

Algorithms and Data Structures for Data Science

Set Cardinality and Minhash

CS 277

February 27, 2023

Brad Solomon



UNIVERSITY OF
ILLINOIS
URBANA - CHAMPAIGN

Department of Computer Science

Learning Objectives

Introduce the concept of cardinality and cardinality estimation

Demonstrate the relationship between cardinality and similarity

Introduce the MinHash Sketch for set similarity detection

Bloom Filters

A probabilistic data structure storing a set of values

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

Has three key properties:

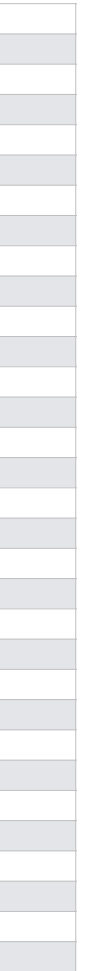
k , number of hash functions

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



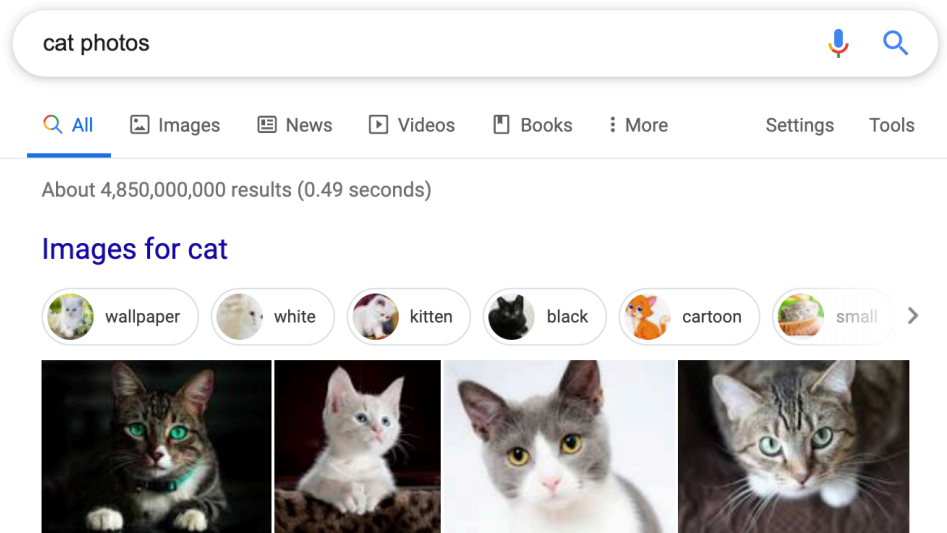
Cardinality

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

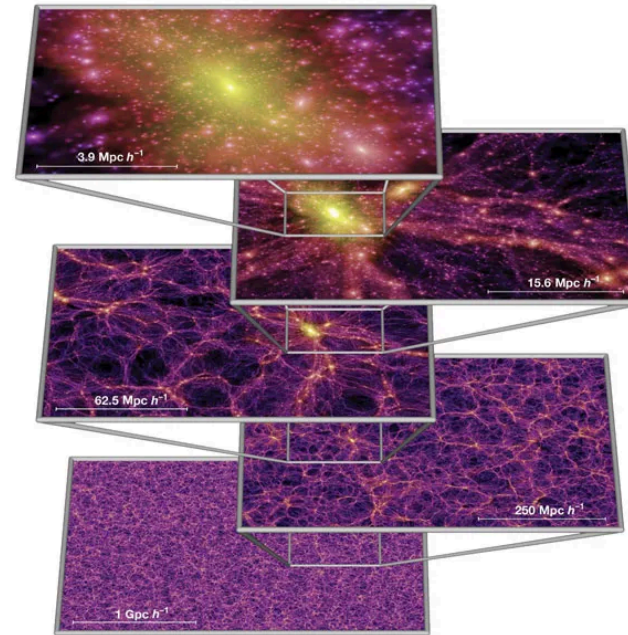


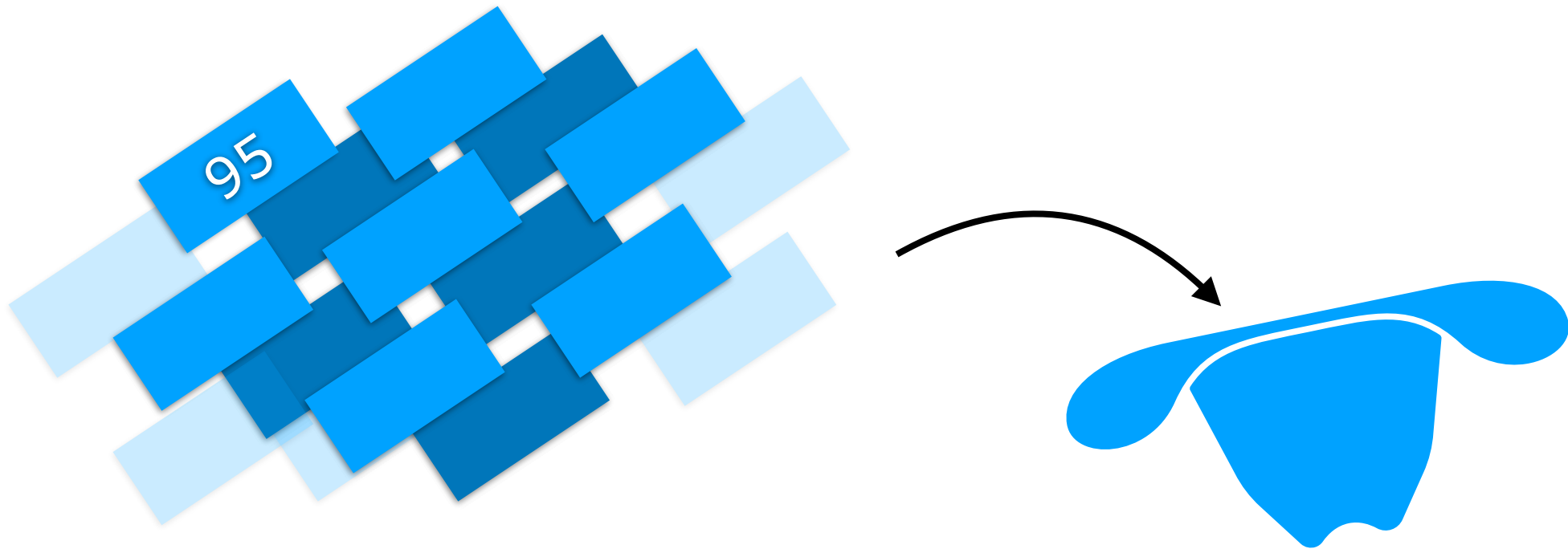
Image: <https://doi.org/10.1038/nature03597>

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.

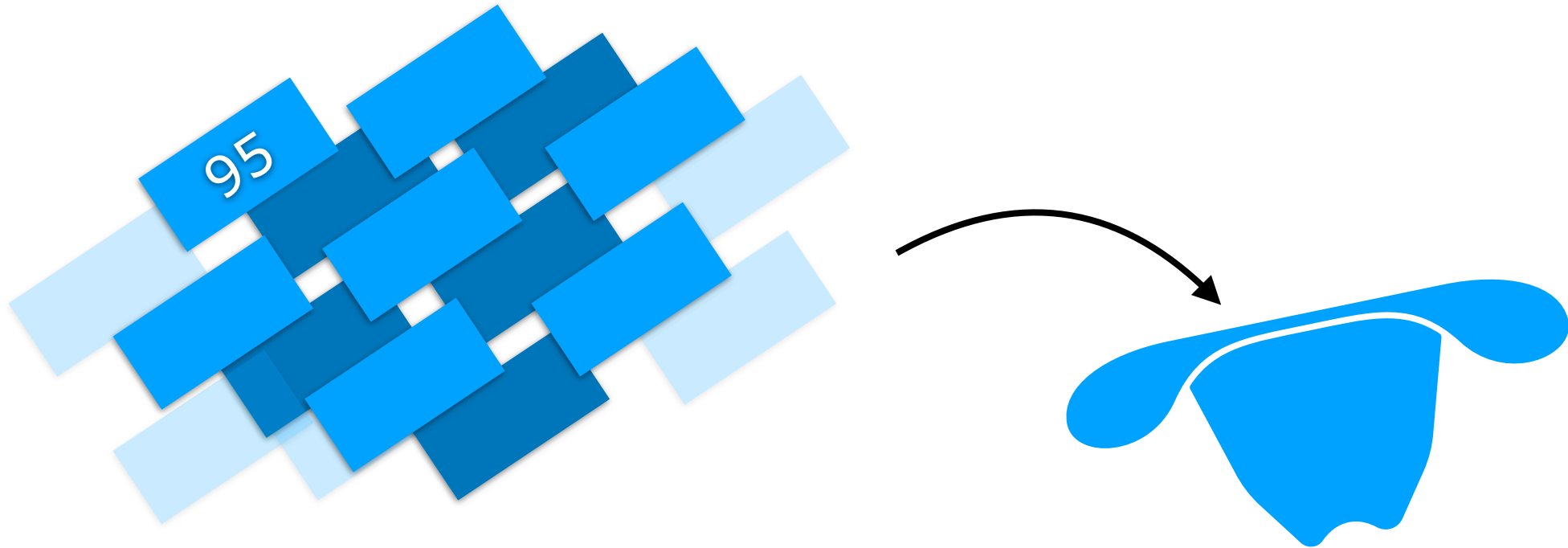
If I told you the value of the card was 95, what have we learned?



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 sets (multiple minimums).



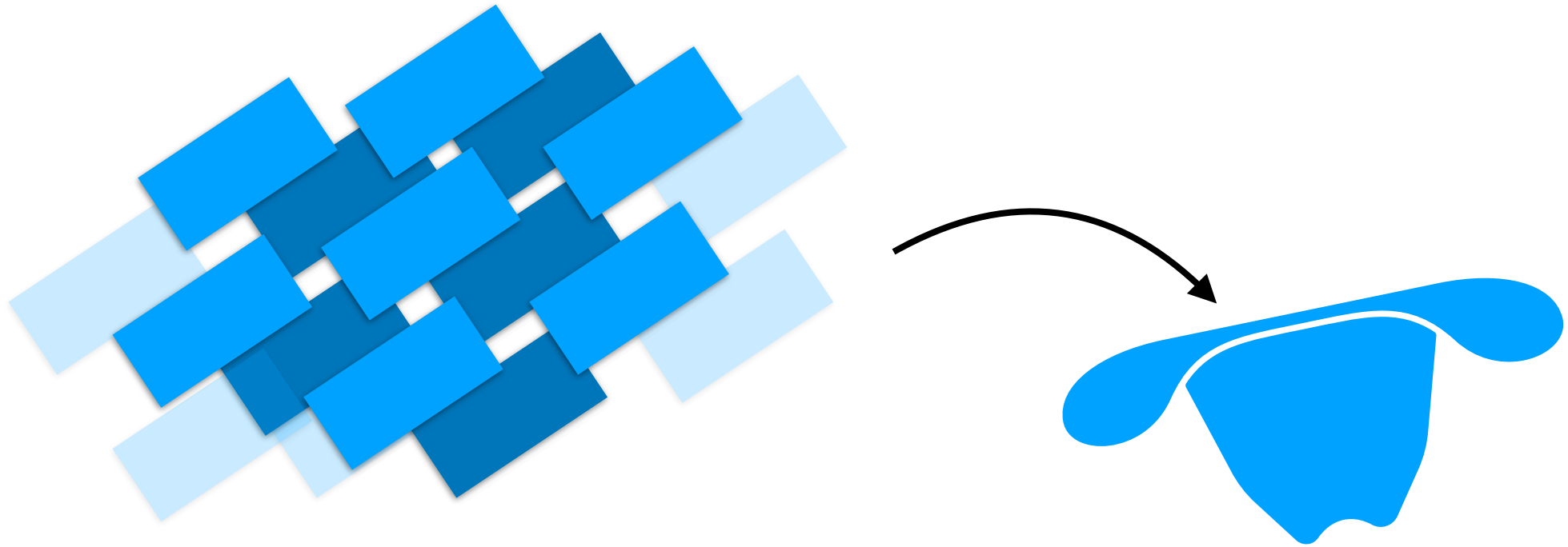
Cardinality Estimation

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



Cardinality Estimation

Why do we care about “the hat problem”?



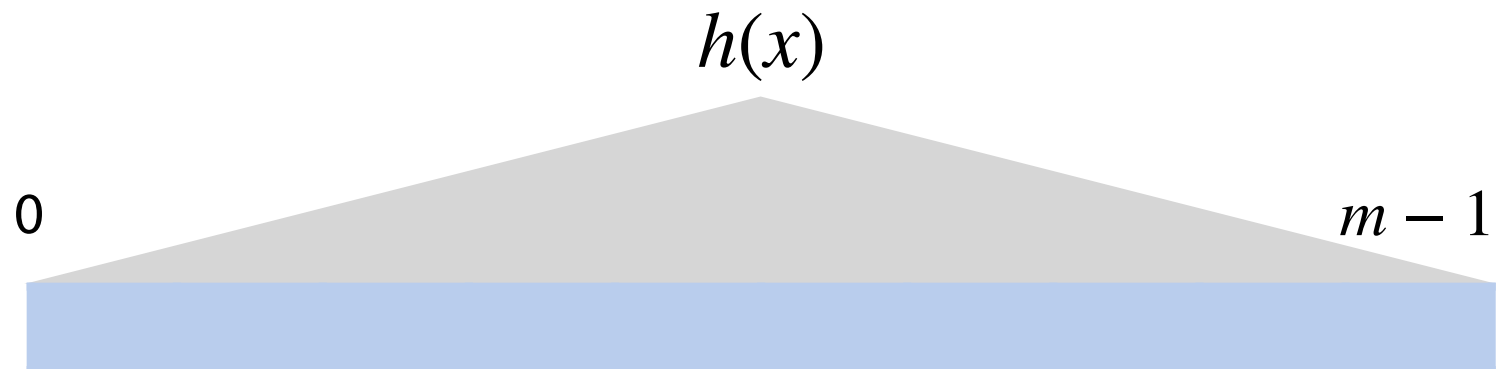


Cardinality Estimation

Now imagine we have a SUHA hash h over a range m .

Here a hash insert is equivalent to adding a card to our hat!

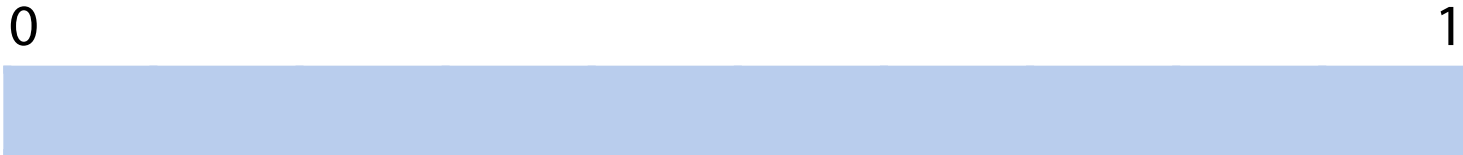
Now storing only the minimum hash value is a **sketch**!



Cardinality Sketch

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

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



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

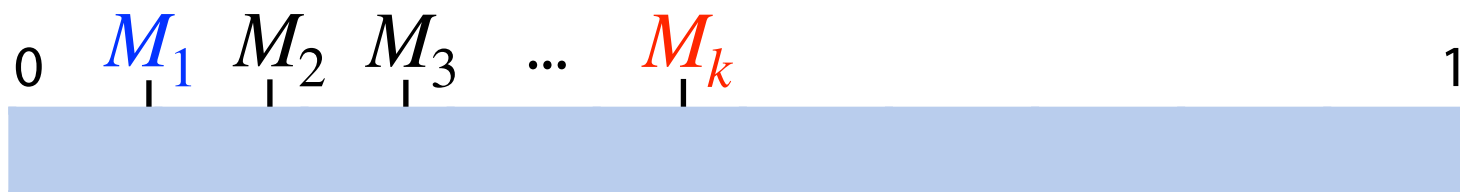
1



Cardinality Sketch

Claim: Taking the k^{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!

$$\begin{aligned} \text{Claim: } \mathbf{E}[M_k] &= \frac{k}{N+1} \\ &= \left[\mathbf{E}[M_1] + (\mathbf{E}[M_2] - \mathbf{E}[M_1]) + \dots + (\mathbf{E}[M_k] - \mathbf{E}[M_{k-1}]) \right] \cdot \frac{1}{k} \end{aligned}$$

M_1
|

M_2
|

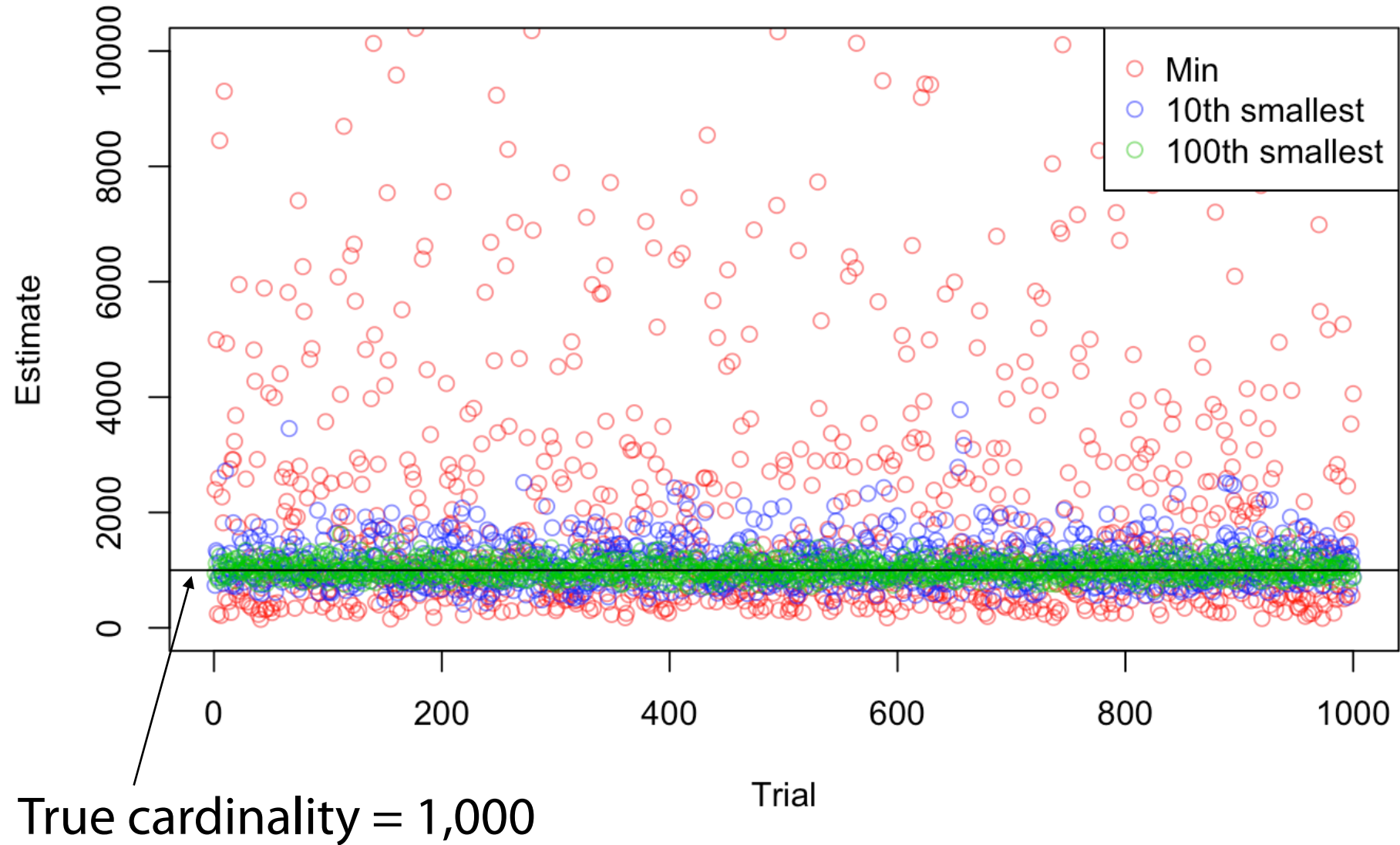
M_3
|

...

M_{k-1}
|

M_k
|

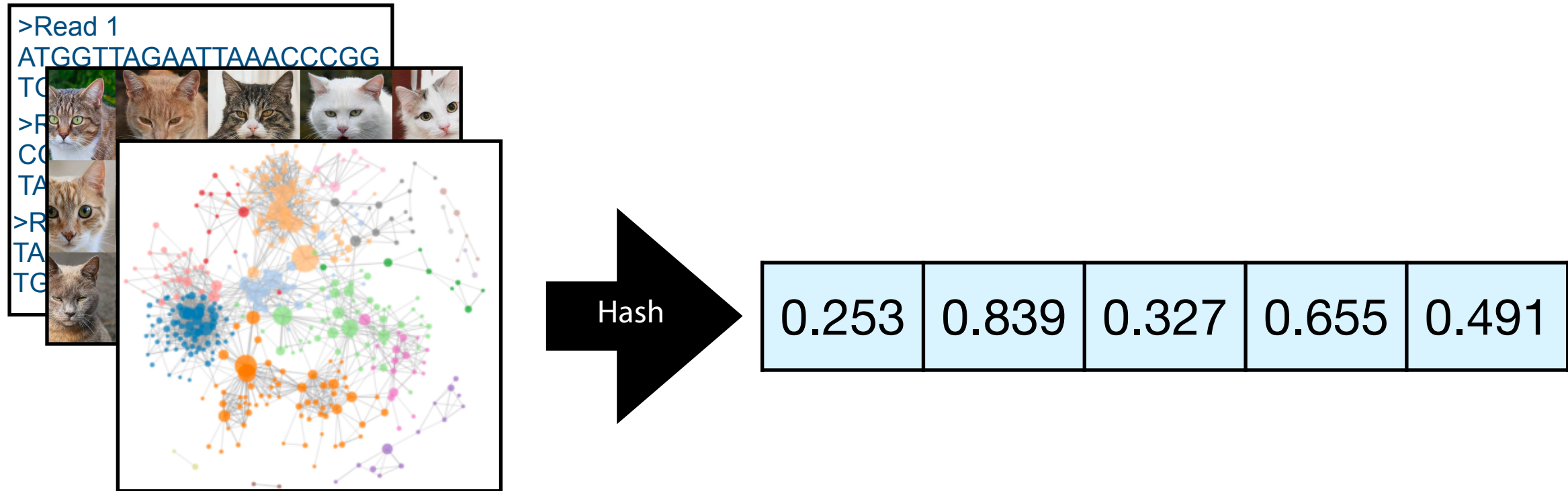
Cardinality



Cardinality



Given any dataset and a SUHA hash function, we can estimate the number of unique items by tracking the minimum hash values.



Applied Cardinalities

Cardinalities

$|A|$

$|B|$

$|A \cup B|$

$|A \cap B|$

Set similarities

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

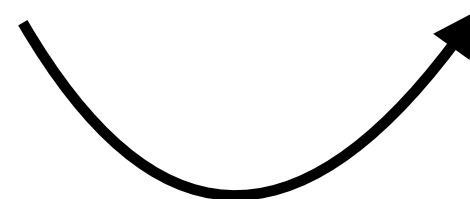
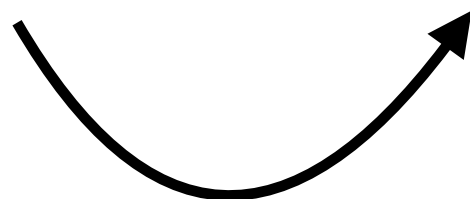
$$J = \frac{|A \cap B|}{|A \cup B|}$$

Real-world
Meaning

AGGCCACAGTGTATTATGACTG
||||| |||||
AGGCCACAGTGAGTTATGACTG

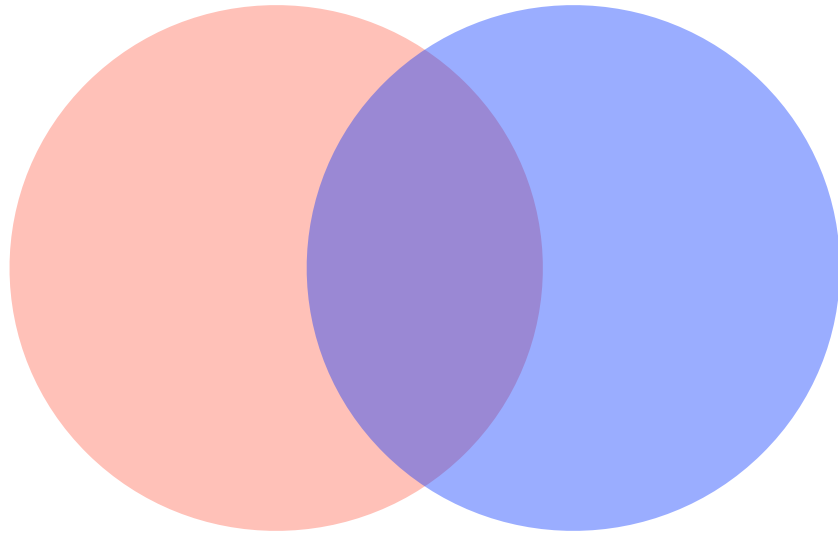
AAAAAAAAAAGATGT-AAGTA
||||| |||||
AAAAAAAAAAGATGTAAAGTA

GAGG--TCAGATTCACAGCCAC
|||| |
GAGGGGTCAGATTCACAGCCAC



Set Similarity

To measure **similarity** of A & B , we need both a measure of how similar the sets are but also the total size of both sets.

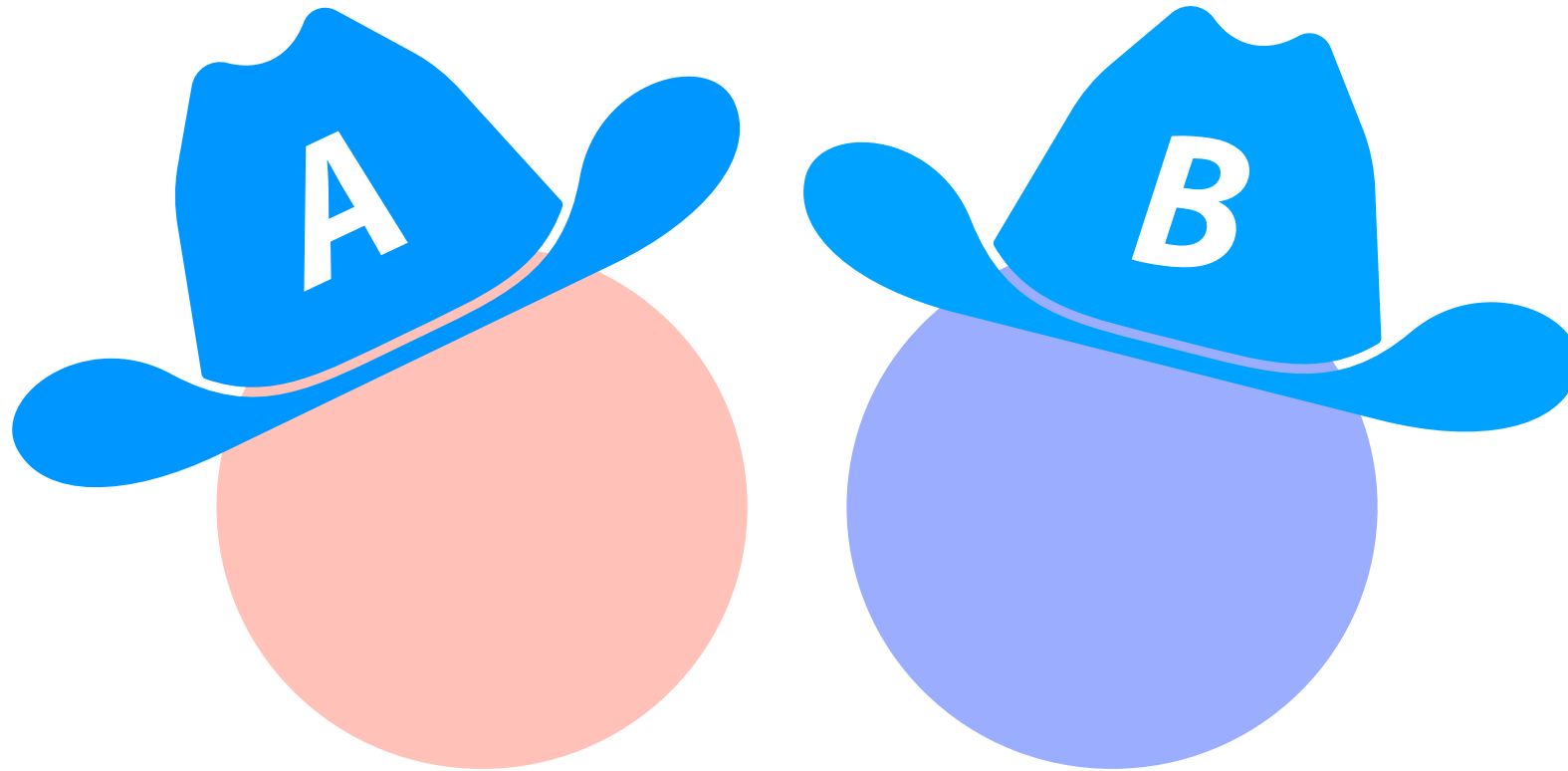


$$J = \frac{|A \cap B|}{|A \cup B|}$$

J is the **Jaccard coefficient**

Similarity Sketches

But what do we do when we only have a sketch?



Similarity Sketches

Imagine we have two datasets represented by their k th minimum values

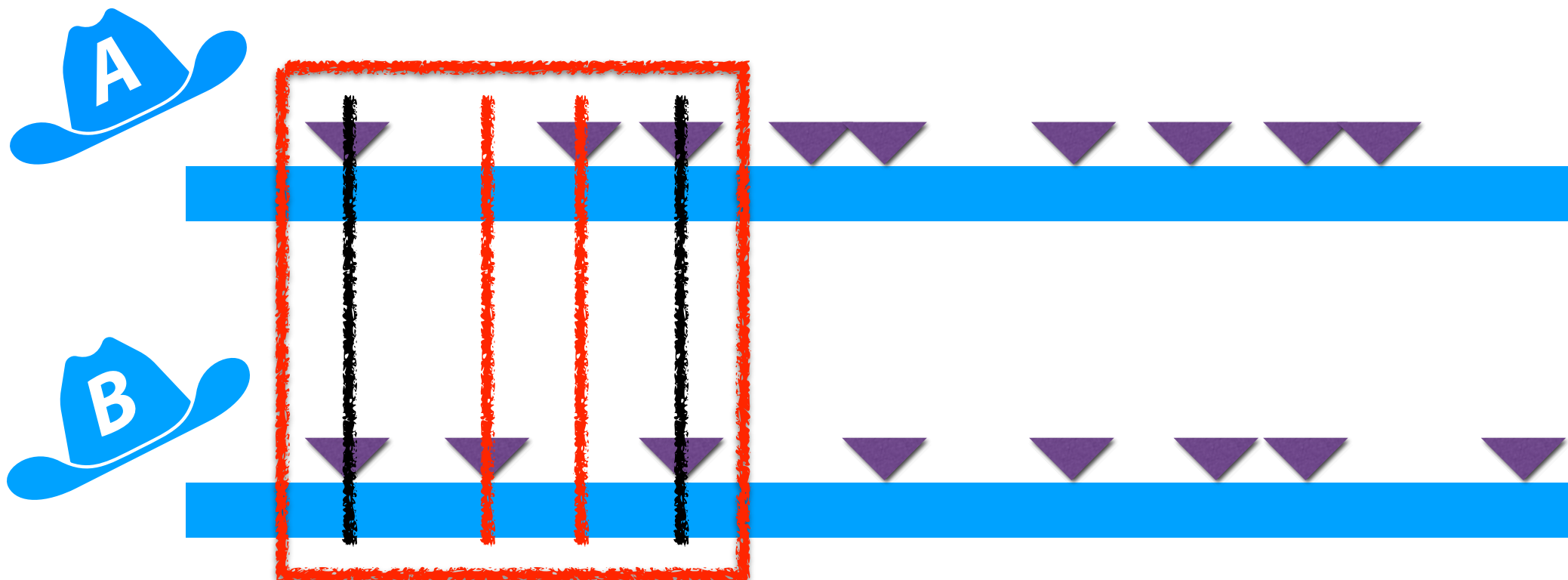


Image inspired by: Ondov B, Starrett G, Sappington A, Kostic A, Koren S, Buck CB, Phillippy AM. **Mash Screen: high-throughput sequence containment estimation for genome discovery.** *Genome Biol* 20, 232 (2019)

Similarity Sketches

Claim: Under SUHA, set similarity can be estimated by sketch similarity!

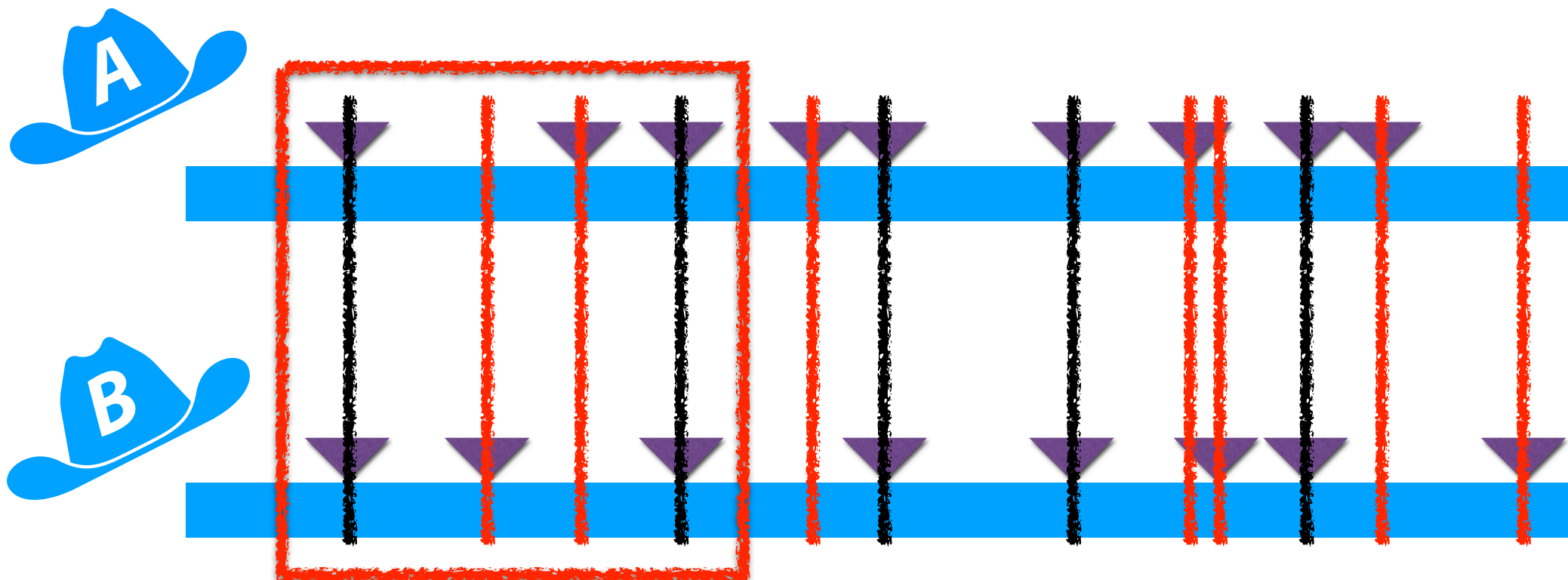


Image inspired by: Ondov B, Starrett G, Sappington A, Kostic A, Koren S, Buck CB, Phillippy AM. **Mash Screen: high-throughput sequence containment estimation for genome discovery.** *Genome Biol* 20, 232 (2019)

The MinHash Sketch

Let sets A and B be two arbitrary sets of at least 8 elements

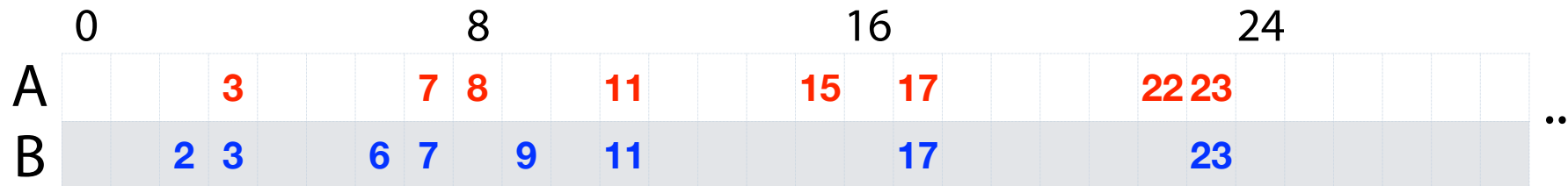
The eight minimum hash values for sets A and B is a **MinHash Sketch**

Sketch A

3	15
7	17
8	22
11	23

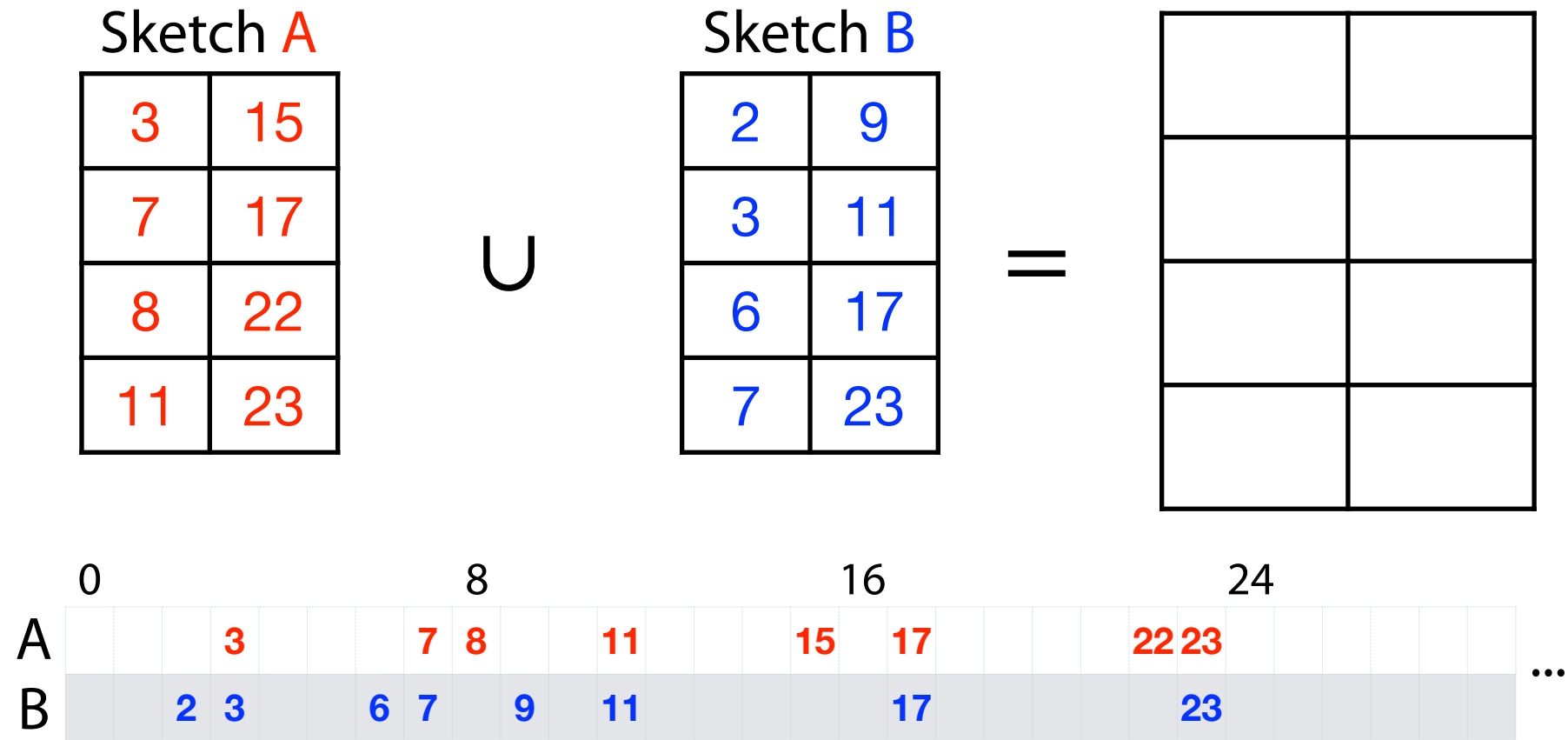
Sketch B

2	9
3	11
6	17
7	23



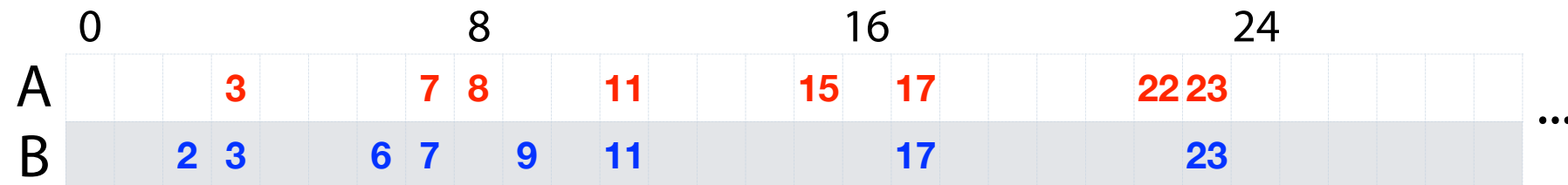
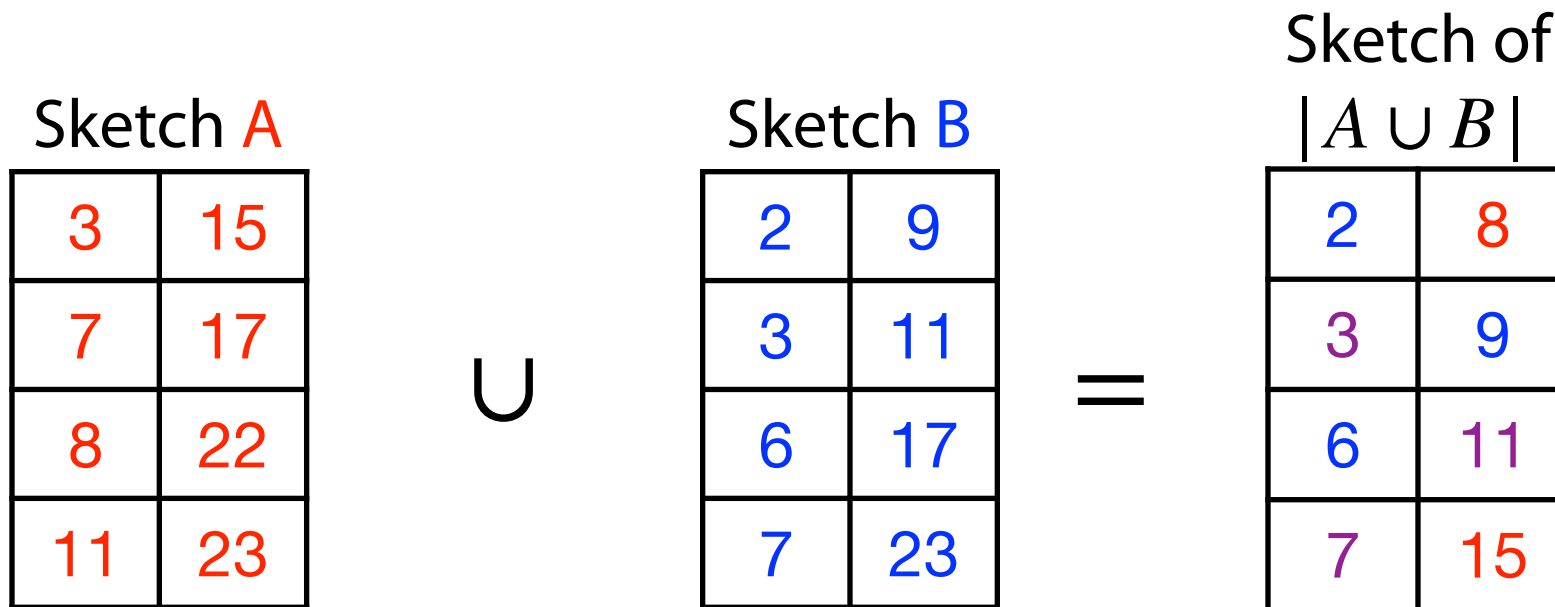
The MinHash Sketch

To get similarity, we want to estimate $|A \cup B|$ and $|A \cap B| \dots$



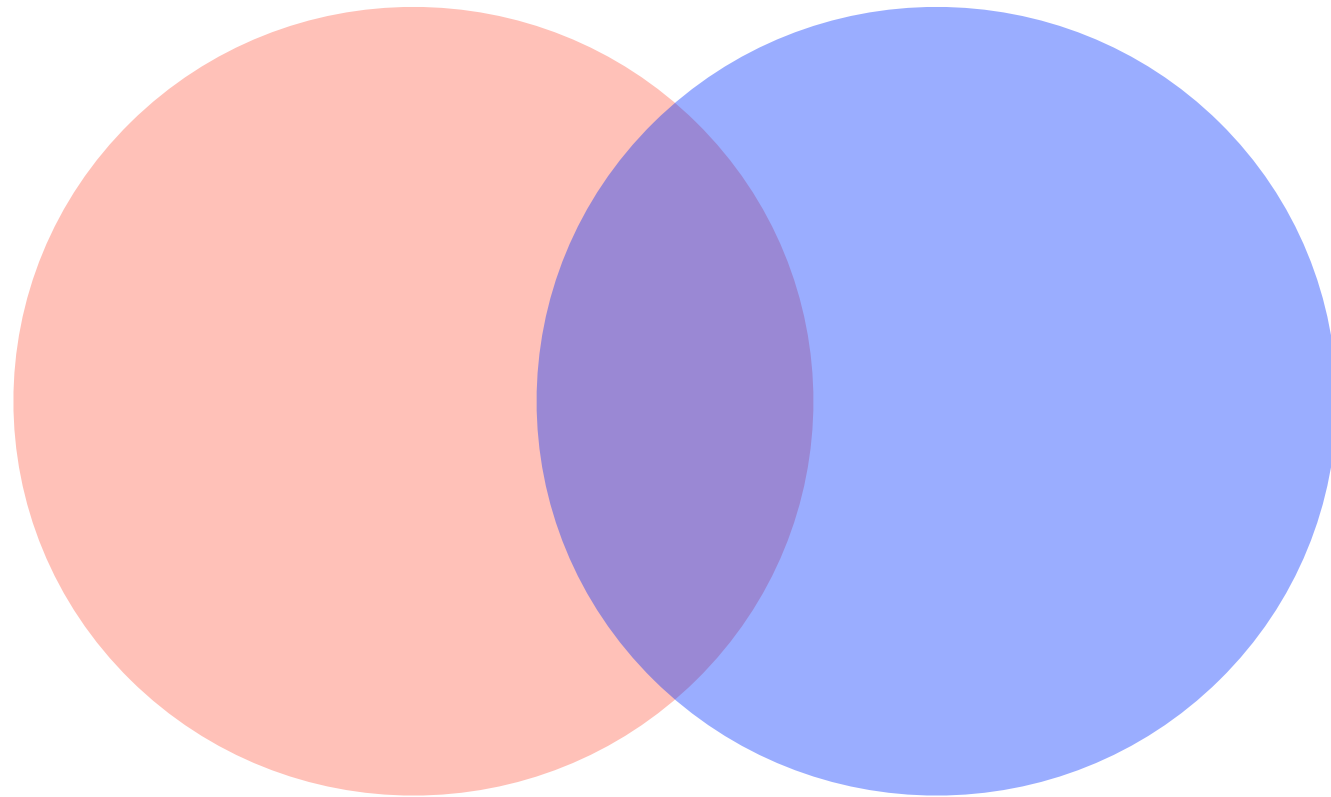
The MinHash Sketch

To get similarity, we want to estimate $|A \cup B|$ and $|A \cap B| \dots$



Inclusion-Exclusion Principle

$$|A \cap B| =$$



Set Similarity

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

$$= \frac{|A \cap B|}{|A \cap B| + |A \triangle B|}$$

$$= \frac{|A \cap B|}{|A| + |B| - |A \cap B|}$$

$$= \frac{|A| + |B| - |A \cup B|}{|A \cup B|}$$

The MinHash Sketch

Using **inclusion-exclusion** principle and KMV, we can estimate similarity!

Sketch A

3	15
7	17
8	22
11	23

Sketch B

2	9
3	11
6	17
7	23

Sketch of
 $|A \cup B|$

2	8
3	9
6	11
7	15

k th minimum value (KMV) with $k = 8$,
assuming hash range is integers in $[0, 100)$:

$$\begin{aligned}
 &= \frac{800/23 - 1 + 800/23 - 1 - 800/15 - 1}{800/15 - 1} \\
 &= \frac{34.782 + 34.782 - 53.333 - 1}{53.333 - 1} \\
 &\approx 0.29
 \end{aligned}$$

$$\frac{|A| + |B| - |A \cup B|}{|A \cup B|}$$

The MinHash Sketch

Claim: Cardinality of the intersection can also be estimated directly!

Sketch A

3	15
7	17
8	22
11	23

Sketch B

2	9
3	11
6	17
7	23

Sketch of
 $|A \cup B|$

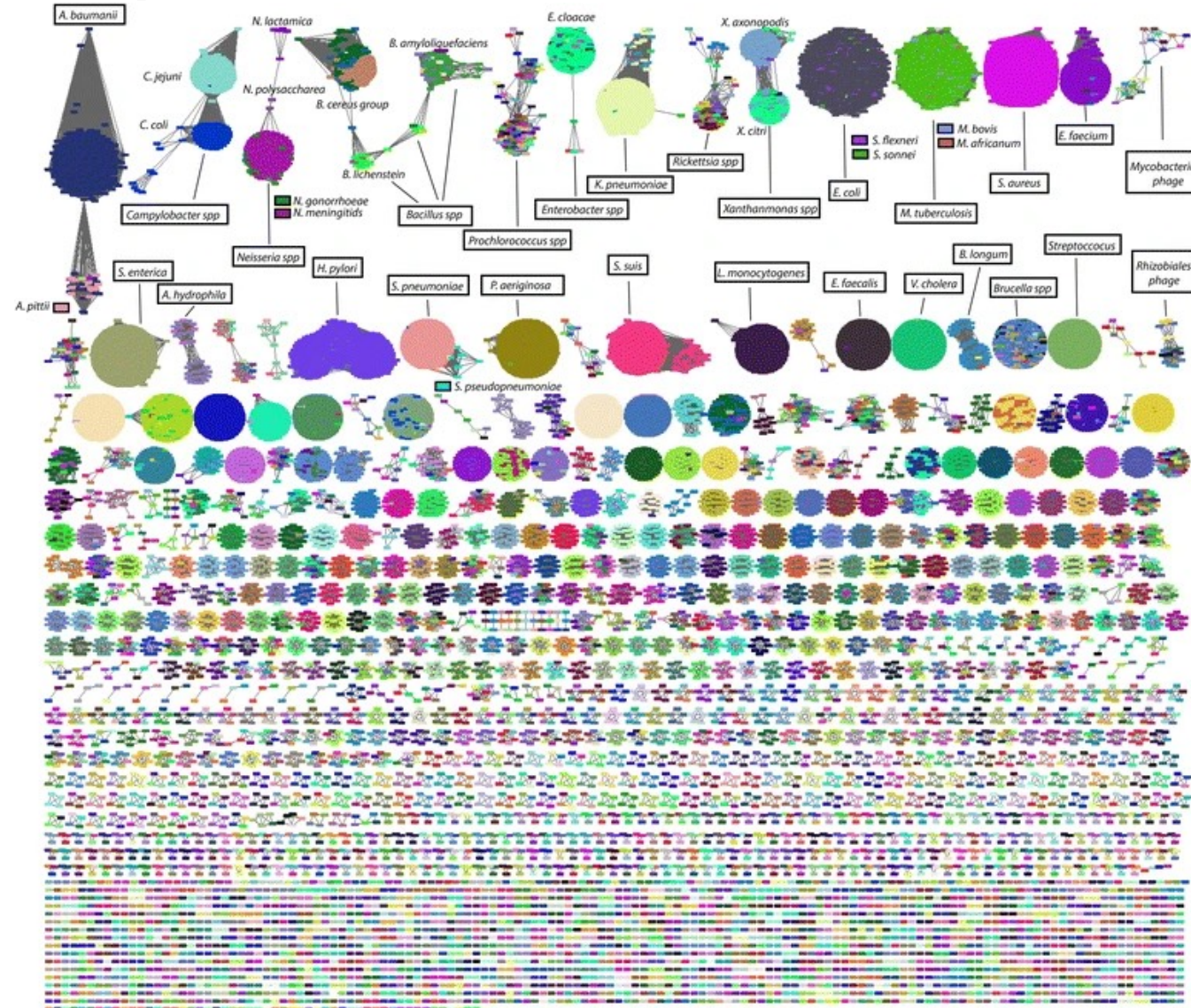
2	8
3	9
6	11
7	15

1) Sequence decomposed into **kmers**

S_1 : CATGGACCGACCAG
CAT GAC GAC
ATG ACC ACC
TGG CCG CCA
GGA CGA CAG

GCAGTACCGATCGT : S_2
GTA CGA CGT
AGT CCG TCG
CAG ACC ATC
GCA TAC GAT

MinHash in practice



Mash: fast genome and metagenome distance estimation using MinHash

Ondov et al (2016) *Genome Biology*

Reviewing probabilistic data sketches



What sketch would I use for the following:

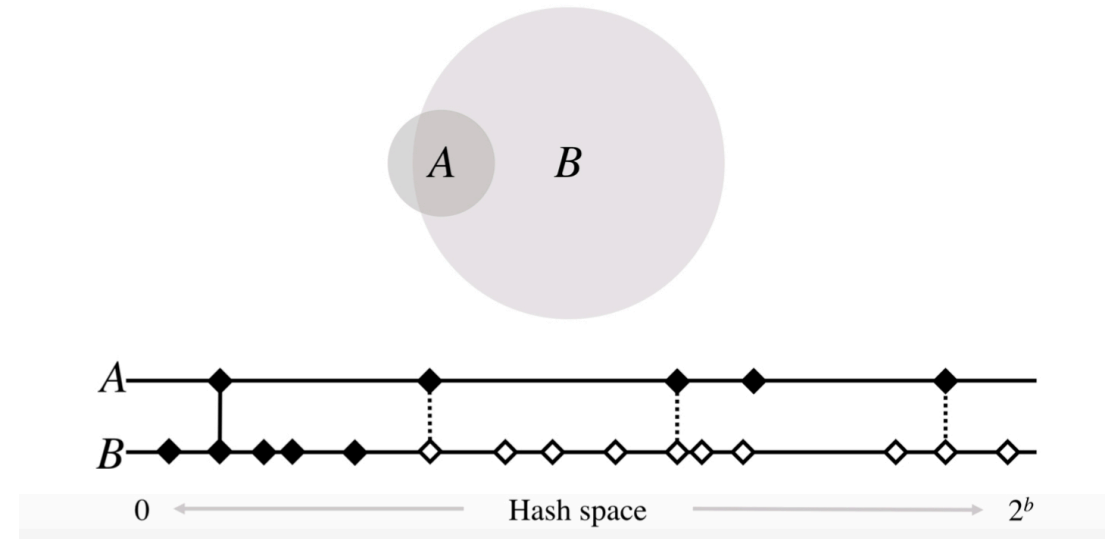
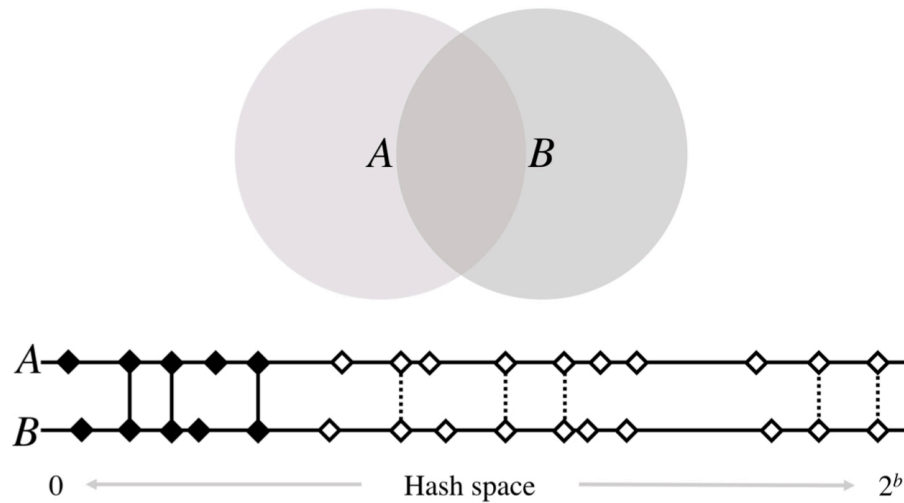
Does a *specific* object exist in my data?

How many unique objects do I have in my set?

How similar are two datasets?

Bonus Slides (Taking it one step further...)

Bottom-k minhash has low accuracy if the cardinality of sets are skewed



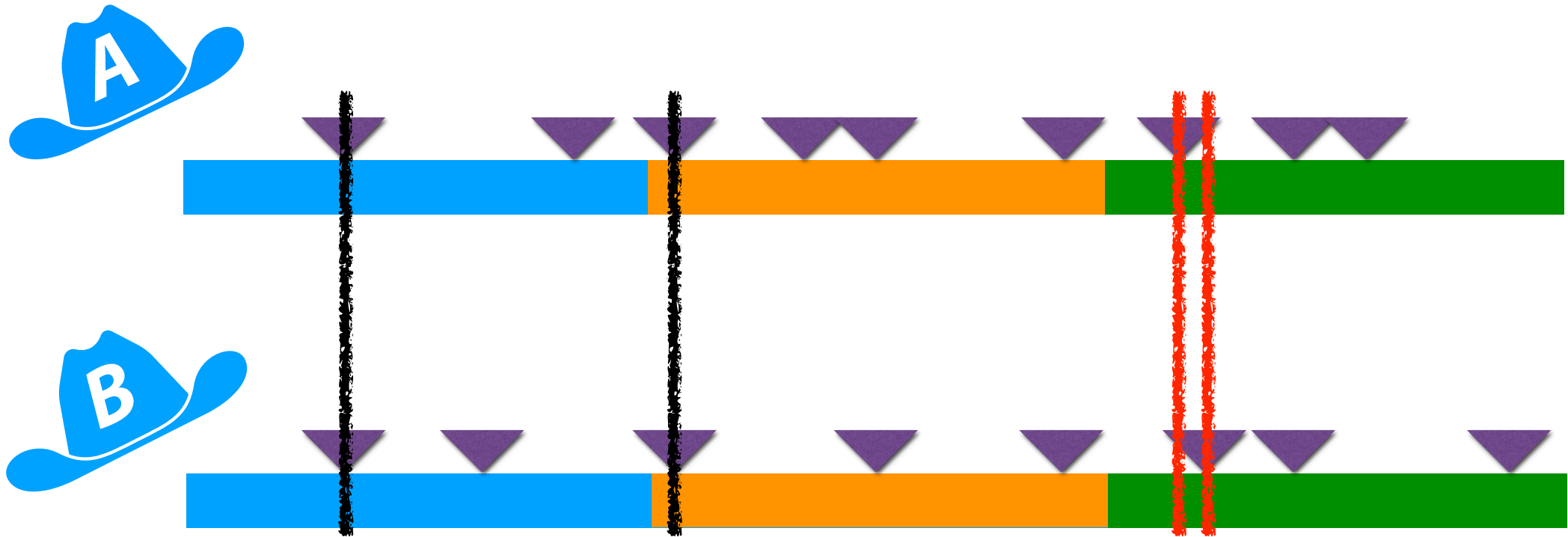
K-Hash Minhash

What if instead we used k different hashes and took the min each time?

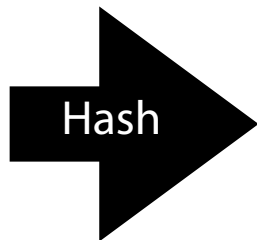
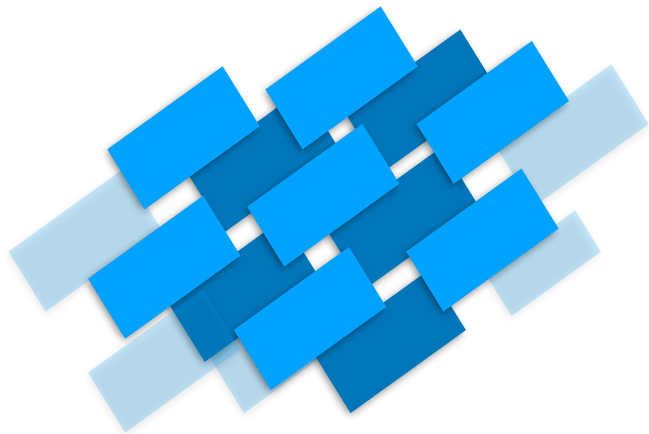


K-Partition Minhash

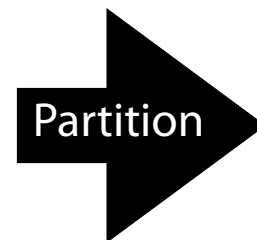
What if we instead took the minimum of k-partitions?



K-Partition Minhash



1010110101
0001111010
1101101011
1011010110
0101100000
0010001101



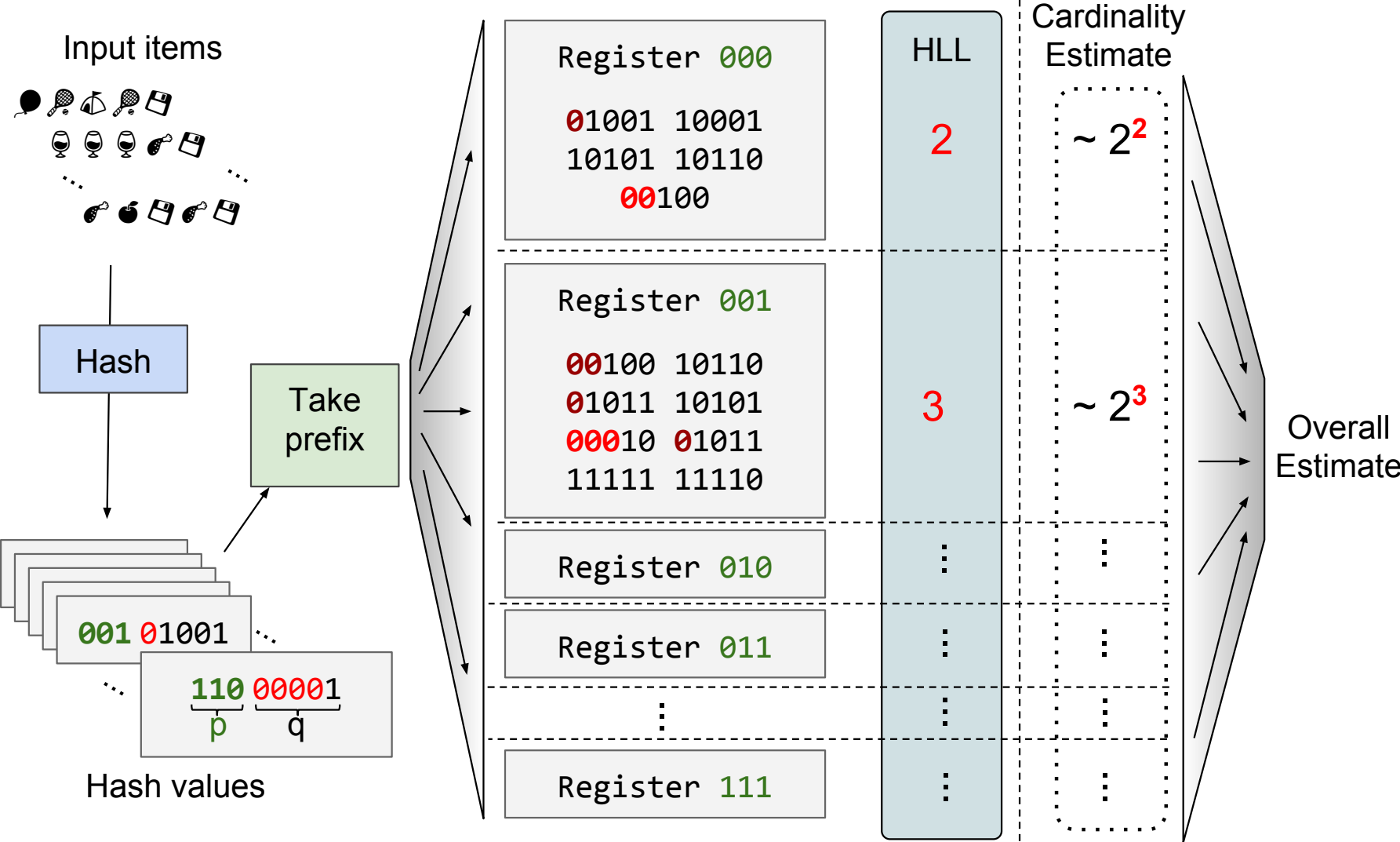
00
01111010
10001101

01
01100000

10
10110101
11010110

11
01101011

HyperLogLog



Where do we go from here?

Hashing (and hash adjacent data structures) seem optimal!

Is there something they can't do that we might want?

The Sorting Problem

Given a collection of objects, C , with comparable values, order the objects such that $\forall x \in C, x_i \leq x_{i+1}$

Input:

8	4	3	1	2	5	6	9	0	7
---	---	---	---	---	---	---	---	---	---

Output:

0	1	2	3	4	5	6	7	8	9
---	---	---	---	---	---	---	---	---	---

SelectionSort

4	3	6	7	1
---	---	---	---	---

--	--	--	--	--

--	--	--	--	--

--	--	--	--	--

--	--	--	--	--

1. Find the i -th smallest value
2. Place it at position i via swap
3. Repeat for $0 \leq i \leq n - 1$

SelectionSort Efficiency

(large n)