

Data Structures and Algorithms

Cardinality

CS 225

Brad Solomon

May 1, 2026



UNIVERSITY OF
ILLINOIS
URBANA - CHAMPAIGN

Department of Computer Science

Learning Objectives

Review bloom filters and identify the 'weakness' of BFs

Introduce the concept of cardinality and cardinality estimation

Bloom Filters



A probabilistic data structure storing a set of values

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

Has three key properties:

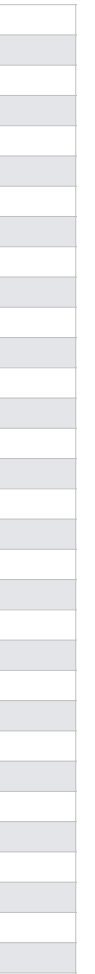
k , number of hash functions

n , expected number of insertions

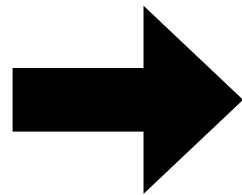
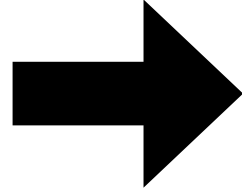
m , filter size in bits

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



The hidden problem with (most) sketches...



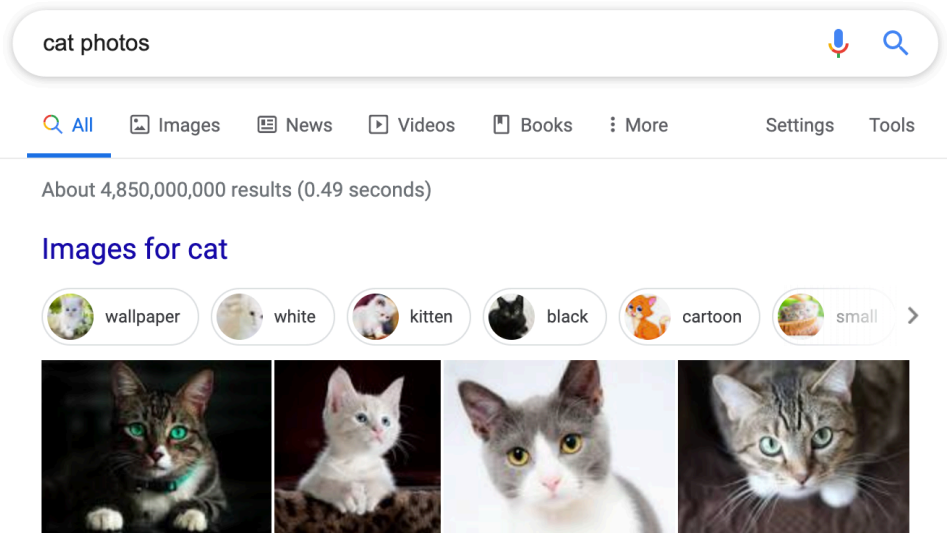
Cardinality

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!



Estimate: 60 billion — 130 trillion

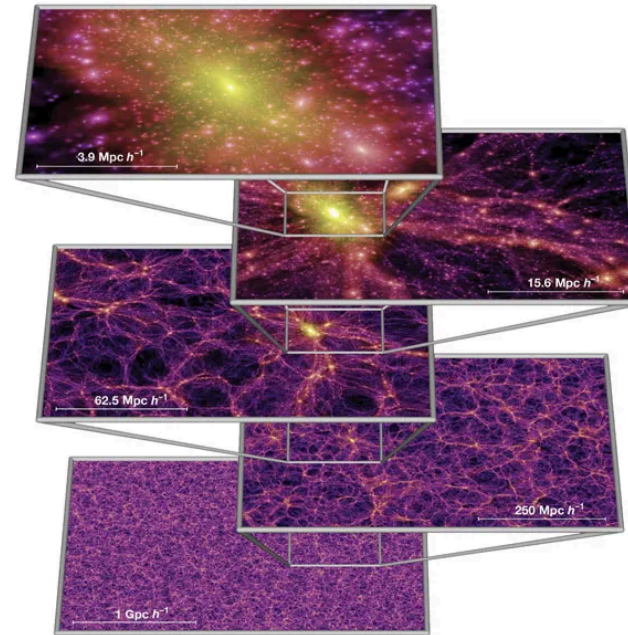


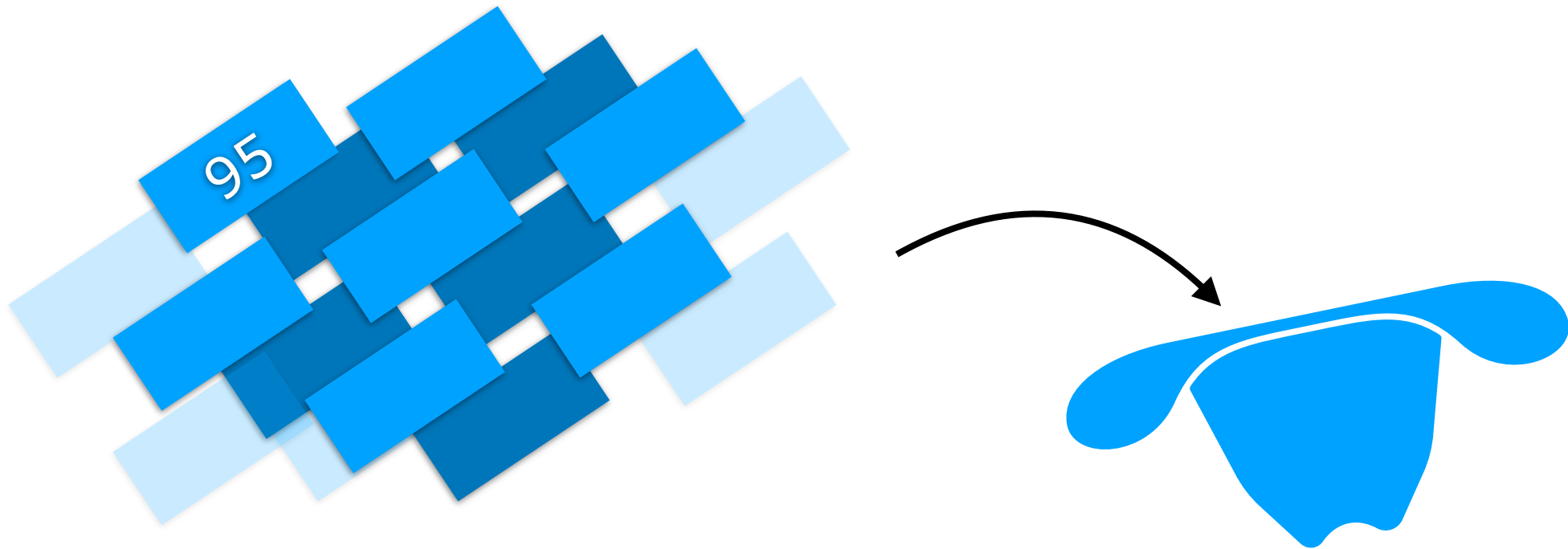
Image: <https://doi.org/10.1038/nature03597>

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

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

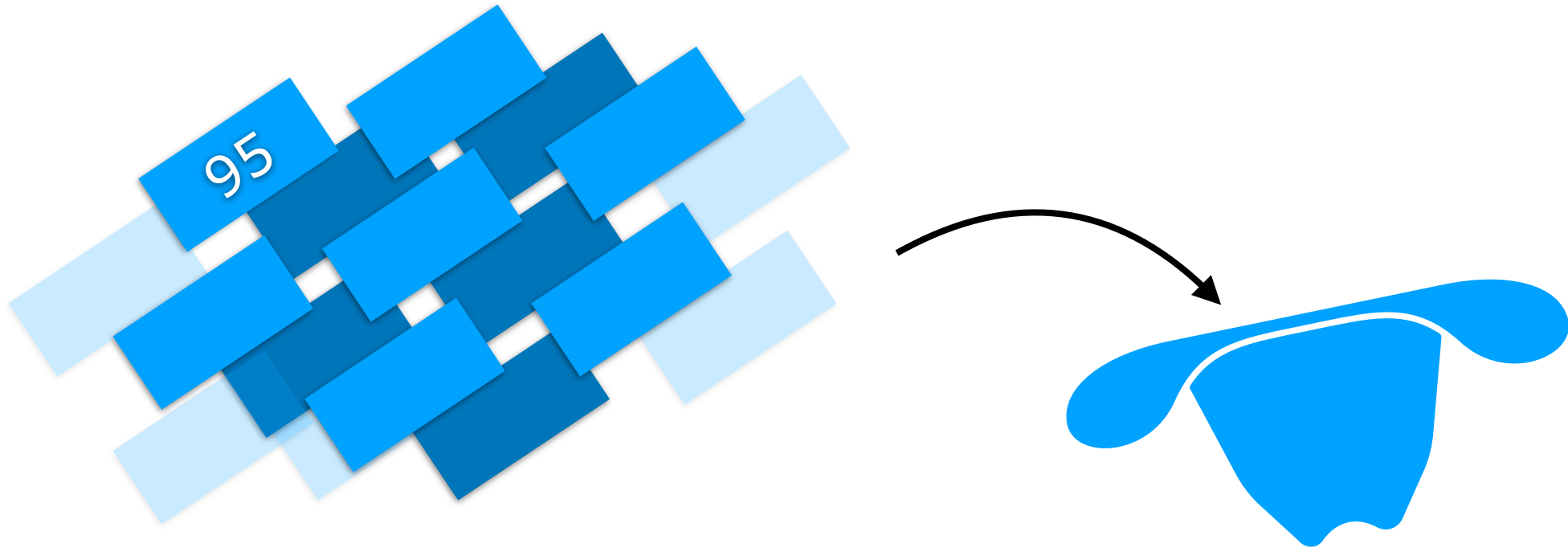
If I told you the value of the card was 95, what have we learned?



Cardinality Estimation

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

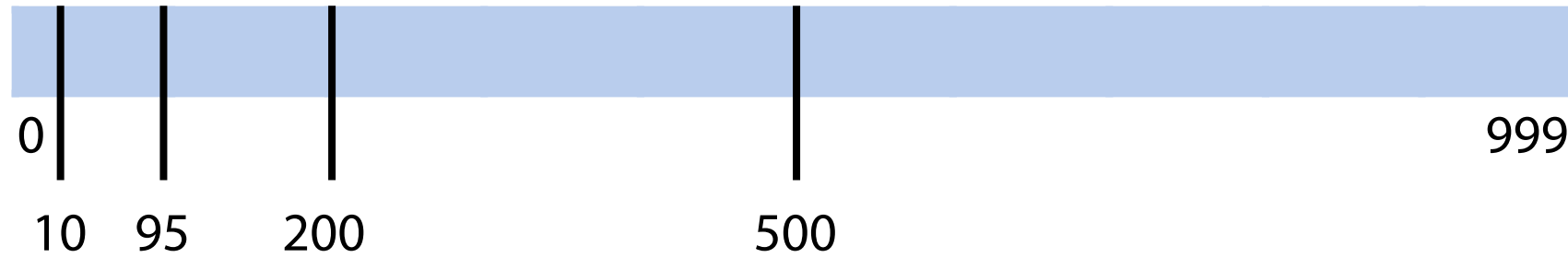
If I told you that the **minimum** value was 95, what have we learned?



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 $\text{min} = 95$. 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

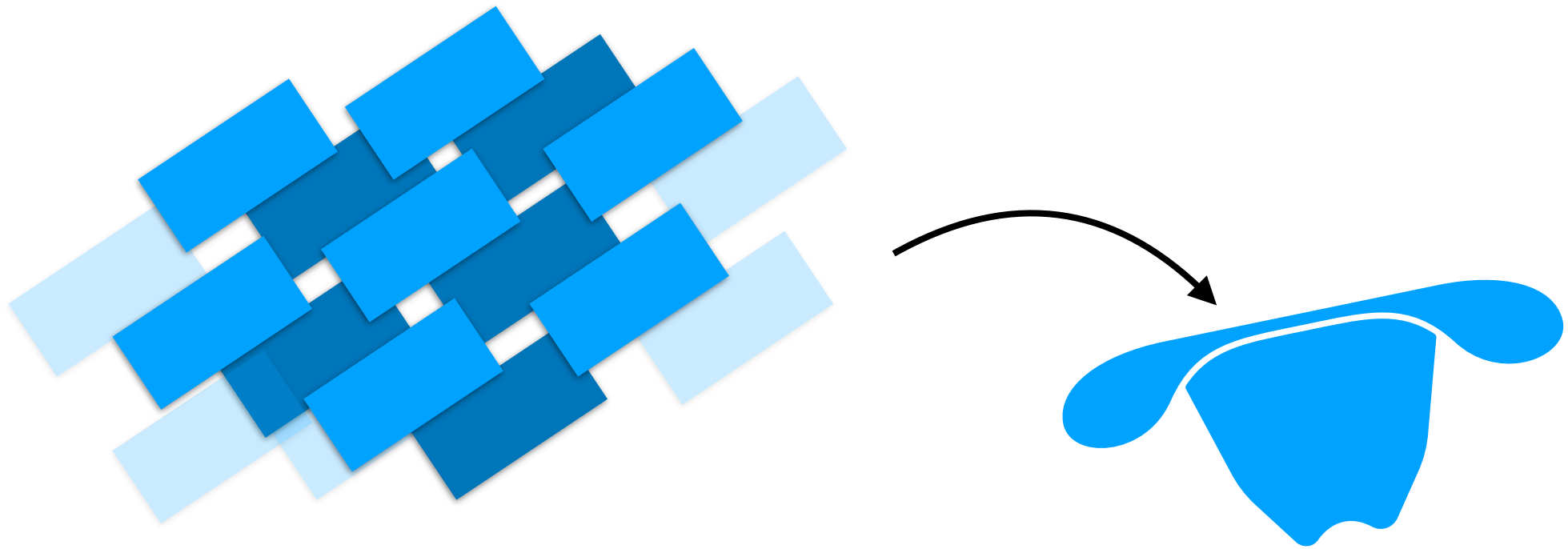
Assuming our first 'partition' is about average: $95 \approx 1000/(N + 1)$

$$N + 1 \approx 10.5$$

$$N \approx 9.5$$

Cardinality Estimation

Why do we care about “the hat problem”?



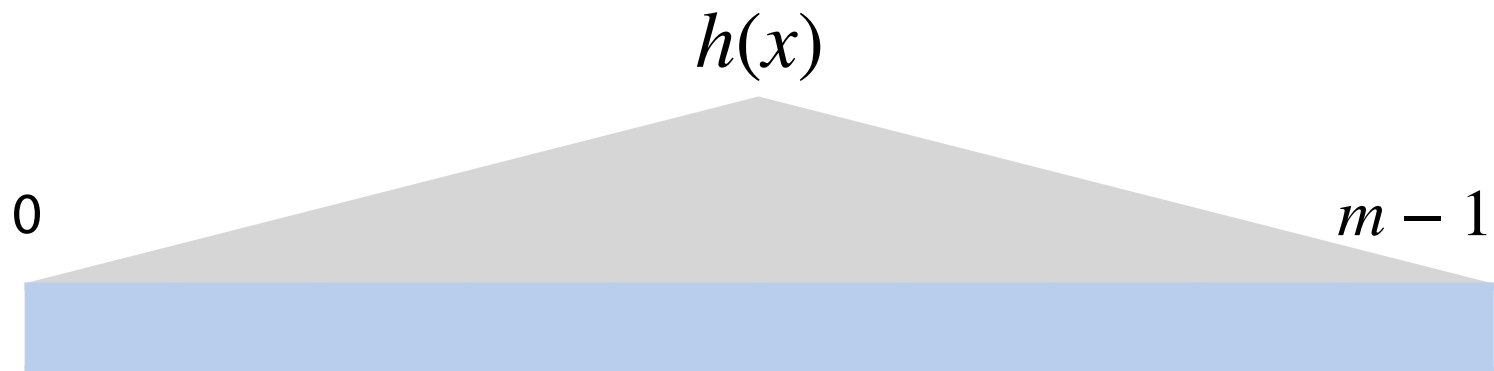
Cardinality Estimation



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

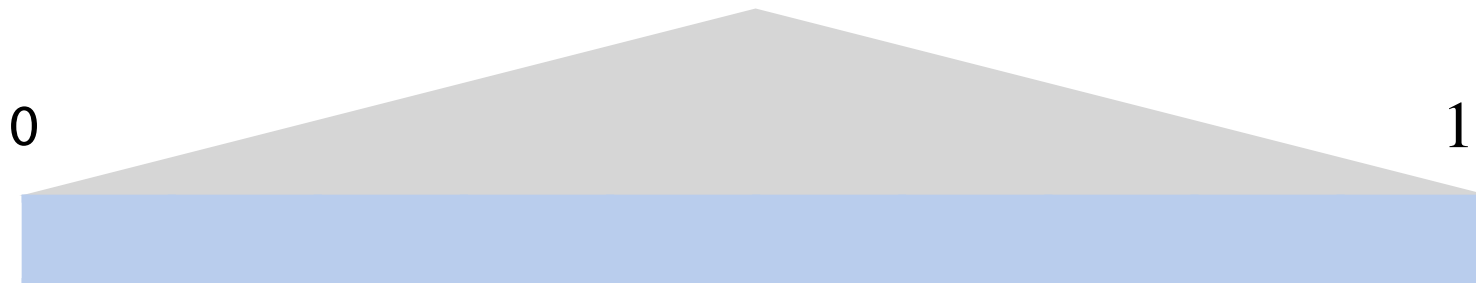
Imagine we have a SUHA hash h over a range m .

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

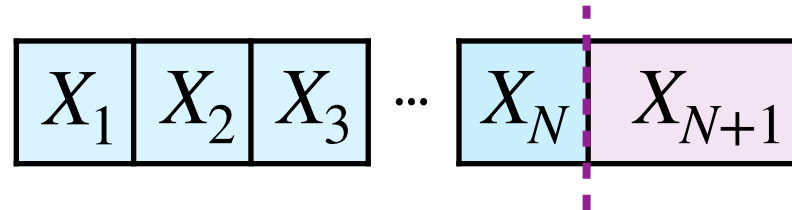
0

1



Cardinality Sketch

Consider an $N + 1$ draw:



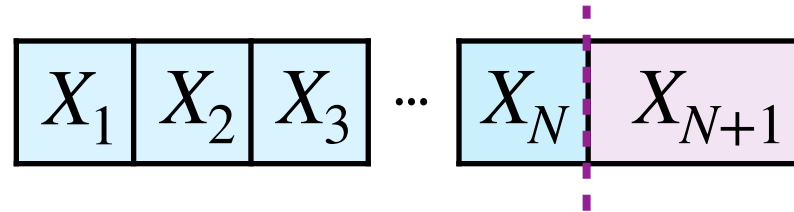
$$M = \min_{1 \leq i \leq N} X_i$$

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}

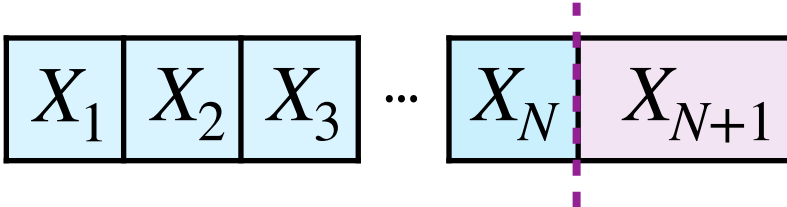
$$M = \min_{1 \leq i \leq N} X_i$$

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

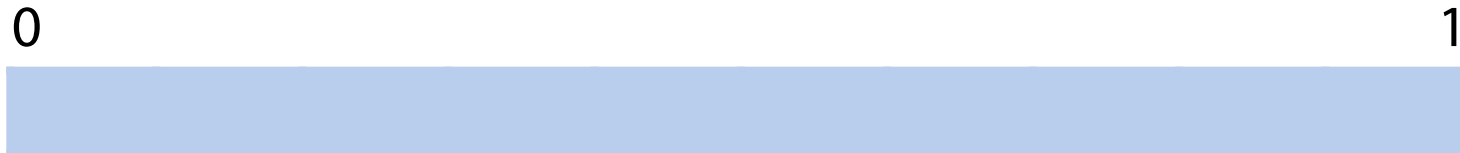
0.253	0.839	0.327	0.655	0.491
-------	-------	-------	-------	-------

Attempt 3

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

Cardinality Sketch

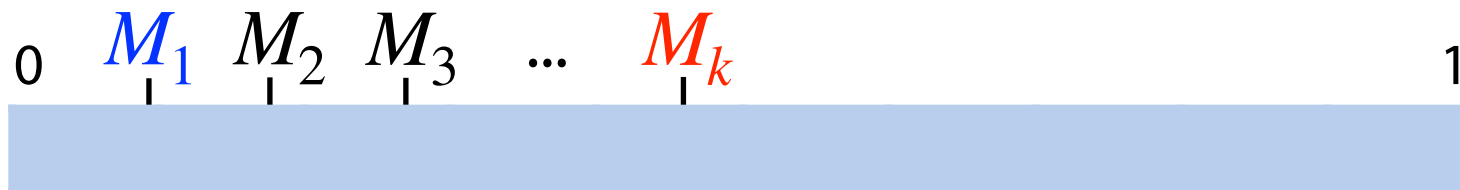
The minimum hash is a valid sketch of a dataset but can we do better?



Cardinality Sketch

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

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



Cardinality Sketch

Claim: Taking the k^{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
|

...

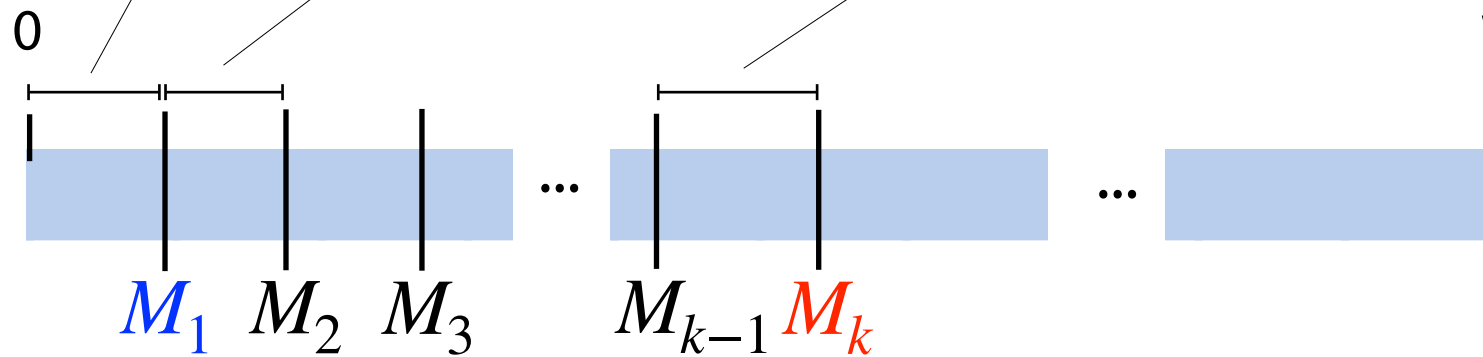
M_{k-1}
|

M_k
|

Cardinality Sketch

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

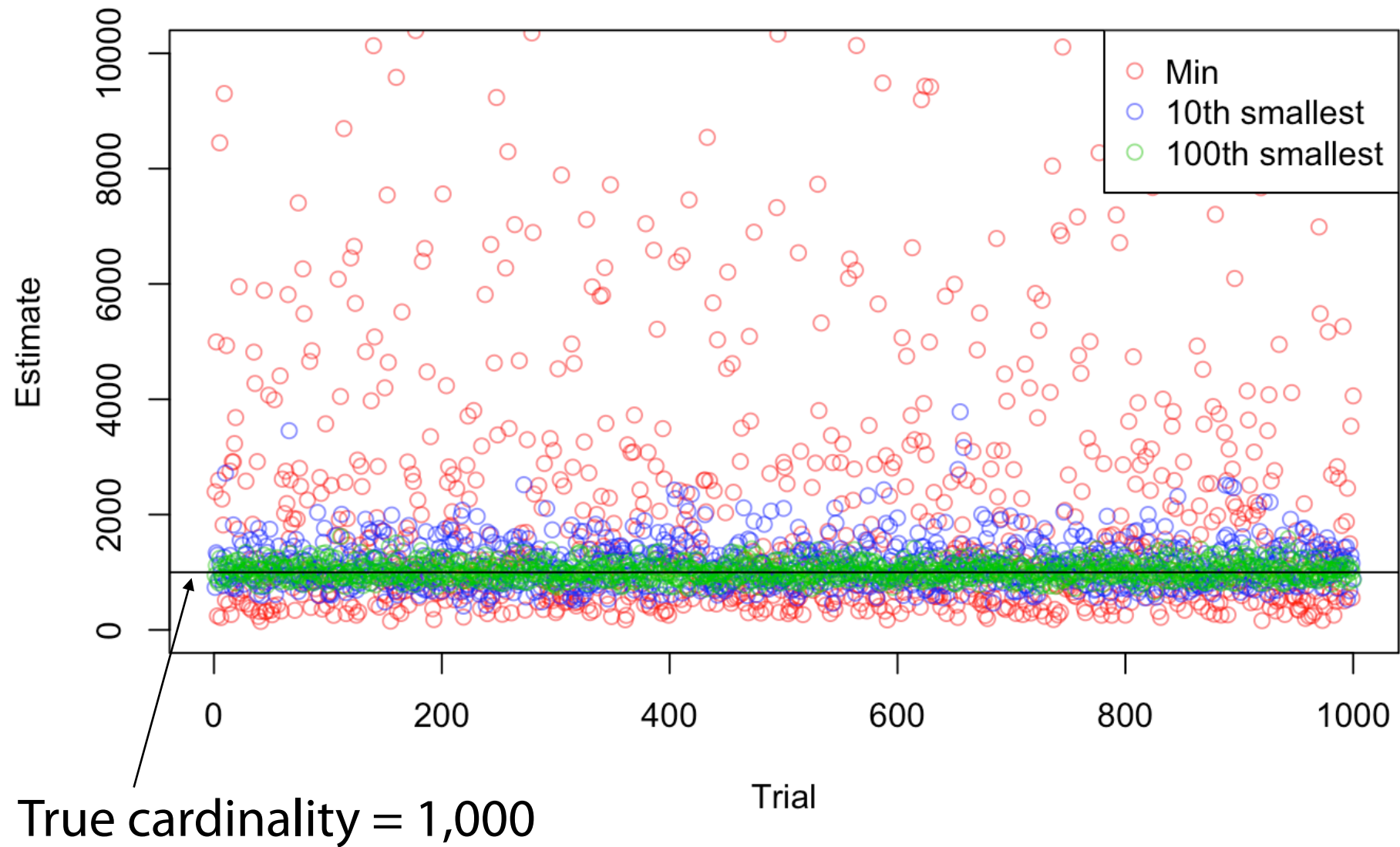
$$= \left[\underbrace{\mathbf{E}[M_1]} + \underbrace{(\mathbf{E}[M_2] - \mathbf{E}[M_1])} + \dots + \underbrace{(\mathbf{E}[M_k] - \mathbf{E}[M_{k-1}])} \right] \cdot \frac{1}{k}$$



k^{th} minimum
value (KMV)

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

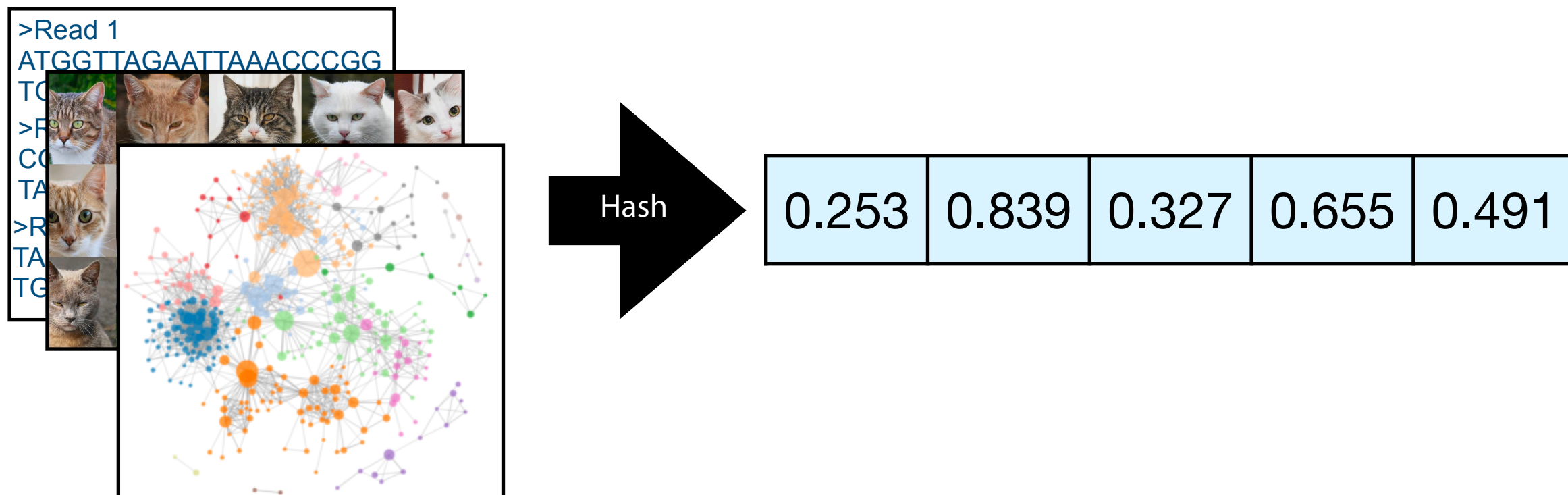
Cardinality Sketch



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?**

