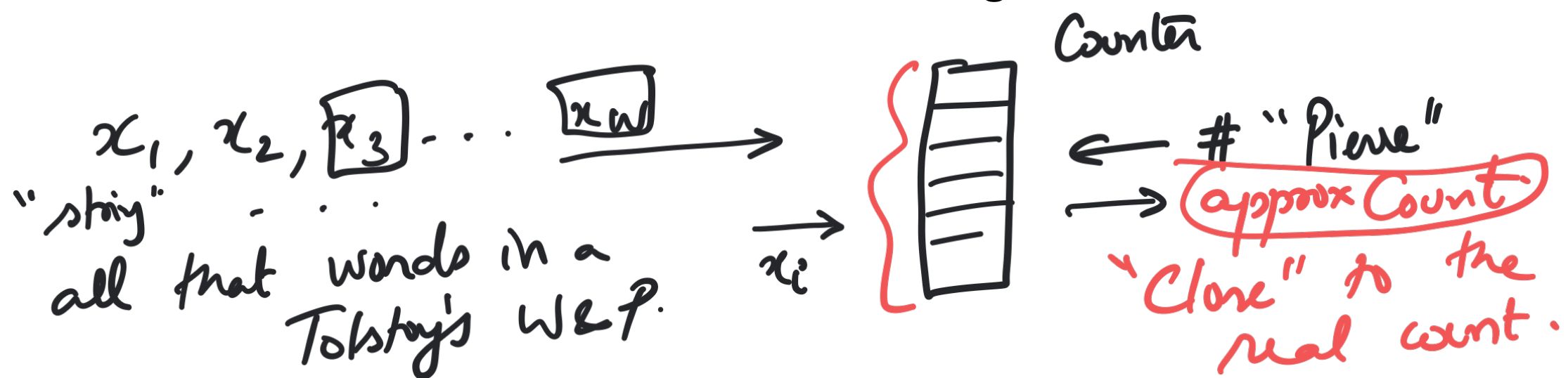


hash functions

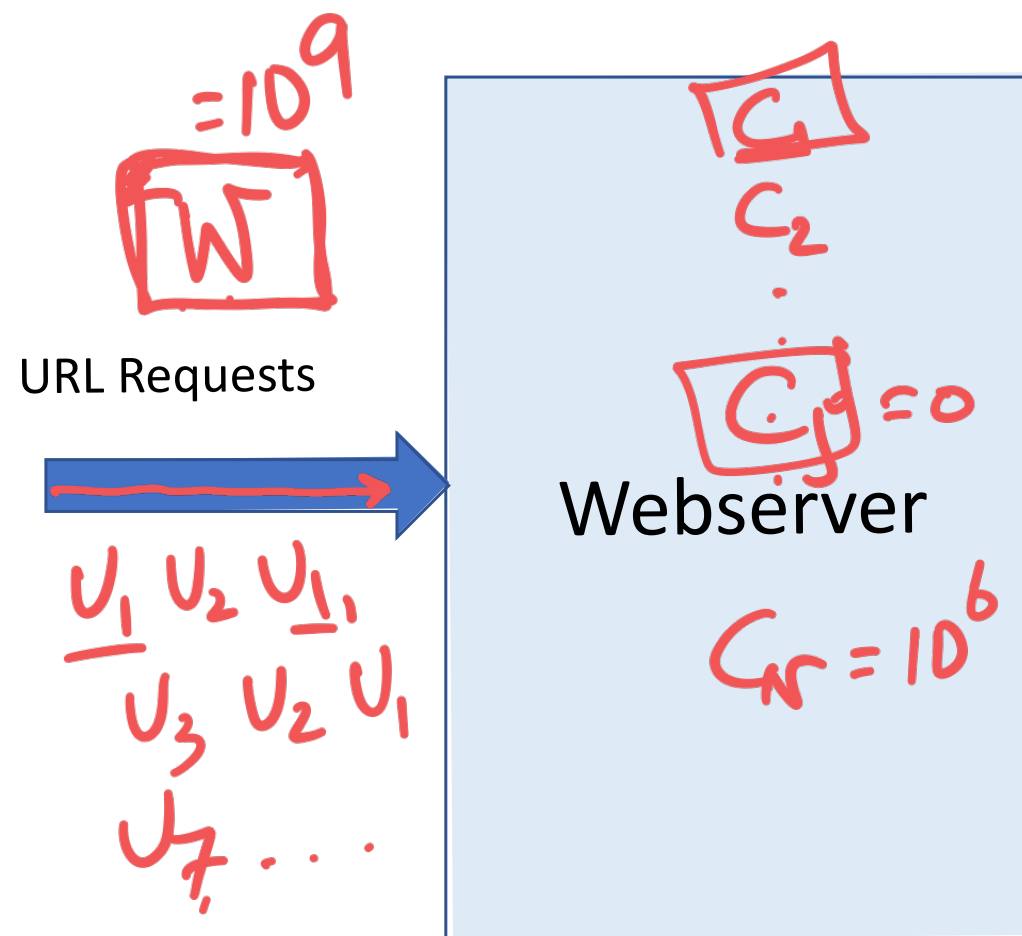
Count-Min Sketches

Sriram Sankaranarayanan

Data Structures and Algorithms



Problem: Count Items in a Data Stream



Problem: Keep count of how often each URL is requested.

Unique URLs: $U_1 U_2 U_3 U_4 U_5 U_6 \dots U_N$ $N \approx 10^6$

Have N distinct counters: $C(1), \dots, C(N)$

Each time URL U_j requested: Increment counter $C(j)$

Problem: N is humungous. not known in advance

Most URLs are requested very few times.

A few URLs are requested a lot of times.

Advantage: Approximate count within 10% of actual answer (say) is acceptable.

Approximate Counting Data-Structure

Stream of data: $x_1 x_2 x_3 x_4, \dots, x_W$

Each element of the stream: $x_j \in \{1, \dots, N\}$

Return $approxCount(j)$: must be within ϵW of true count
with probability at least δ

$W = \text{size of stream}$

$\epsilon = 10^{-6} \times 10^9 \approx 1000$

- Typical Numbers: $N \sim 10^8$, $W \sim 10^9$, $\epsilon \sim 10^{-6}$, $\delta = 0.99$
- From a stream of nearly 1 billion items each having a number between 1 and 100 Million, count how often each item occurs where the count is within 1000 of the true count at least 99% of the time.

how to choose m .

Basic Idea of Count-Min Sketch

Use m counters: $C(1) \dots C(m)$

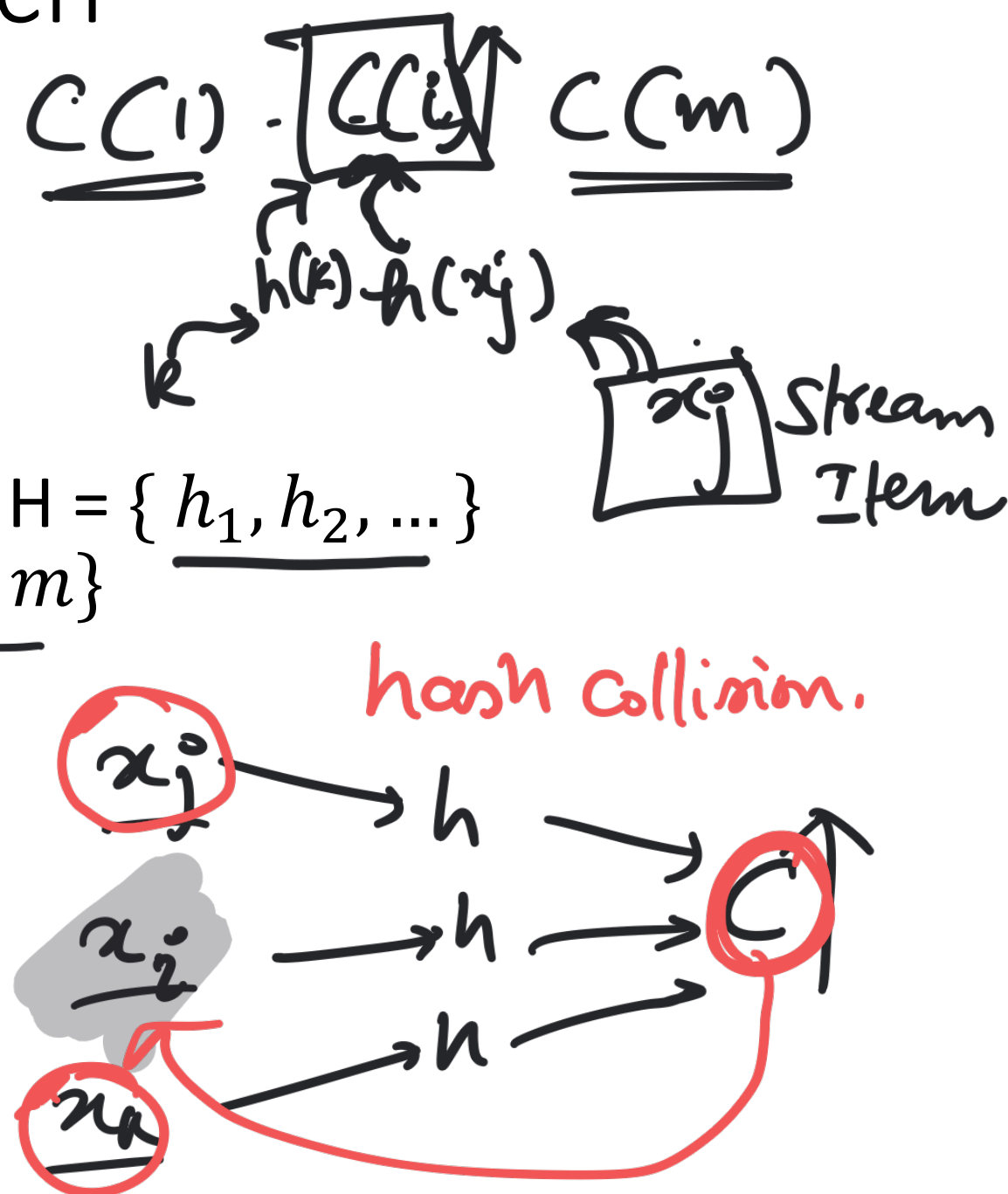
We will choose m later (expect $m \ll \cancel{M}$)

Draw a hash function at random from family $H = \{h_1, h_2, \dots\}$

$$h_i: \{1, \dots, \cancel{M}\} \rightarrow \{1, \dots, m\}$$

Stream Item x_j Increment($C(h(x_j))$)

$$\text{approxCount}(k) = C(h(k))$$



Count-Min Sketch Error Analysis

- approxCount(j) \geq count(j)

$$\Pr_{h \in H} (h(i) = h(j)) \leq \frac{C}{m}$$

$$E(\text{approxCount}(j)) = \text{count}(j) + \sum_{\substack{i=1 \\ i \neq j}}^N \frac{\Pr(h(i) = h(j))}{\text{count}(i)}$$

$$\leq \text{count}(j) + \frac{C}{m} \cdot \sum_{i=1}^N \text{count}(i)$$

$$\leq \text{count}(j) + \frac{C}{m} W$$

$$E(\text{approxCount}(j) - \text{count}(j)) \leq \frac{C}{m} W$$

Count-Min Sketch: Choosing m

$$\mathbb{E}(\text{error}(j)) \leq \boxed{\frac{CW}{m}}$$

$$\text{Pr.}(\underbrace{\text{error}(j)}_{\substack{\text{approxCount} - \text{count} \\ \geq 0}} \geq \frac{\epsilon W}{t^{\uparrow}}) \leq \frac{\cancel{CW}}{m \cdot \epsilon W} \quad (\text{M. Ineq.})$$

$$\therefore \boxed{\underline{\underline{m}} = \frac{2.72 \dots e C}{\epsilon}}$$

$$\leq \frac{C}{m\epsilon} \leq \boxed{\frac{1}{e}}$$

Not good enough

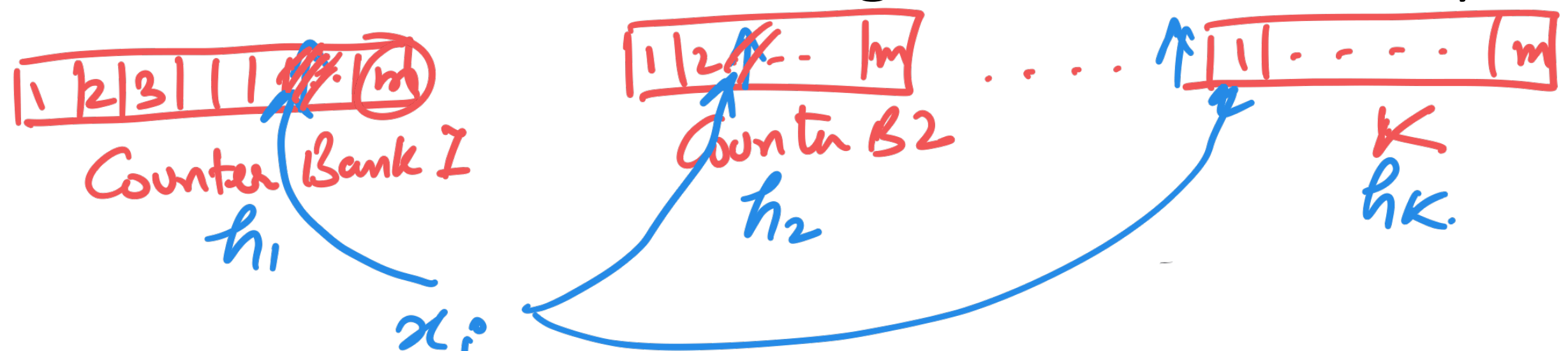
$E(x)$



Markov Ineq.

$$\text{Pr}(x \geq t) \leq \frac{\boxed{E(x)}}{\boxed{t}}$$

Count-Min Sketch: Reducing Error Probability



$$\text{approx Count}(k) = \min_i (C_1(h_1(k)), C_2(h_2(k)), \dots, C_k(h_k(k)))$$

$$\Pr(\text{error}_j(h_j(k)) \geq \lceil \epsilon N \rceil) \leq 1/e$$

All of them make

$$\left(\frac{1}{e} \right)^K \leq 1 - \delta \quad \left\{ \begin{array}{l} K \geq \frac{-\ln(1-\delta)}{-\ln(0.01)} \end{array} \right.$$

Count-Min Sketch: Overall Algorithm

- Initialize K counter-banks with hash functions : $\underline{h_1}, \underline{h_2}, \dots, \underline{h_K}$ *randomly chosen*
- Stream item x_j
 - Increment $\underline{C_1(h_1(x_j))} \uparrow \underline{C_2(h_2(x_j))} \uparrow \dots, \underline{C_K(h_K(x_j))} \uparrow$ $\Theta(K)$
- Query count of k
 - $\underset{\uparrow}{\text{approxCount}(k)} = \min(\underline{C_1(h_1(k))}, \underline{C_2(h_2(k))}, \dots, \underline{C_K(h_K(k))})$ $\Theta(K)$

$$K \approx S$$

Count-Min Sketch: Some actual numbers

- Stream of 1 Billion Items 10^9
- $\epsilon = 10^{-6}$ (tolerate error of upto 1000)
- $\delta = 0.9$ for 90% of my queries I wish the answer to be within 1000 of actual answer
- $m = \frac{e}{\epsilon} \approx 3 \times 10^6$, $K = -\ln(1 - \delta) \approx 3$
- Use 3 banks of 3 million counters, each. $10M \sim 10^7$
 - Guarantees that approxCount will be within 1000 of true at least answer 90% of the time.

- Merkle : Blockchain.

- Locality Sensitive . .

. . .