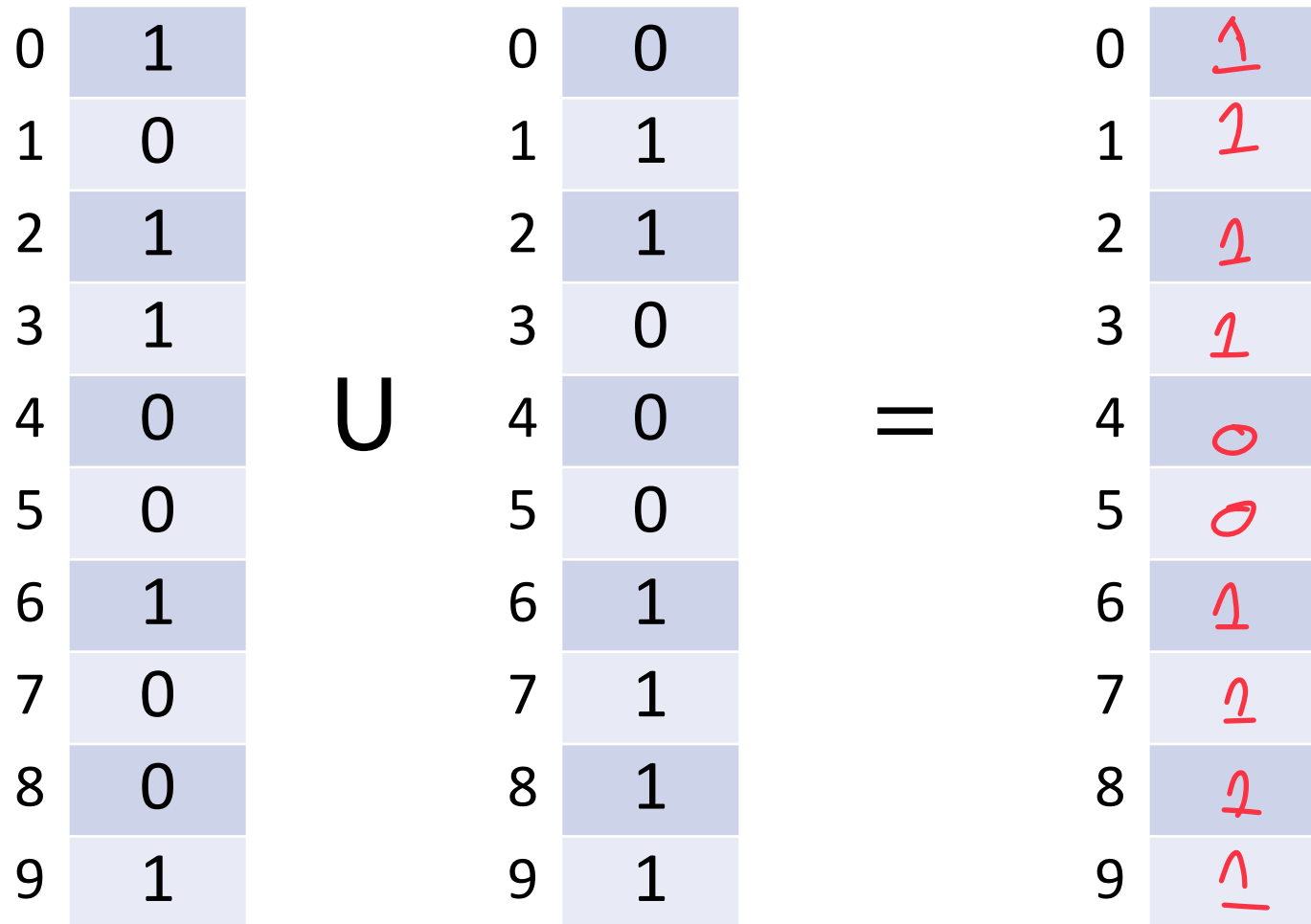
(discard least sig bits and replace most w 0)

B << 2: 11000

(discard most sig bits & replace least w 0)

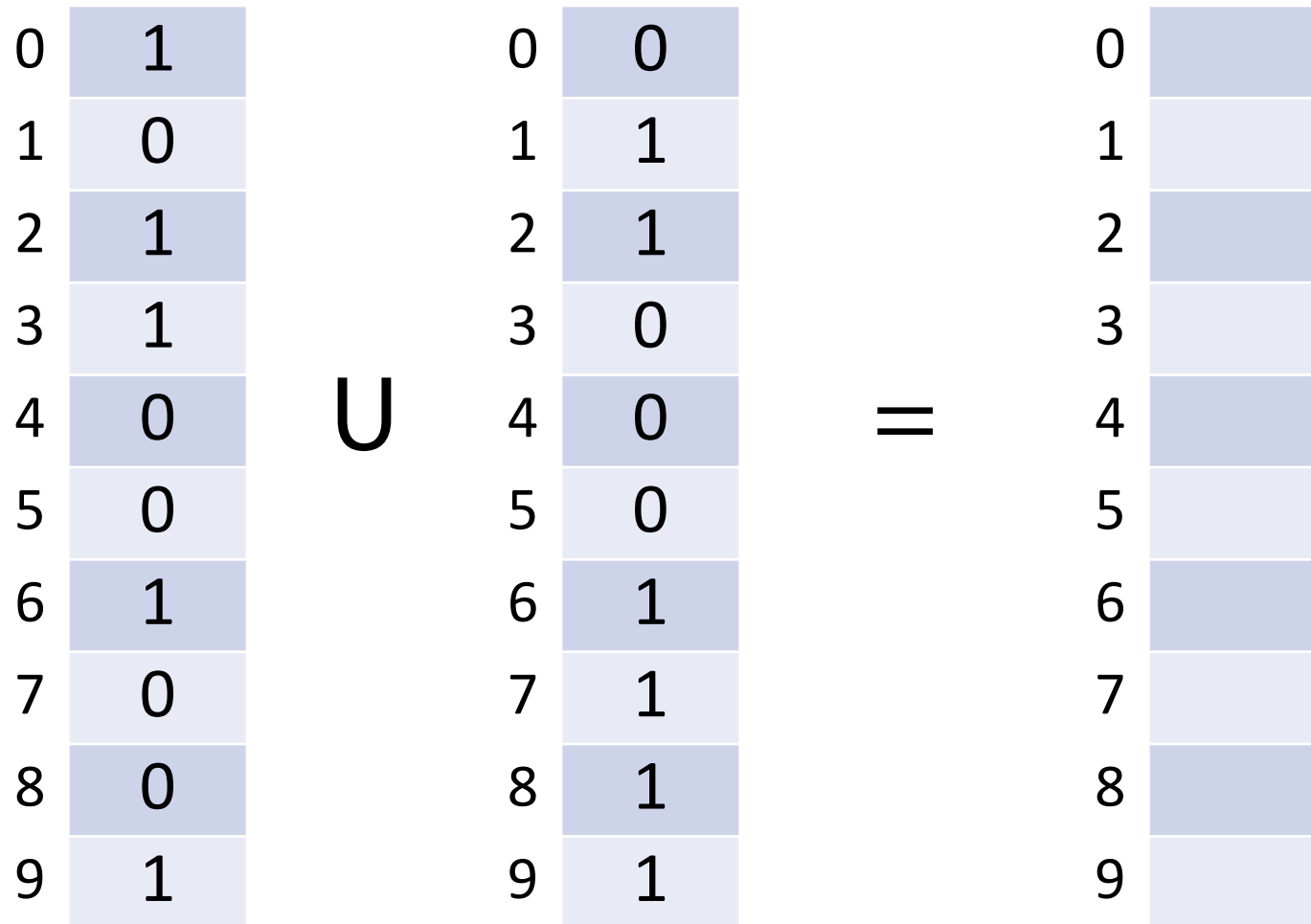
# Bit Vectors: Unioning

Bit Vectors can be trivially merged using bit-wise union.



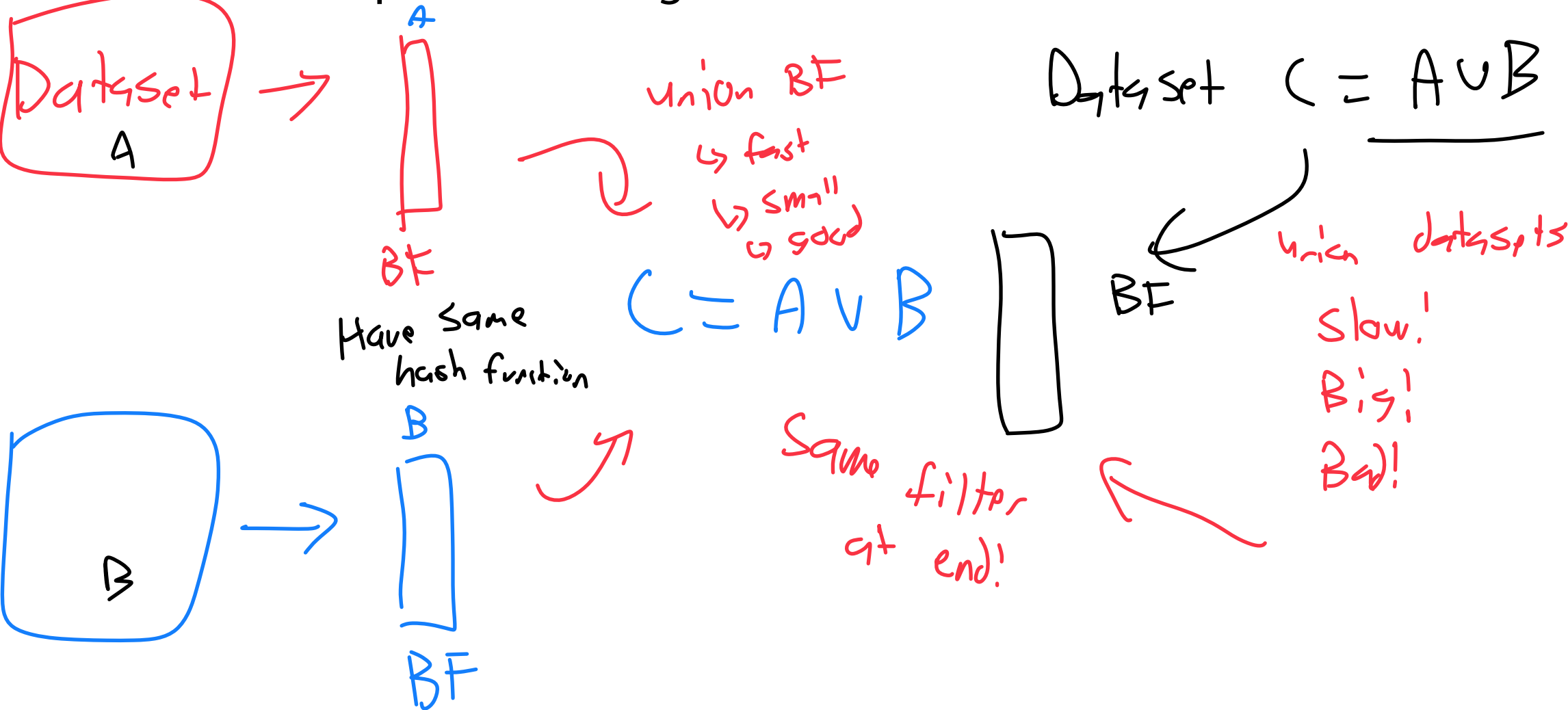
# Bit Vectors: Intersection

Bit Vectors can be trivially merged using bit-wise intersection.



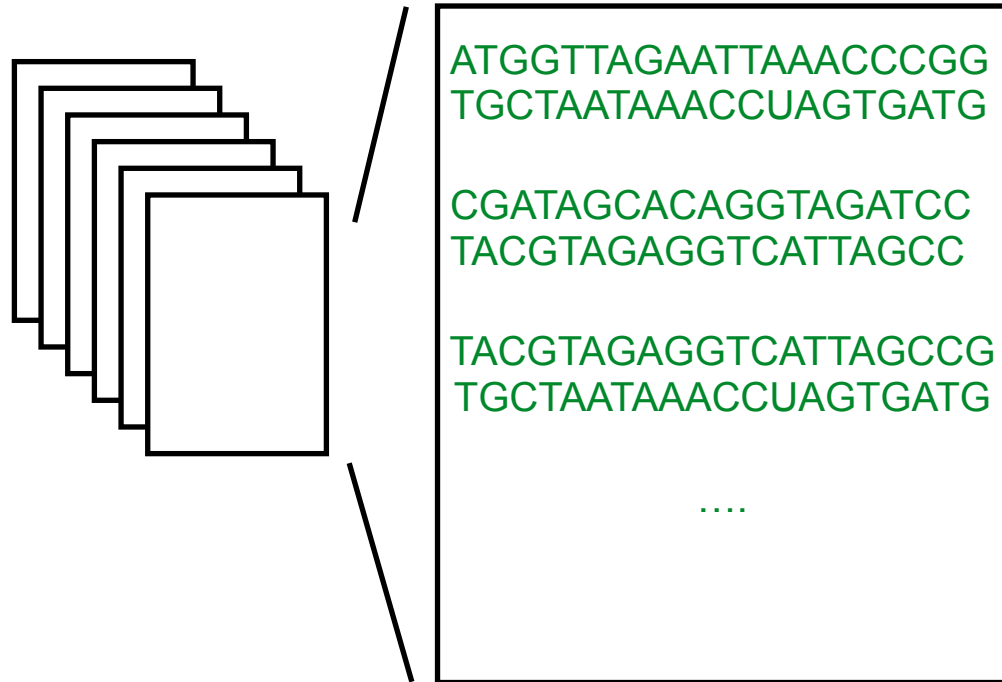
# Bit Vector Merging

What is the conceptual meaning behind **union** and **intersection**?

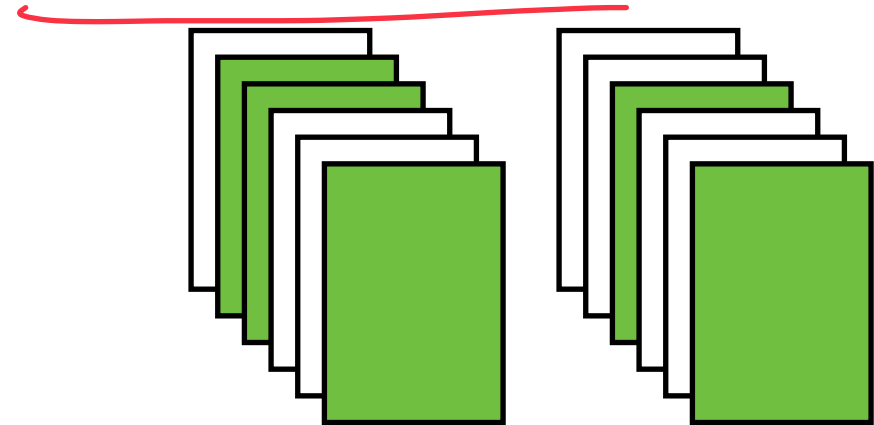


# Sequence Bloom Trees

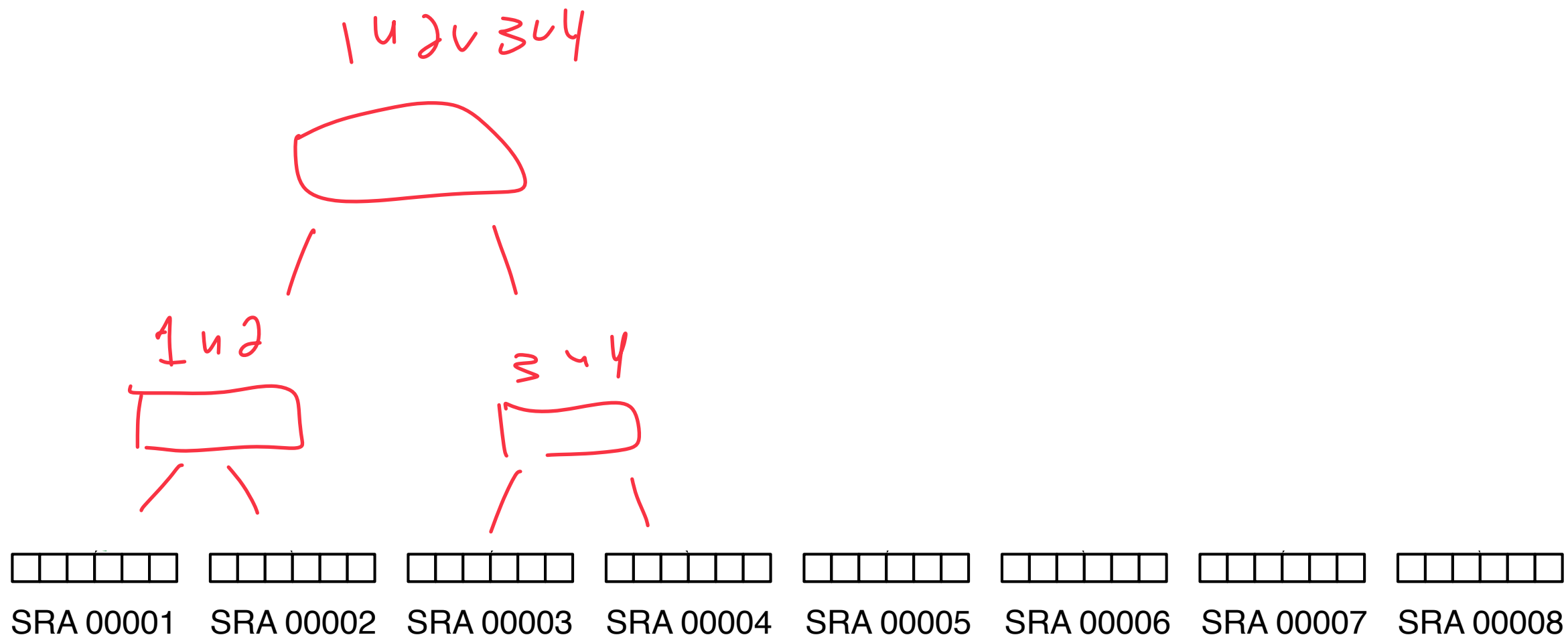
Imagine we have a large collection of text...



And our goal is to search these files for a query of interest...

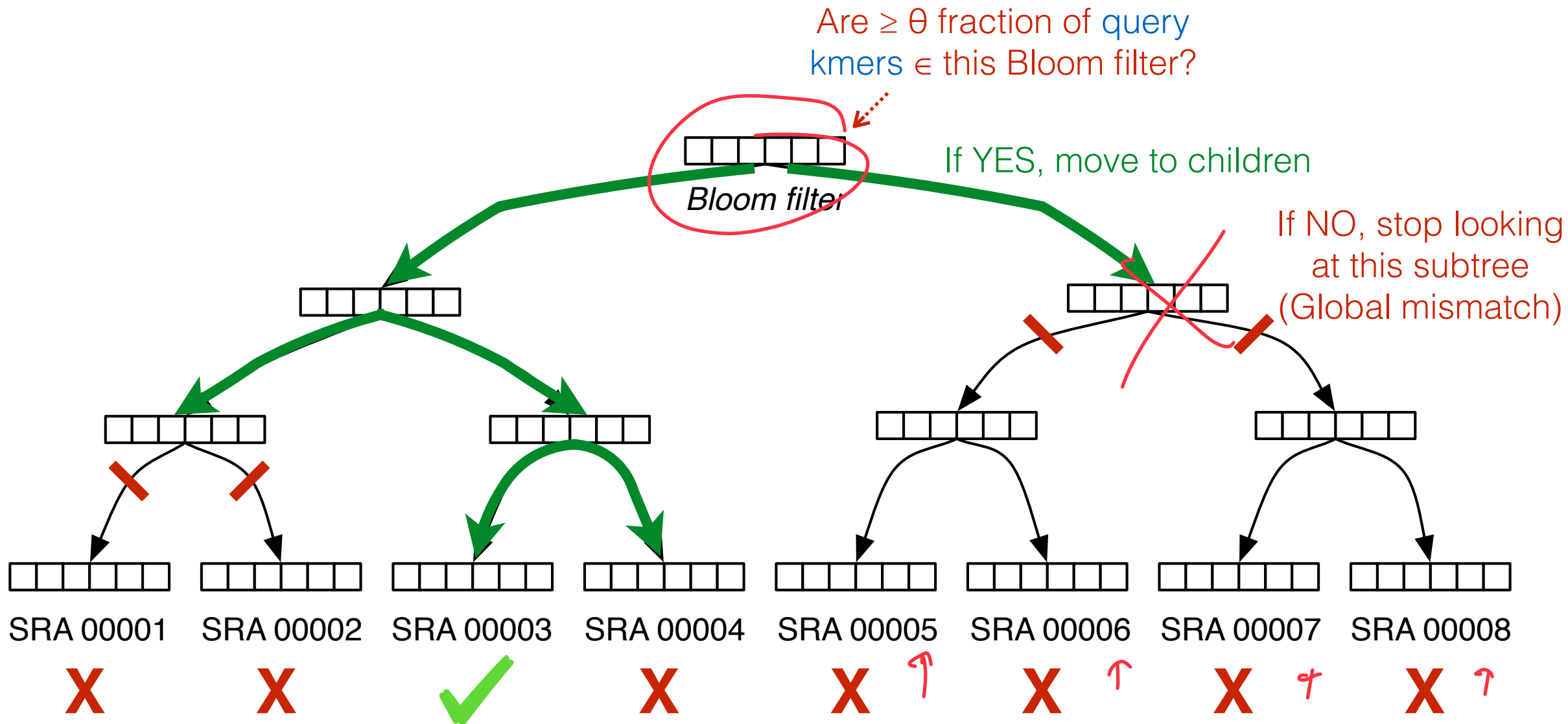


# Sequence Bloom Trees

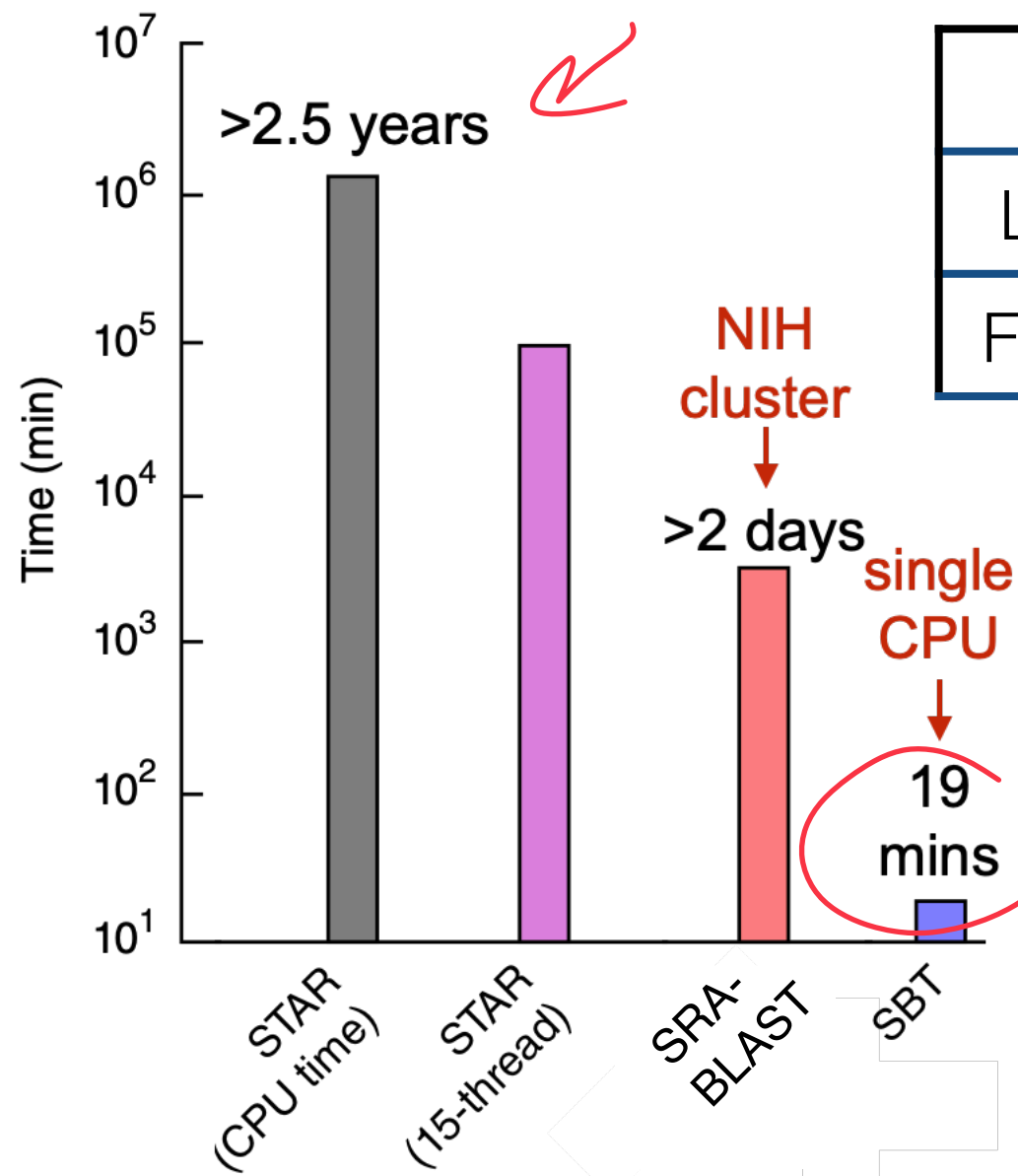




# Sequence Bloom Trees



# Sequence Bloom Trees



	SRA	FASTA.gz	SBT
Leaves	4966 GB	2692 GB	63 GB
Full Tree	-	-	200 GB

Solomon, Brad, and Carl Kingsford. "Fast search of thousands of short-read sequencing experiments." *Nature biotechnology* 34.3 (2016): 300-302.

Solomon, Brad, and Carl Kingsford. "Improved search of large transcriptomic sequencing databases using split sequence bloom trees." *International Conference on Research in Computational Molecular Biology*. Springer, Cham, 2017.

Sun, Chen, et al. "Allsome sequence bloom trees." *International Conference on Research in Computational Molecular Biology*. Springer, Cham, 2017.

Harris, Robert S., and Paul Medvedev. "Improved representation of sequence bloom trees." *Bioinformatics* 36.3 (2020): 721-727.

# Bloom Filters: Tip of the Iceberg



Cohen, Saar, and Yossi Matias. "Spectral bloom filters." *Proceedings of the 2003 ACM SIGMOD international conference on Management of data*. 2003.

Fan, Bin, et al. "Cuckoo filter: Practically better than bloom." *Proceedings of the 10th ACM International on Conference on emerging Networking Experiments and Technologies*. 2014.

Nayak, Sabuzima, and Ripon Patgiri. "countBF: A General-purpose High Accuracy and Space Efficient Counting Bloom Filter." *2021 17th International Conference on Network and Service Management (CNSM)*. IEEE, 2021.

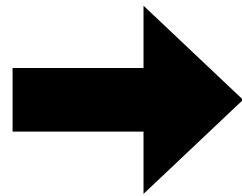
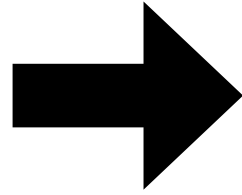
Mitzenmacher, Michael. "Compressed bloom filters." *IEEE/ACM transactions on networking* 10.5 (2002): 604-612.

Crainiceanu, Adina, and Daniel Lemire. "Bloofi: Multidimensional bloom filters." *Information Systems* 54 (2015): 311-324.

Chazelle, Bernard, et al. "The bloomier filter: an efficient data structure for static support lookup tables." *Proceedings of the fifteenth annual ACM-SIAM symposium on Discrete algorithms*. 2004.

There are many more than shown here...

# The hidden problem with (most) sketches... *S4MA!*



*Super Saturated!*

↳ Center

→ New strings not found in healthy people



# Cardinality

**Cardinality** is a measure of how many unique items are in a set

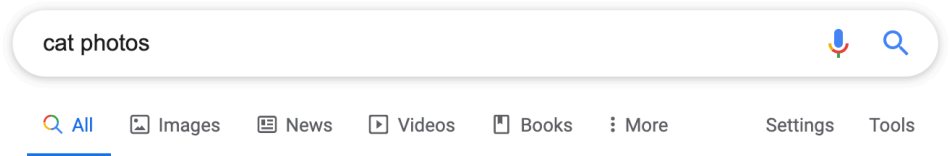
---

2
4
9
3
7
9
7
8
5
6

# Cardinality

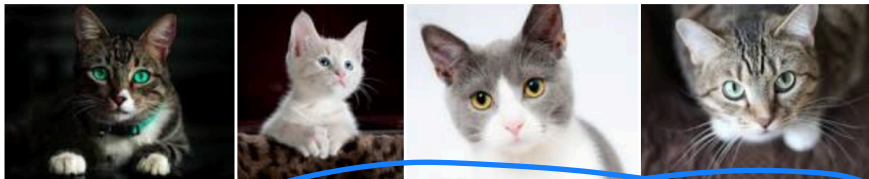
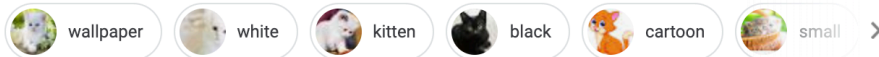
Sometimes its not possible or realistic to count all objects!

↳ Don't know data size!



About 4,850,000,000 results (0.49 seconds)

## Images for cat



Estimate: 60 billion — 130 trillion

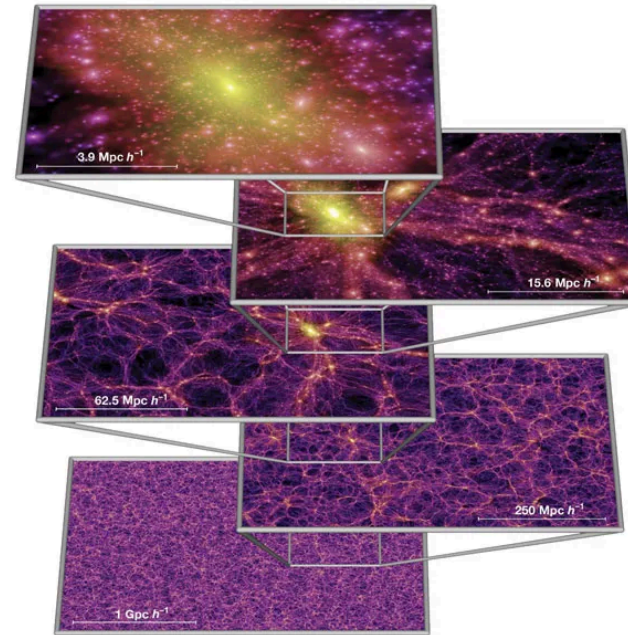


Image: <https://doi.org/10.1038/nature03597>

5581
8945
6145
8126
3887
8925
1246
8324
4549
9100
5598
8499
8970
3921
8575
4859
4960
42
6901
4336
9228
3317
399
6925
2660
2314

# Cardinality Estimation

1 - 1000

Draw is minimum  
↑

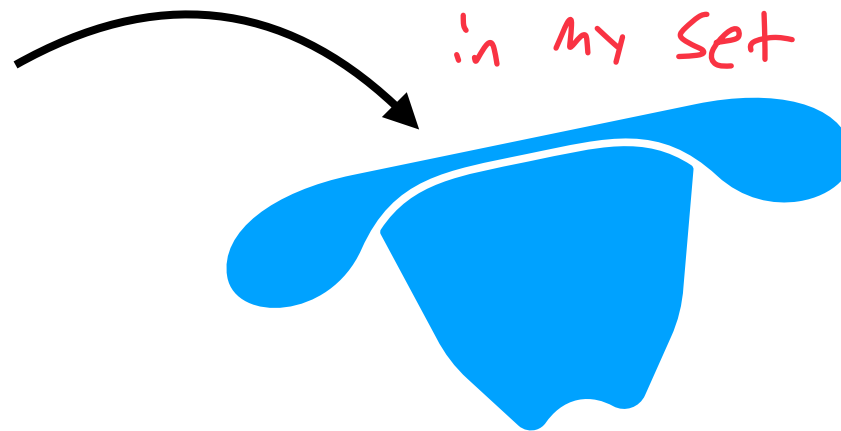
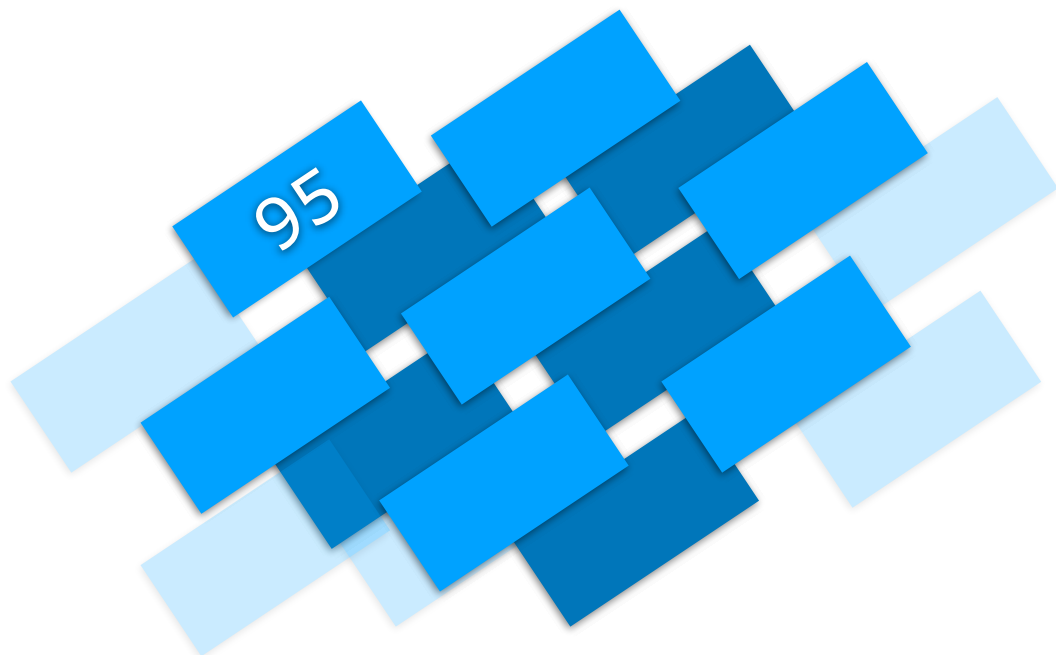
Imagine I fill a hat with numbered cards and draw one card out at random.

If I told you the value of the card was 95, what have we learned?

↳ 95 is in my set

All other cards  $\leq 95$

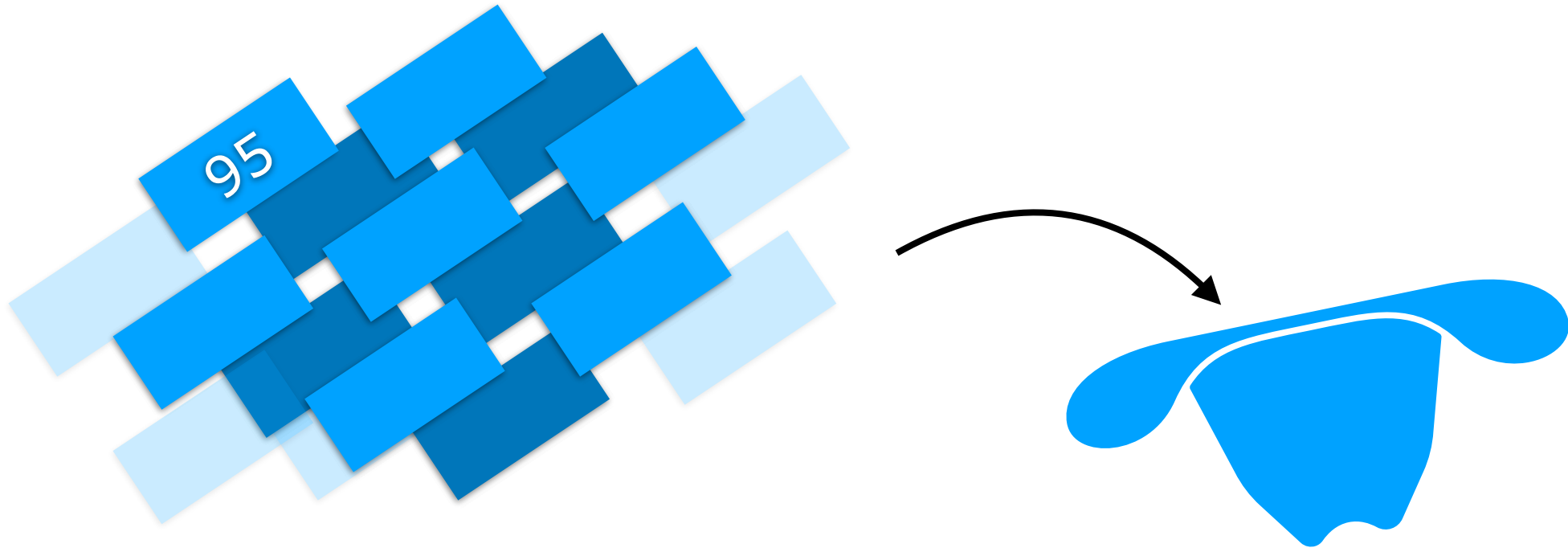
Claim: There are 10 items in my set



# Cardinality Estimation

Imagine I fill a hat with a **random subset** of numbered cards **from 0 to 999**

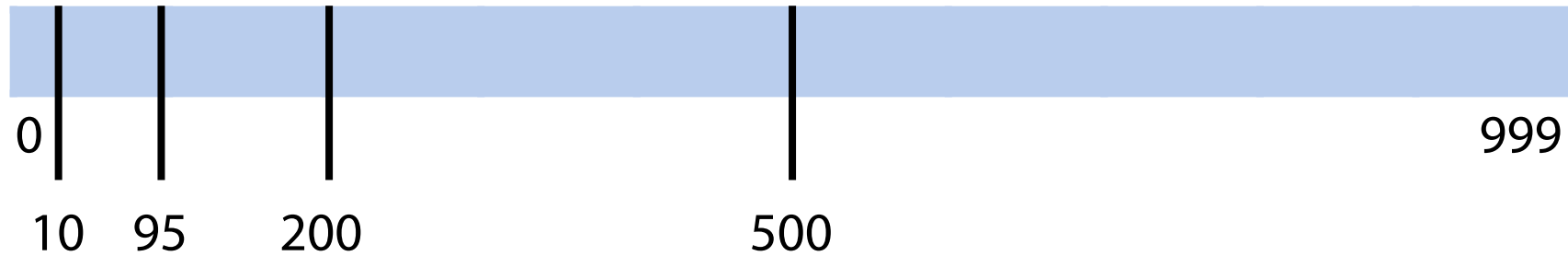
If I told you that the **minimum** value was 95, what have we learned?





# Cardinality Estimation

Imagine we have multiple uniform random sets with different minima.



# Cardinality Estimation

Let  $\min = 95$ . Can we estimate  $N$ , the cardinality of the set?



# Cardinality Estimation

Let  $\min = 95$ . Can we estimate  $N$ , the cardinality of the set?



**Claim:**  $95 \approx \frac{1000}{(N + 1)}$

# Cardinality Estimation



Let  $\min = 95$ . Can we estimate  $N$ , the cardinality of the set?



Conceptually: If we scatter  $N$  points randomly across the interval, we end up with  $N + 1$  partitions, each about  $1000/(N + 1)$  long

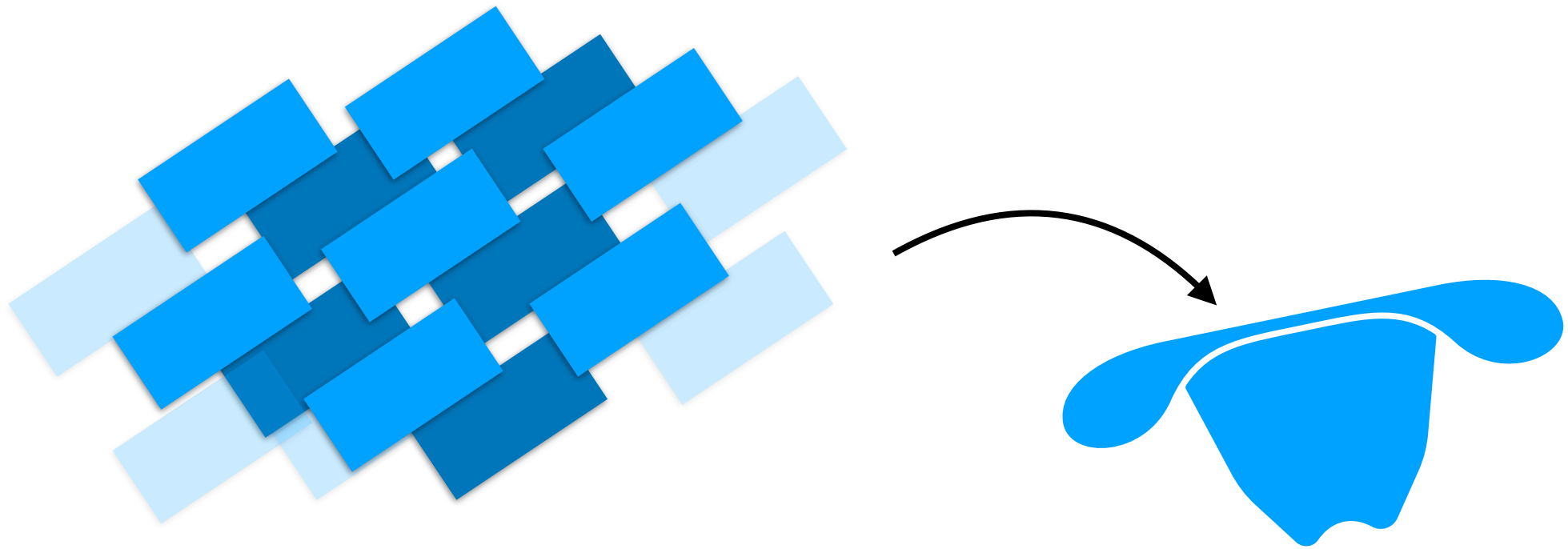
Assuming our first 'partition' is about average:  $95 \approx 1000/(N + 1)$

$$N + 1 \approx 10.5$$

$$N \approx 9.5$$

# Cardinality Estimation

Why do we care about “the hat problem”?





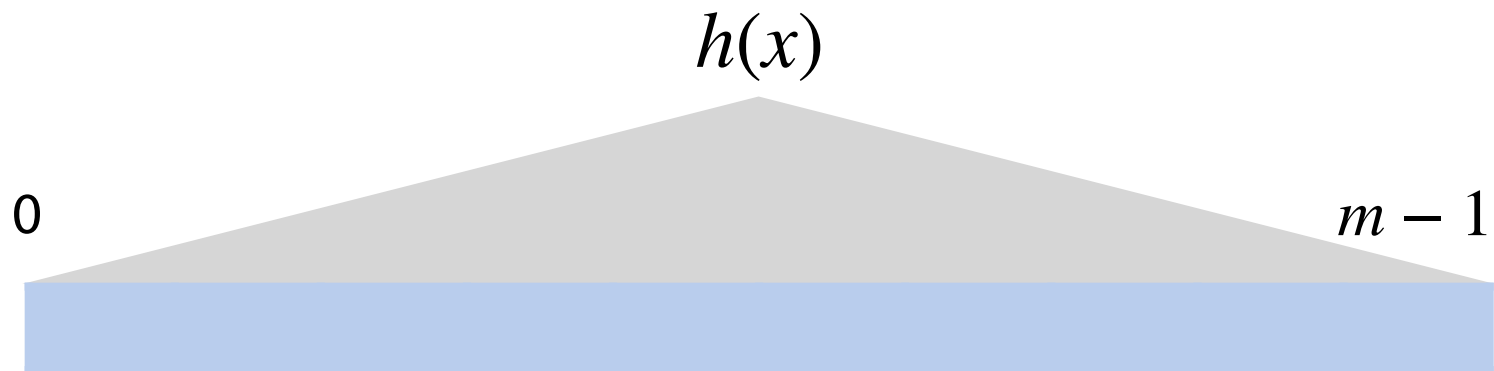
# Cardinality Estimation



Imagine we have a SUHA hash  $h$  over a range  $m$ .

Inserting a new key is equivalent to adding a card to our hat!

Tracking only the minimum value is a **sketch** that estimates the cardinality!



# Cardinality Estimation

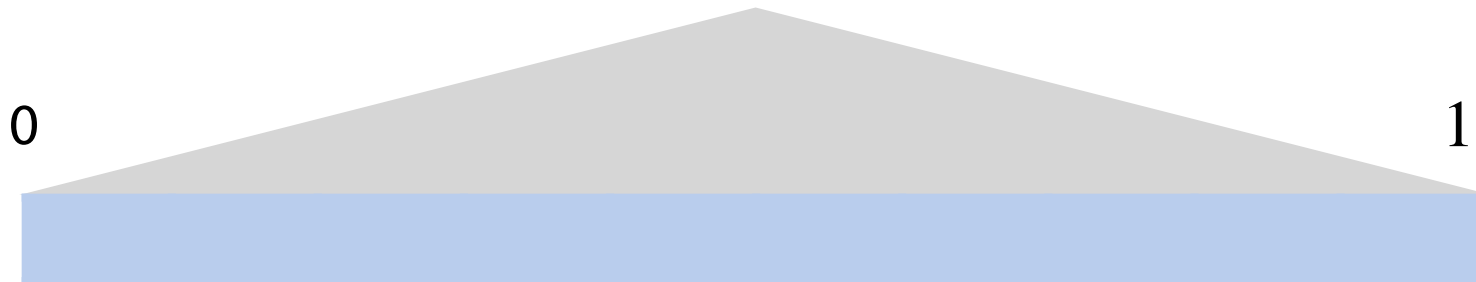
Imagine we have a SUHA hash  $h$  over a range  $m$ .

Inserting a new key is equivalent to adding a card to our hat!

Tracking only the minimum value is a **sketch** that estimates the cardinality!

To make the math work out, lets normalize our hash...

$$h'(x) = h(x) / (m - 1)$$





# Cardinality Sketch

Let  $M = \min(X_1, X_2, \dots, X_N)$  where each  $X_i \in [0, 1]$  is an uniform independent random variable

**Claim:**  $\mathbf{E}[M] = \frac{1}{N + 1}$

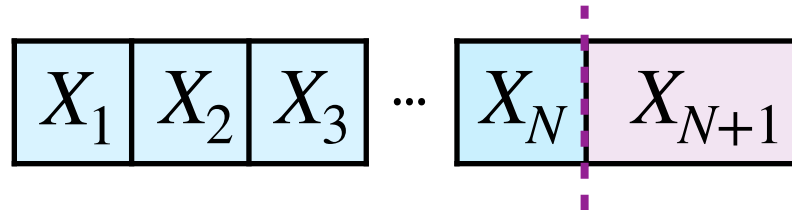
0

1



# Cardinality Sketch

Consider an  $N + 1$  draw:

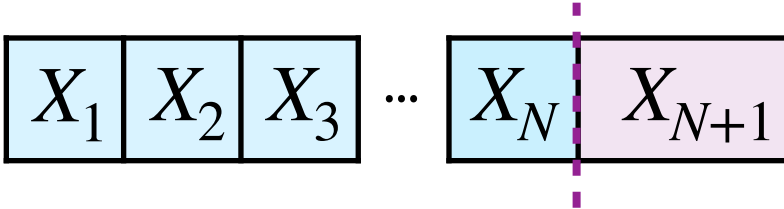


$$M = \min_{1 \leq i \leq N} X_i$$

$X_{N+1}$  can end up in one of two ranges:



# Cardinality Sketch

Consider an  $N + 1$  draw: 

$$M = \min_{1 \leq i \leq N} X_i$$

$X_{N+1}$  can end up in one of two ranges:

$X_{N+1}$  will be the new minimum with probability  $M$



# Cardinality Sketch

Consider an  $N + 1$  draw:  $X_1$   $X_2$   $X_3$  ...  $X_N$   $X_{N+1}$

$$M = \min_{1 \leq i \leq N} X_i$$

$X_{N+1}$  can end up in one of two ranges:

$X_{N+1}$  will be the new minimum with probability  $M$

$X_{N+1}$  will not change minimum with probability  $1 - M$



# Cardinality Sketch

Consider an  $N + 1$  draw:  $X_1$   $X_2$   $X_3$  ...  $X_N$   $X_{N+1}$

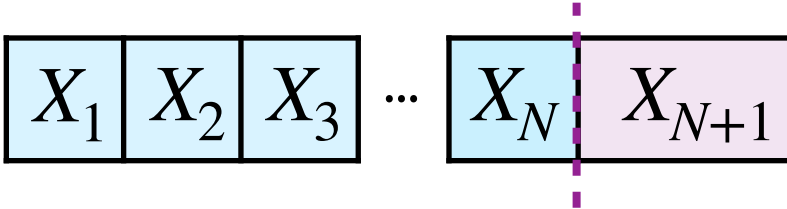
$$M = \min_{1 \leq i \leq N} X_i$$

$X_{N+1}$  will be the new minimum with probability  $M$

By definition of SUHA,  $X_{N+1}$  has a  $\frac{1}{N+1}$  chance of being smallest item



# Cardinality Sketch

Consider an  $N + 1$  draw: 

$$M = \min_{1 \leq i \leq N} X_i$$

$X_{N+1}$  will be the new minimum with probability  $M$

By definition of SUHA,  $X_{N+1}$  has a  $\frac{1}{N+1}$  chance of being smallest item

$$\text{Thus, } \mathbf{E}[M] = \frac{1}{N+1}$$



# Cardinality Sketch

**Claim:**  $\mathbf{E}[M] = \frac{1}{N+1}$        $N \approx \frac{1}{M} - 1$

**Attempt 1**

0.962	0.328	0.771	0.952	0.923
-------	-------	-------	-------	-------

**Attempt 2**

0.253	0.839	0.327	0.655	0.491
-------	-------	-------	-------	-------

**Attempt 3**

0.134	0.580	0.364	0.743	0.931
-------	-------	-------	-------	-------

# Cardinality Sketch

The minimum hash is a valid sketch of a dataset but can we do better?

0

1

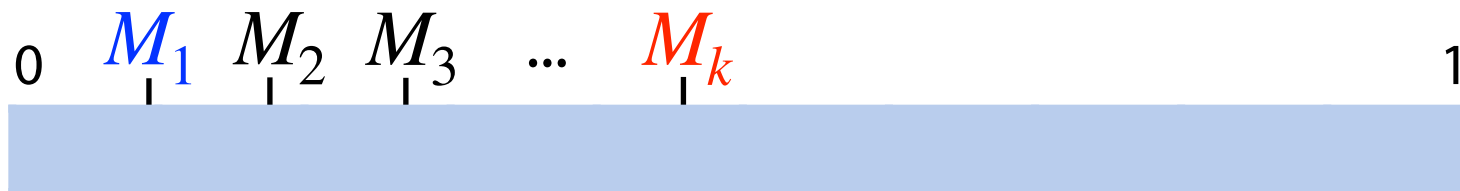




# Cardinality Sketch

**Claim:** Taking the  $k^{\text{th}}$ -smallest hash value is a better sketch!

**Claim:**  $\mathbf{E}[M_k] = \frac{k}{N + 1}$



# Cardinality Sketch

**Claim:** Taking the  $k^{\text{th}}$ -smallest hash value is a better sketch!

**Claim:** 
$$\frac{\mathbf{E}[M_k]}{k} = \frac{1}{N+1}$$

$$= \left[ \mathbf{E}[M_1] + (\mathbf{E}[M_2] - \mathbf{E}[M_1]) + \dots + (\mathbf{E}[M_k] - \mathbf{E}[M_{k-1}]) \right] \cdot \frac{1}{k}$$

$M_1$   
|

$M_2$   
|

$M_3$   
|

...

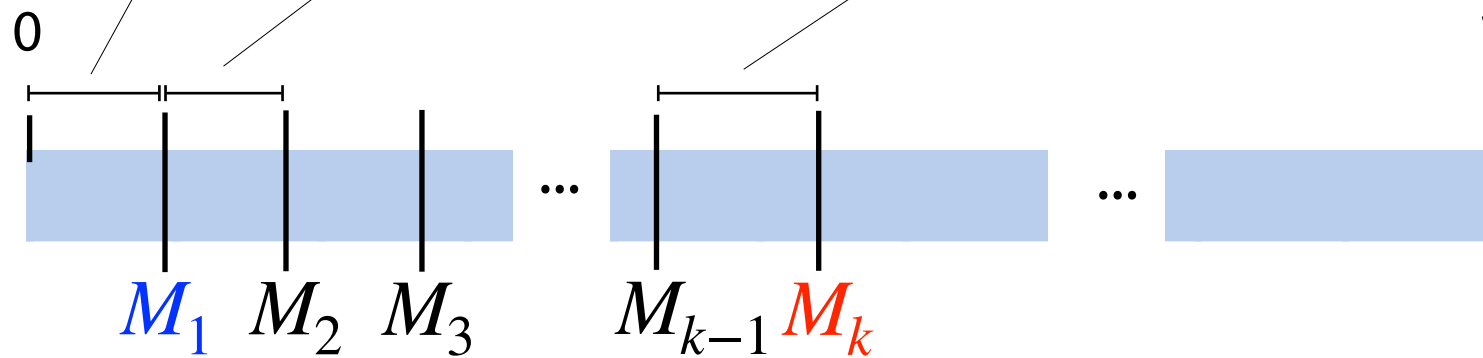
$M_{k-1}$   
|

$M_k$   
|

# Cardinality Sketch

$$\frac{1}{N+1} = \frac{\mathbf{E}[M_k]}{k}$$

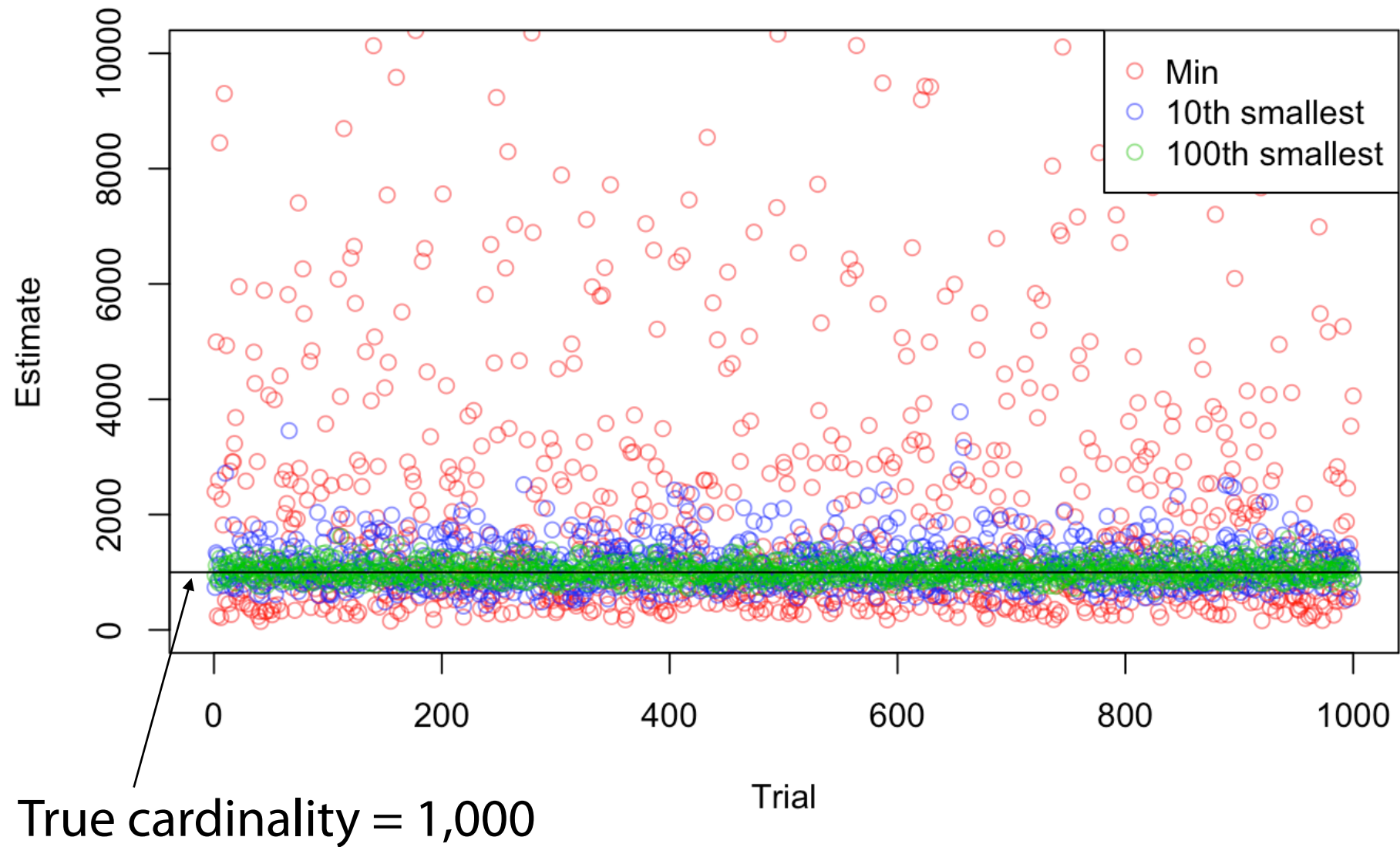
$$= \left[ \underbrace{\mathbf{E}[M_1]} + \underbrace{(\mathbf{E}[M_2] - \mathbf{E}[M_1])} + \dots + \underbrace{(\mathbf{E}[M_k] - \mathbf{E}[M_{k-1}])} \right] \cdot \frac{1}{k}$$



$k^{\text{th}}$  minimum  
value (KMV)

Averages  $k$  estimates for  $\frac{1}{N+1}$

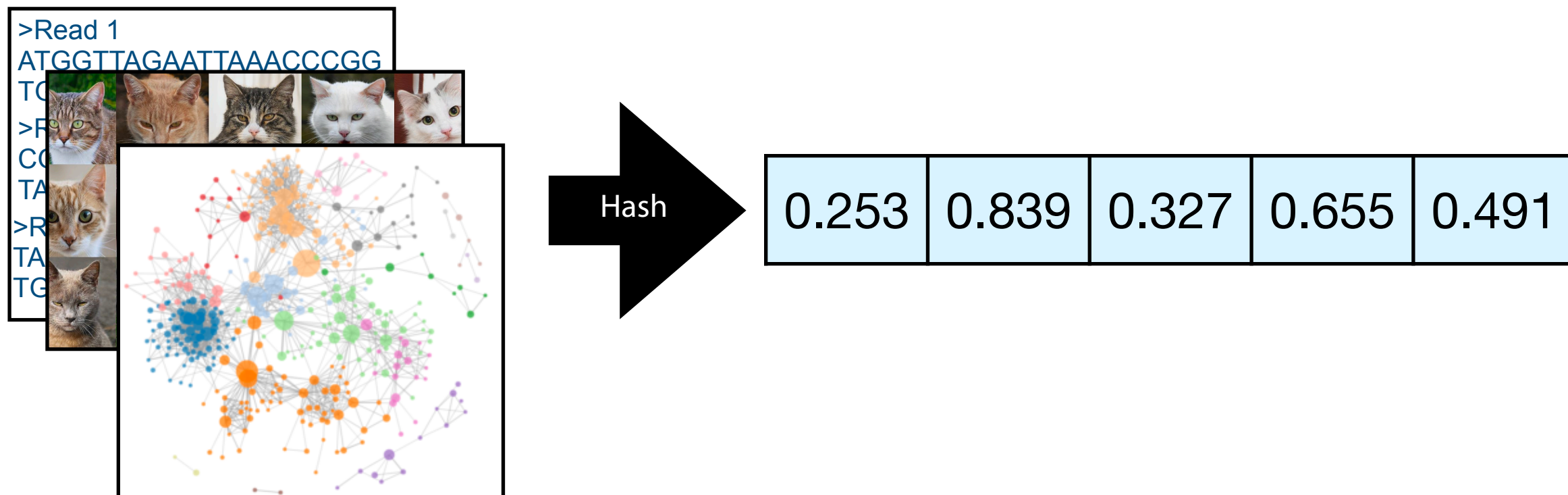
# Cardinality Sketch



# Cardinality Sketch



Given any dataset and a SUHA hash function, we can **estimate the number of unique items** by tracking the **k-th minimum hash value**.



To use the k-th min, we have to track k minima. **Can we use ALL minima?**

# Applied Cardinalities

Cardinalities

$$\frac{|A|}{|B|}$$

$$\frac{|A \cup B|}{|A \cap B|}$$

Set similarities

$$O = \frac{|A \cap B|}{\min(|A|, |B|)}$$

$$J = \frac{|A \cap B|}{|A \cup B|}$$

Real-world  
Meaning

```
AGGCCACAGTGTATTATGACTG
|||||          |||||
AGGCCACAGTGAGTTATGACTG
```

```
AAAAAAAAAAAGATGT-AAGTA
|||||          |||||
AAAAAAAAAAAGATGTAAAGTA
```

```
GAGG--TCAGATTCACAGCCAC
||||  |||||
GAGGGGTCAGATTCACAGCCAC
```

