# Data Structures and Algorithms
# Review and Return to Cardinality

CS 225

December 2, 2024

Brad Solomon

UNIVERSITY OF ILLINOIS
URBANA-CHAMPAIGN

Department of Computer Science

# Course Announcements

This week's lab is **optional.** Will be worth the equivalent value in EC

Part 2 of **External Research Survey** releases tomorrow! Worth 2 EC

Reminder: Exam 5 is this week!

Reminder: Final exam starts as early as Thursday December 12th

**Please fill out ICES evaluations!**
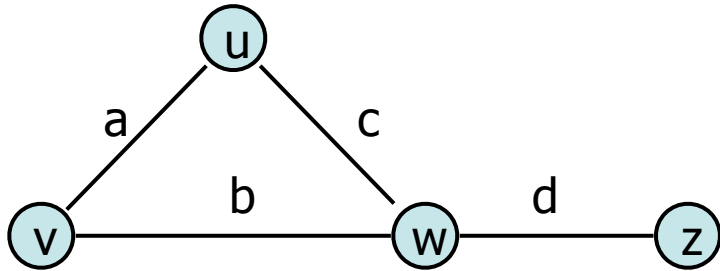
# Learning Objectives

A brief review of exam 5 content

Review high level motivation behind sketching data structure

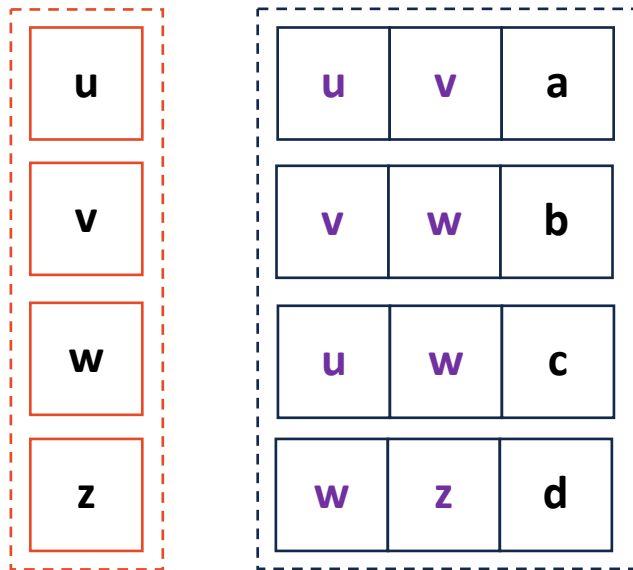Introduce the concept of cardinality and cardinality estimation

# Graph Implementation: Edge List  $|V| = n, |E| = m$

*The equivalent of an 'unordered' data structure*

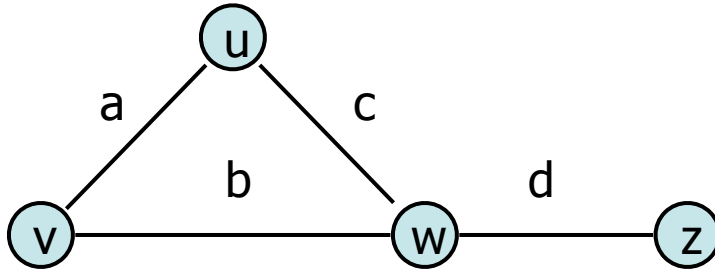**Vertex Storage:**

An optional list of vertices

**Edge Storage:**

A list storing edges as (V1, V2, Weight)

**Most graphs are stored as just an edge list!**

# Graph Implementation: Adjacency Matrix

$|V| = n, |E| = m$



**Vertex Storage:**

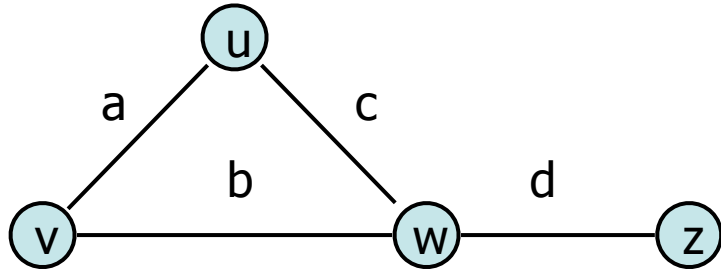A hash table of vertices

Implicitly or explicitly store index

| | u | 0 |
|---|---|---|
| | v | 1 |
| | w | 2 |
| | z | 3 |

| | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| **0** | - | a | c | 0 |
| **1** | | - | b | 0 |
| **2** | | | - | d |
| **3** | | | | - |

**Edge Storage:**

A |V| x |V| matrix of edges

Weight is stored at position (u, v)
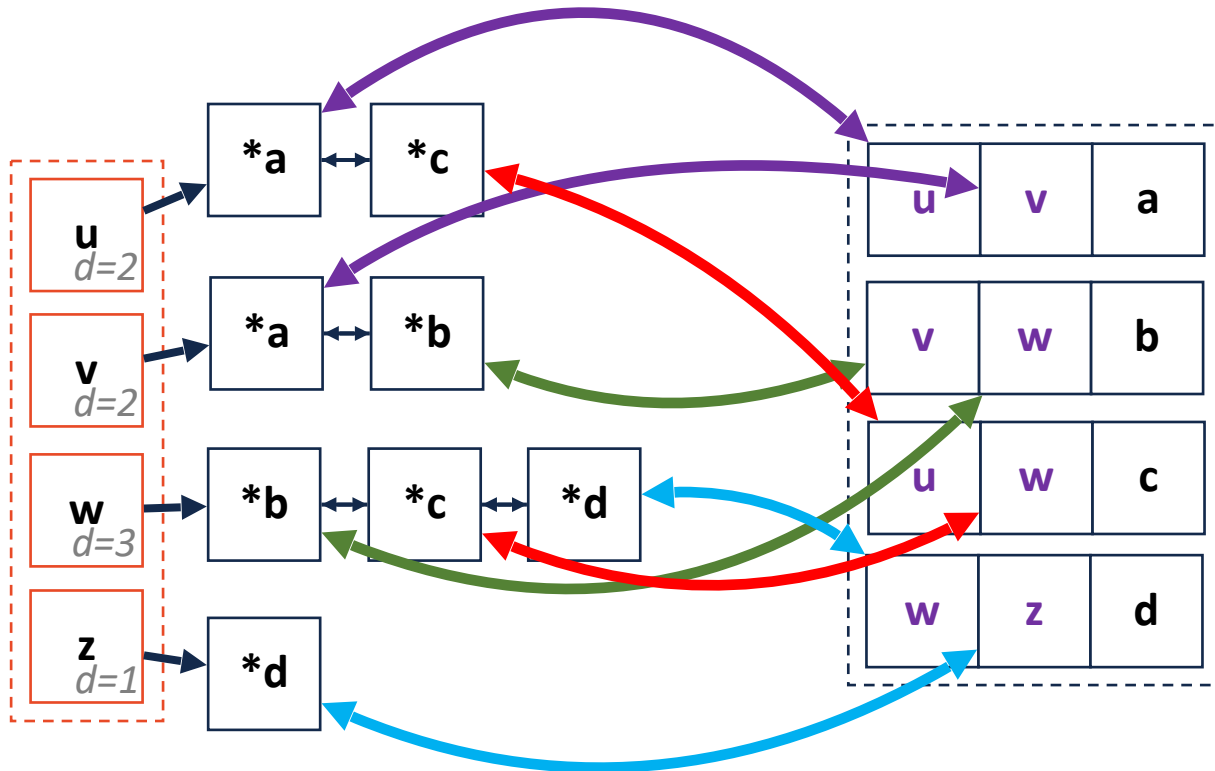
# Adjacency List



**Vertex Storage:**

A bidirectional linked list with size variable

Each node is a pointer to edge in edge list

**Edge Storage:**

A list of (v1, v2, weight) edges

Also store pointers back to nodes

$$|V| = n, |E| = m$$

| Expressed as O(f) | Edge List | Adjacency Matrix | Adjacency List |
|---|---|---|---|
| Space | n+m | n² | n+m |
| insertVertex(v) | 1* | n* | 1* |
| removeVertex(v) | n+m | n | deg(v) |
| insertEdge(u, v) | 1 | 1 | 1* |
| removeEdge(u, v) | m | 1 | min( deg(u), deg(v) ) |
| incidentEdges(v) | m | n | deg(v) |
| areAdjacent(u, v) | m | 1 | min( deg(u), deg(v) ) |

# Summary: DFS and BFS     $|V| = n, |E| = m$

Both are **O(n+m)** traversals! They label every edge and every node

**BFS**

Solves unweighted MST

Solves shortest path

Solves cycle detection

Memory bounded by width

**DFS**

Solves unweighted MST

Solves cycle detection

Memory bounded by longest path

# Kruskal's Algorithm

1) Build a **priority queue** on edges

```
1  KruskalMST(G):
2    DisjointSets forest
3    foreach (Vertex v : G.vertices()):
4      forest.makeSet(v)
5
6    PriorityQueue Q    // min edge weight
7    Q.buildFromGraph(G.edges())
8
9    Graph T = (V, {})
10
11   while |T.edges()| < n-1:
12     Vertex (u, v) = Q.removeMin()
13     if forest.find(u) != forest.find(v):
14       T.addEdge(u, v)
15       forest.union( forest.find(u),
16                     forest.find(v) )
17
18   return T
19
```
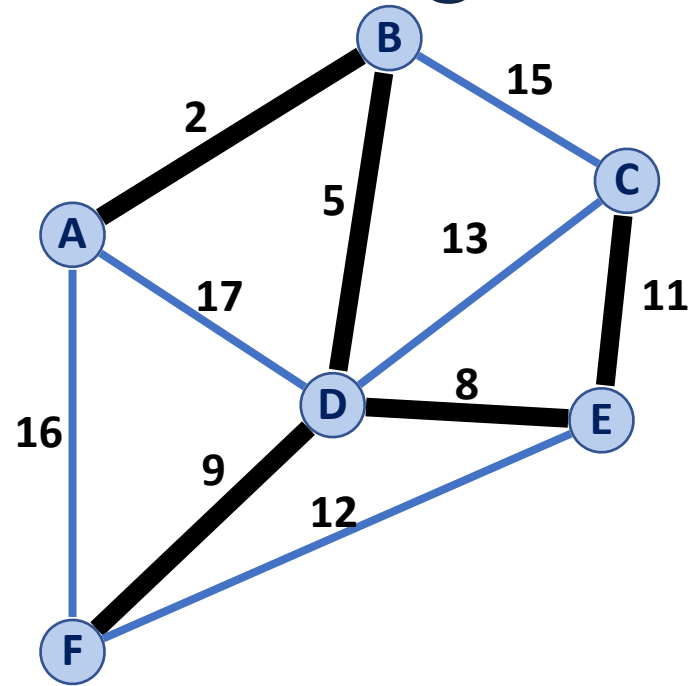
2) Build a **disjoint set** on vertices

3) Repeatedly find min edge

    If edge connects two sets

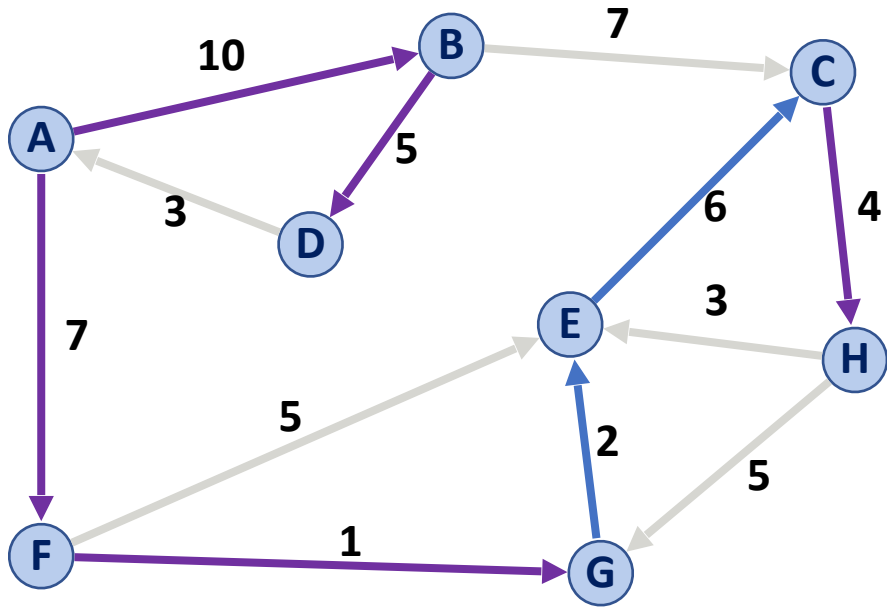    Union and record edge

4) Stop after n-1 edges recorded

# Prim's Algorithm



| A | B | C | D | E | F |
|---|---|---|---|---|---|
| 0, — | 2, A | 11, E | 5, B | 8, D | 9, D |

```
1   PrimMST(G, s):
2      Input: G, Graph;
3             s, vertex in G, starting vertex
4      Output: T, a minimum spanning tree (MST) of G
5
6      foreach (Vertex v : G.vertices()):
7         d[v] = +inf
8         p[v] = NULL
9      d[s] = 0
10
11     PriorityQueue Q    // min distance, defined by d[v]
12     Q.buildHeap(G.vertices())
13     Graph T            // "labeled set"
14
15     repeat n times:
16        Vertex m = Q.removeMin()
17        T.add(m)
18        foreach (Vertex v : neighbors of m not in T):
19           if cost(v, m) < d[v]:
20              d[v] = cost(v, m)
21              p[v] = m
22
23     return T
```

# Dijkstra's Algorithm (SSSP)



```
DijkstraSSSP(G, s):
6    foreach (Vertex v : G.vertices()):
7        d[v] = +inf
8        p[v] = NULL
9    d[s] = 0
10
11   PriorityQueue Q // min distance, defined by d[v]
12   Q.buildHeap(G.vertices())
13   Graph T           // "labeled set"
14
15   repeat n times:
16       Vertex u = Q.removeMin()
17       T.add(u)
18       foreach (Vertex v : neighbors of u not in T):
19           if cost(u, v) + d[u] < d[v]:
20               d[v] = cost(u, v) + d[u]
21               p[v] = u
```

| A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|
| -- | A | E | B | G | A | F | C |
| 0 | 10 | 16 | 15 | 10 | 7 | 8 | 20 |

# Floyd-Warshall Algorithm

Floyd-Warshall's Algorithm is an alternative to Dijkstra in the presence of negative-weight edges (not negative weight cycles).

```
1   FloydWarshall(G):
2     Let d be a adj. matrix initialized to +inf
3     foreach (Vertex v : G):
4       d[v][v] = 0
5     foreach (Edge (u, v) : G):
6       d[u][v] = cost(u, v)
7
8     foreach (Vertex u : G):
9       foreach (Vertex v : G):
10        foreach (Vertex w : G):
11          if (d[u, v] > d[u, w] + d[w, v])
12            d[u, v] = d[u, w] + d[w, v]
```

# A Hash Table based Dictionary

**User Code (is a map):**

```
1 Dictionary<KeyType, ValueType> d;
2 d[k] = v;
```

A **Hash Table** consists of three things:

1. A hash function

2. A data storage structure
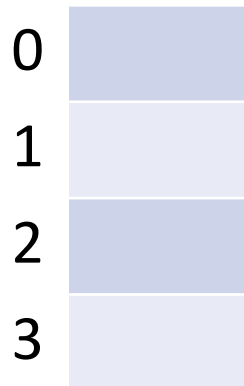
3. **A method of addressing** *hash collisions*

# Open vs Closed Hashing

Addressing hash collisions depends on your storage structure.

- **Open Hashing:** store *k,v* pairs externally



- **Closed Hashing:** store *k,v* pairs in the hash table

# Separate Chaining Under SUHA

**Claim:** Under SUHA, expected length of chain is $\dfrac{n}{m}$

**Table Size:** $m$

**Num objects:** $n$

$\alpha_j =$ expected # of items hashing to position j

$$\alpha_j = \sum_i H_{i,j}$$

$$H_{i,j} = \begin{cases} 1 \text{ if item i hashes to j} \\ 0 \text{ otherwise} \end{cases}$$

$$E[\alpha_j] = E\left[\sum_i H_{i,j}\right]$$

$$Pr[H_{i,j} = 1] = \frac{1}{m}$$

$$E[\alpha_j] = n * Pr(H_{i,j} = 1)$$

$$\mathbf{E}[\alpha_\mathbf{j}] = \frac{\mathbf{n}}{\mathbf{m}}$$

# Separate Chaining Under SUHA

**Under SUHA, a hash table of size $m$ and $n$ elements:**

Find runs in: $O(1 + \alpha)$

Insert runs in: $O(1)$

Remove runs in: $O(1 + \alpha)$

# Running Times *(Don't memorize these equations, no need.)*

*The expected number of probes for find(key) under SUHA*

**Linear Probing:**
- Successful: $\frac{1}{2}(1 + 1/(1-\alpha))$
- Unsuccessful: $\frac{1}{2}(1 + 1/(1-\alpha))^2$

**Double Hashing:**
- Successful: $1/\alpha * \ln(1/(1-\alpha))$
- Unsuccessful: $1/(1-\alpha)$

**Separate Chaining:**
- Successful: $1 + \alpha/2$
- Unsuccessful: $1 + \alpha$

**Instead, observe:**

**- As $\alpha$ increases:**

Runtime approaches infinity!

**- If $\alpha$ is constant:**

Runtime is a constant!

# Resizing a hash table

When and how do you resize?

# Any (review) questions?

# Memory-Constrained Data Structures

What method would you use to build a search index on a collection of objects *in a memory-constrained environment*?
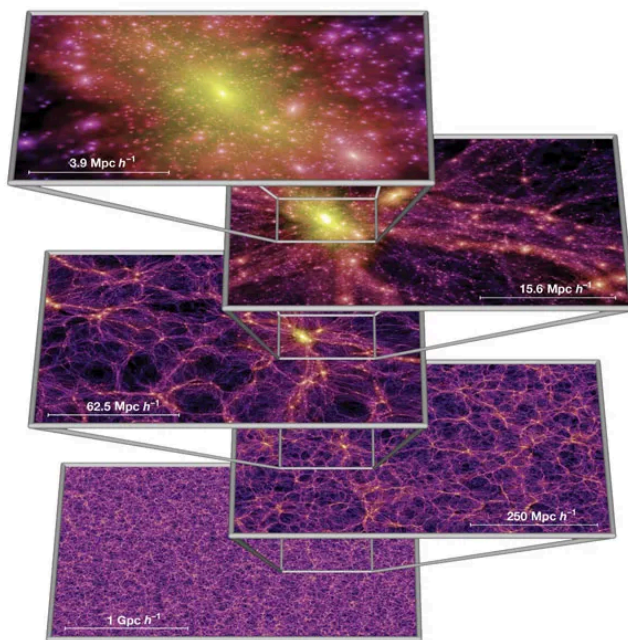
## Constrained by Big Data (Large $N$)



| Sky Survey Projects | Data Volume |
| --- | --- |
| DPOSS (The Palomar Digital Sky Survey) | 3 TB |
| 2MASS (The Two Micron All-Sky Survey) | 10 TB |
| GBT (Green Bank Telescope) | 20 PB |
| GALEX (The Galaxy Evolution Explorer ) | 30 TB |
| SDSS (The Sloan Digital Sky Survey) | 40 TB |
| SkyMapper Southern Sky Survey | 500 TB |
| PanSTARRS (The Panoramic Survey Telescope and Rapid Response System) | ~ 40 PB expected |
| LSST (The Large Synoptic Survey Telescope) | ~ 200 PB expected |
| SKA (The Square Kilometer Array) | ~ 4.6 EB expected |

Table: http://doi.org/10.5334/dsj-2015-011

Estimated total volume of one array: 4.6 EB

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

# Memory-Constrained Data Structures

What method would you use to build a search index on a collection of objects *in a memory-constrained environment*?

**Constrained by resource limitations**

| | |
|---|---|
| cache | < 1 second |
| RAM | Hours - Days |
| disk | Months |
| network | Years |

(Estimates are Time x 1 billion courtesy of https://gist.github.com/hellerbarde/2843375)

# 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 Use Cases

Which of the following problems can be solved with a bloom filter?

**A) Find the closest matching item to a query of interest**

**B) Check if a query exists in a dataset**

**C) Compare the similarity between two datasets**

**D) Count the number of unique items in a dataset**

# Cardinality

Sometimes its not possible or realistic to count all objects!
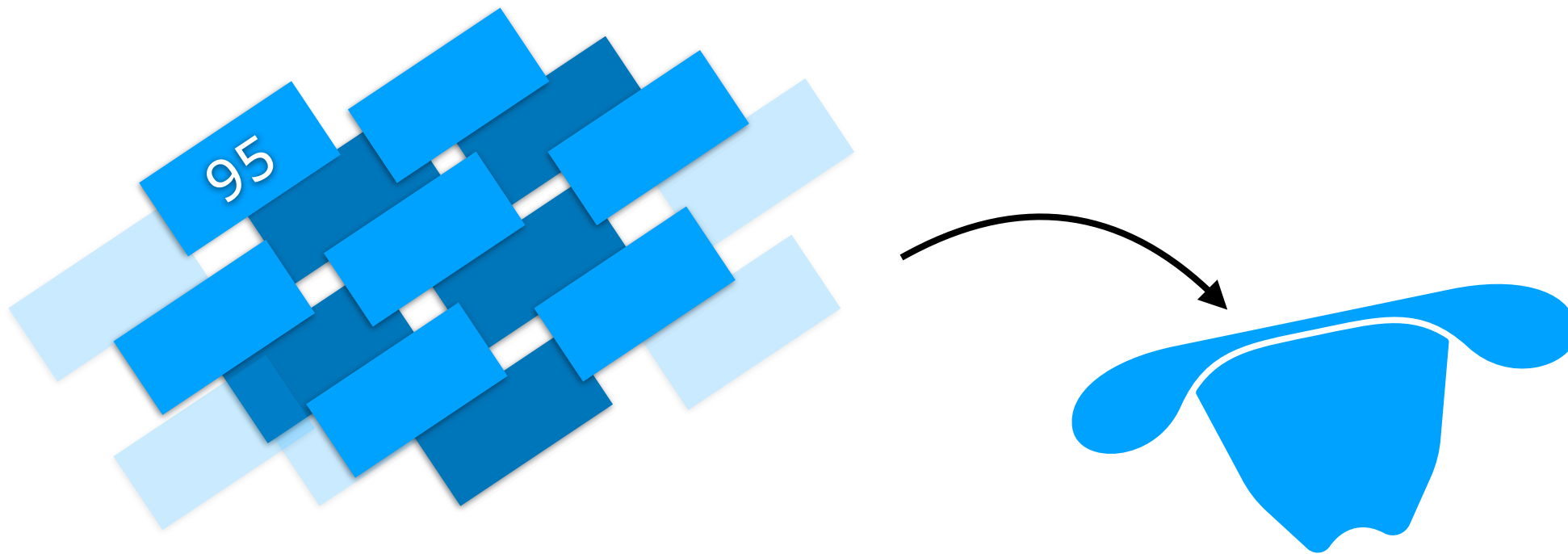


Estimate: 60 billion — 130 trillion

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?

95

# 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?



0                                                                          999

95

# Cardinality Estimation

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



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

# 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

Why do we care about "the hat problem"?



**_m_ possible minima**

| Key | Value |
| --- | --- |
|  |  |
|  |  |
|  |  |
|  |  |
|  |  |
|  |  |
|  |  |
|  |  |

# 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!

$h(x)$

0

$m - 1$

# 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…

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

0                                                                    1

# Cardinality Sketch

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

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

0                                                                                              1

# Cardinality Sketch

Consider an $N + 1$ draw:

$$\boxed{X_1}\boxed{X_2}\boxed{X_3} \cdots \boxed{X_N}\boxed{X_{N+1}}$$

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

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

$$0 \qquad M \qquad\qquad\qquad\qquad\qquad\qquad 1$$

# Cardinality Sketch

Consider an $N+1$ draw:

$$\boxed{X_1}\boxed{X_2}\boxed{X_3} \cdots \boxed{X_N \,\vdots\, X_{N+1}}$$

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

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

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

0         $M$                                   1

# Cardinality Sketch

Consider an $N + 1$ draw:

$$\boxed{X_1}\boxed{X_2}\boxed{X_3} \cdots \boxed{X_N}\boxed{X_{N+1}}$$

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

$X_{N+1}$ 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: $\boxed{X_1 \mid X_2 \mid X_3} \;\cdots\; \boxed{X_N \mid X_{N+1}}$

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

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

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

$$0 \qquad\qquad M \qquad\qquad\qquad\qquad\qquad\qquad\qquad 1$$

# Cardinality Sketch

Consider an $N + 1$ draw: $\boxed{X_1}\boxed{X_2}\boxed{X_3} \cdots \boxed{X_N}\boxed{X_{N+1}}$

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

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

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

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

# Cardinality Sketch

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

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

# Cardinality Sketch

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

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

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

$M_1$　　　　　$M_2$　　　　$M_3$　...　　$M_{k-1}$　　　$M_k$

# Cardinality Sketch

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

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

0                                                                                                1

$M_1$  $M_2$  $M_3$              $M_{k-1}$  $M_k$

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

*Averages $k$ estimates for* $\dfrac{1}{N+1}$

# Cardinality Sketch



True cardinality = 1,000

# 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**.



| 0.253 | 0.839 | 0.327 | 0.655 | 0.491 |

Hash

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

# 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


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


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

# Set Similarity Review

How can we describe how **similar** two sets are?

# Set Similarity Review

How can we describe how *similar* two sets are?

# Set Similarity Review

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.

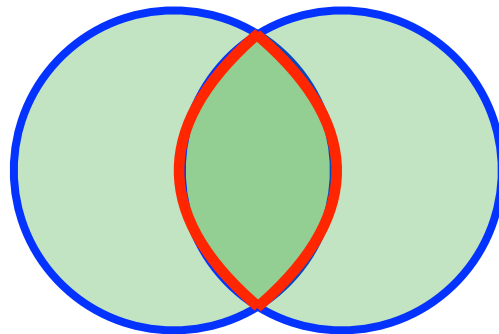$$J = \frac{|A \cap B|}{|A \cup B|}$$

$J$ is the ***Jaccard coefficient***

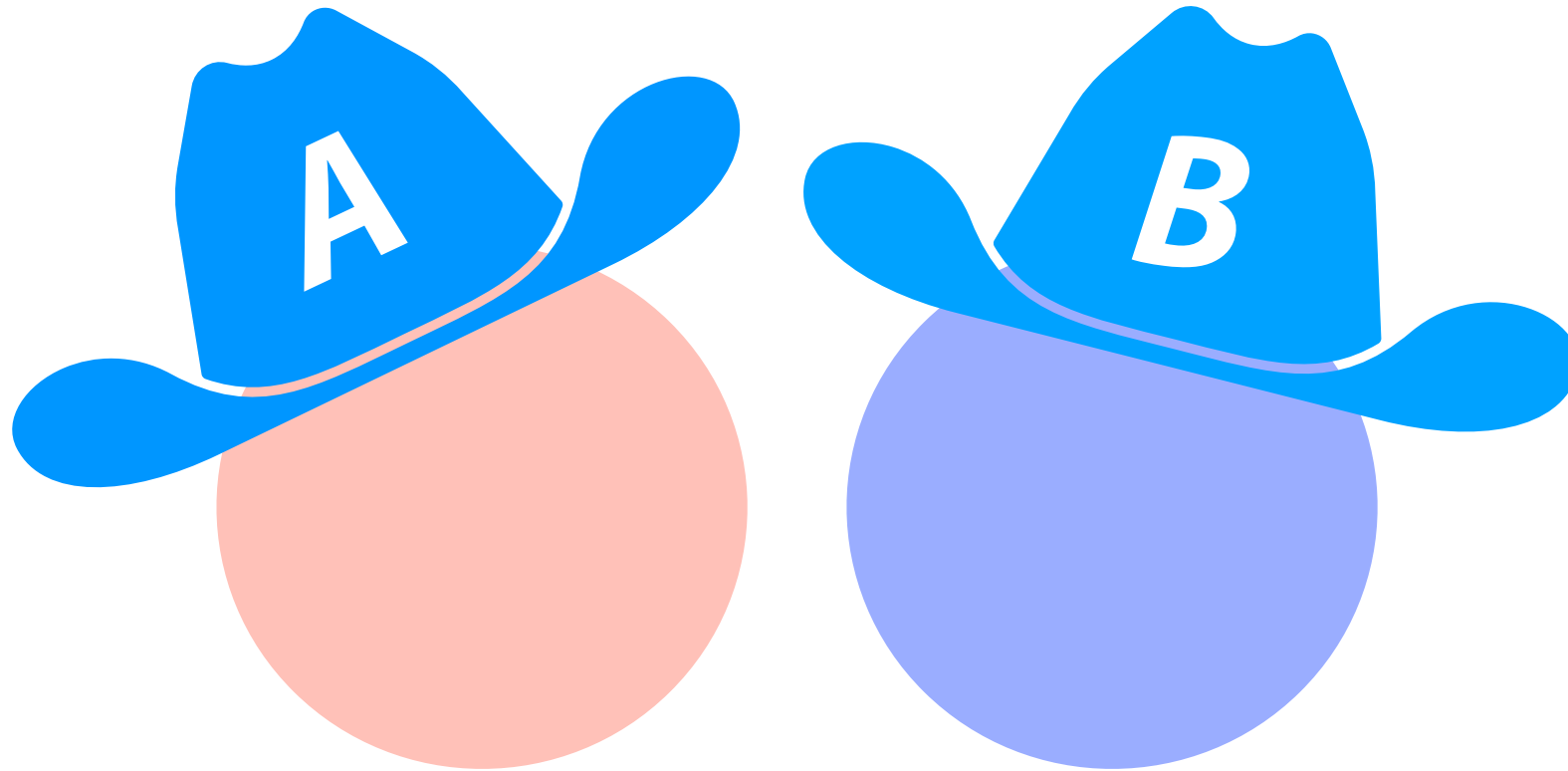# Set Similarity Review

$$\frac{|A \cap B|}{|A \cup B|} = 0$$

$$\frac{|A \cap B|}{|A \cup B|} = 1$$

$$0 < \frac{|A \cap B|}{|A \cup B|} < 1$$

# 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

# Similarity Sketches

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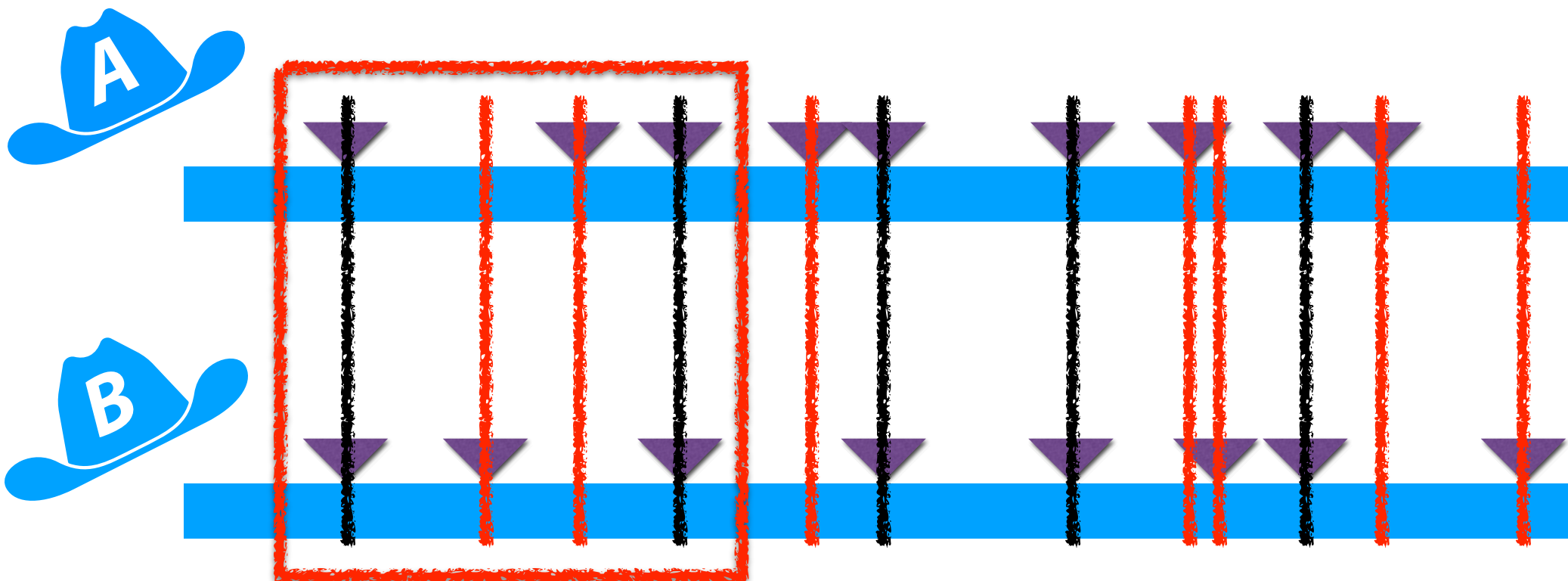