

# Advanced Algorithms

## Hashing and Sketching

尹一通 Nanjing University, 2023 Fall

# Balls into Bins



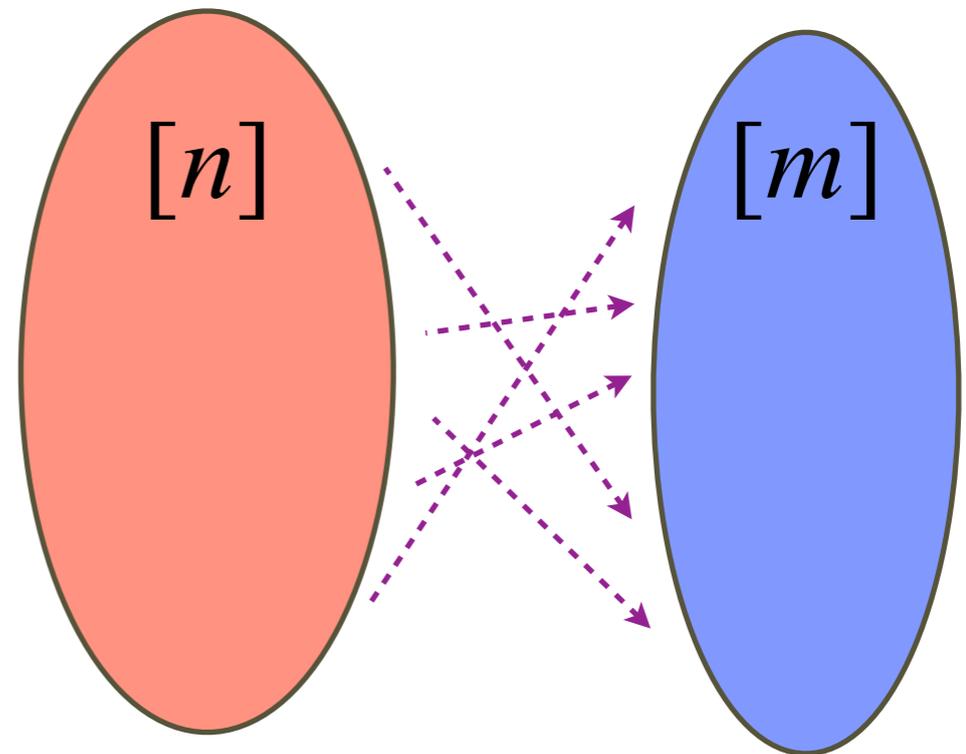
# Balls into Bins (Random Function)

- $n$  balls into  $m$  bins:

$$\Pr[\text{assignment}] = \frac{1}{m} \cdots \frac{1}{m} = \frac{1}{m^n}$$

- uniform random function:

$$\Pr[f] = \frac{1}{|[n] \rightarrow [m]|} = \frac{1}{m^n}$$



uniform random function

$$f : [n] \rightarrow [m]$$

1-1	birthday
on-to	coupon collector
pre-image size	occupancy

# Birthday Paradox

## Paradox:

- (i) a statement that leads to a contradiction;
- (ii) a situation which defies intuition.



In a class of  $m > 57$  students, with  $>99\%$  probability, there are two students with the same birthday.

Assumption: birthdays are uniformly & independently distributed.

$n$  balls are thrown into  $m$  bins:

event  $\mathcal{E}$ : each bin receives  $\leq 1$  balls

# Birthday Paradox

$n$  balls are thrown into  $m$  bins:

event  $\mathcal{E}$ : each bin receives  $\leq 1$  balls

$$\begin{aligned}\Pr[\mathcal{E}] &= \frac{\left| [n] \xrightarrow{1-1} [m] \right|}{\left| [n] \rightarrow [m] \right|} = \frac{m(m-1)\cdots(m-n+1)}{m^n} \\ &= \prod_{i=0}^{n-1} \left( 1 - \frac{i}{m} \right)\end{aligned}$$

# Birthday Paradox

$n$  balls are thrown into  $m$  bins:

event  $\mathcal{E}$ : each bin receives  $\leq 1$  balls

Suppose that balls are thrown one-by-one:

$\Pr[\mathcal{E}] = \Pr[\text{all } n \text{ balls are thrown into distinct bins}]$

chain rule  $= \prod_{i=1}^n \Pr[\text{the } i\text{th ball is thrown into an empty bin} \mid$   
first  $i - 1$  balls are thrown into distinct bins]

$$= \prod_{i=1}^n \left(1 - \frac{i-1}{m}\right) = \prod_{i=0}^{n-1} \left(1 - \frac{i}{m}\right)$$

# Birthday Paradox

$n$  balls are thrown into  $m$  bins:

event  $\mathcal{E}$ : each bin receives  $\leq 1$  balls

(Taylor:  $1 - x \approx e^{-x}$  for  $x = o(1)$ )

$$\Pr[\mathcal{E}] = \prod_{i=0}^{n-1} \left(1 - \frac{i}{m}\right) \approx \prod_{i=0}^{n-1} e^{-\frac{i}{m}} \approx e^{-n^2/2m}$$

**Formally:**  $e^{-(1+o(1))n^2/2m} \leq \prod_{i=0}^{n-1} \left(1 - \frac{i}{m}\right) \leq e^{-(1-o(1))n^2/2m}$   
(assuming  $n \ll m$ )

when  $n = \sqrt{2m \ln \frac{1}{p}}$   $\implies \Pr[\mathcal{E}] = (1 \pm o(1))p$

# Birthday Paradox

$n$  balls are thrown into  $m$  bins:

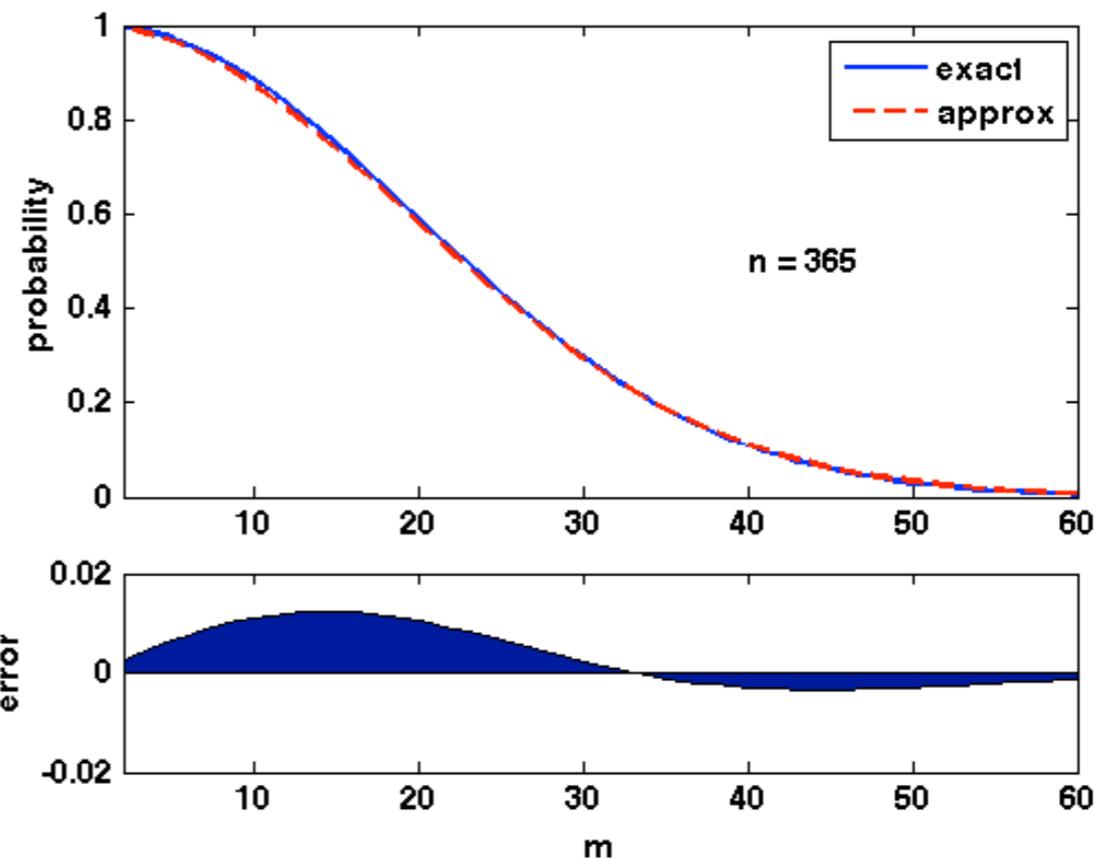
event  $\mathcal{E}$ : each bin receives  $\leq 1$  balls

$$\Pr[\mathcal{E}] = \prod_{i=0}^{n-1} \left(1 - \frac{i}{m}\right)$$

Formally:  $e^{-(1+o(1))n^2/2m}$

(assuming  $n \ll m$ )

$$\text{when } n = \sqrt{2m \ln \frac{1}{p}} \implies \Pr[\mathcal{E}] = (1 \pm o(1))p$$



# Hash Tables & Filters

# Data Structure for Set

**Data:** a set  $S$  of  $n$  items  $x_1, x_2, \dots, x_n \in U = [N]$

**Query:** an item  $x \in U$

Determine whether  $x \in S$ .

- **Space cost:** size of data structure (in bits)
  - **entropy** of a set:  $\log \binom{N}{n} = O(n \log N)$  bits (when  $N \gg n$ )
- **Time cost:** time to answer a query (in memory accesses)
- **Balanced tree:**  $O(n \log N)$  space,  $O(\log n)$  time
- **Perfect hashing:**  $O(n \log N)$  space,  $O(1)$  time

# Perfect Hashing

$S = \{a, b, c, d, e, f\} \subseteq [N]$  of size  $n$

uniform  
random

$h$   $[N] \rightarrow [m]$

no collision

$\Pr[\text{perfect}] \approx e^{-n^2/2m} > 1/2$

Table  $T$ : 

$e$	$b$		$d$		$f$		$c$	$a$	
-----	-----	--	-----	--	-----	--	-----	-----	--

 $m = n^2$   
Birthday

**SUHA:** Simple Uniform Hash Assumption

**Query( $x$ ):**

retrieve hash function  $h$ ;

check whether  $T[h(x)] = x$ ;

# Universal Hashing

## Universal Hash Family (Carter and Wegman 1979):

A family  $\mathcal{H}$  of hash functions in  $U \rightarrow [m]$  is  **$k$ -universal** if for any distinct  $x_1, \dots, x_k \in U$ ,

$$\Pr_{h \in \mathcal{H}} [h(x_1) = \dots = h(x_k)] \leq \frac{1}{m^{k-1}}.$$

Moreover,  $\mathcal{H}$  is **strongly  $k$ -universal** ( $k$ -wise independent) if for any distinct  $x_1, \dots, x_k \in U$  and any  $y_1, \dots, y_k \in [m]$ ,

$$\Pr_{h \in \mathcal{H}} \left[ \bigwedge_{i=1}^k h(x_i) = y_i \right] = \frac{1}{m^k}.$$

# $k$ -Universal Hash Family

hash functions  $h : U \rightarrow [m]$

- **Linear congruential hashing:**

- Represent  $U \subseteq \mathbb{Z}_p$  for sufficiently large prime  $p$

- $h_{a,b}(x) = ((ax + b) \bmod p) \bmod m$

- $\mathcal{H} = \left\{ h_{a,b} \mid a \in \mathbb{Z}_p \setminus \{0\}, b \in \mathbb{Z}_p \right\}$

**Theorem:**

The linear congruential family  $\mathcal{H}$  is 2-wise independent.

- **Degree- $k$  polynomial in finite field with random coefficients**

- Hashing between binary fields:  $GF(2^w) \rightarrow GF(2^l)$

$$h_{a,b}(x) = (a * x + b) \gg (w-1)$$

# Birthday Paradox (pairwise independence)

$n$  balls are thrown into  $m$  bins: **by 2-universal hashing**  
event  $\mathcal{E}$ : each bin receives  $\leq 1$  balls

- Location of  $n$  balls:  $X_1, X_2, \dots, X_n \in [m]$
- Total # of collisions:

$$Y = \sum_{i < j} I[X_i = X_j]$$

- Linearity of expectation:

$$\mathbb{E}[Y] = \sum_{i < j} \Pr[X_i = X_j] \leq \binom{n}{2} \frac{1}{m}$$

**2-universal**

# Markov's Inequality

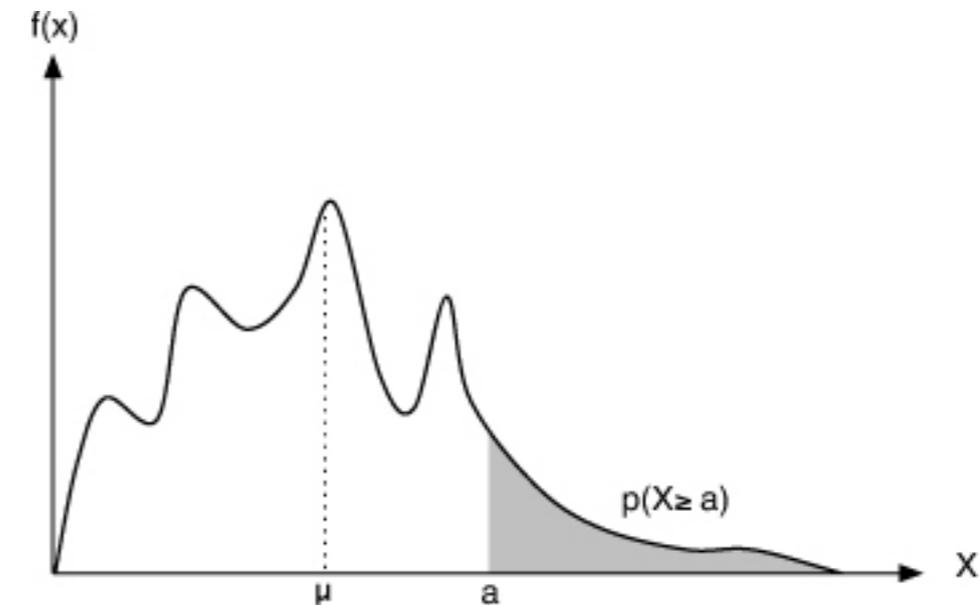
## Markov's Inequality

For *nonnegative* random variable  $X$ , for any  $t > 0$ ,

$$\Pr[X \geq t] \leq \frac{\mathbb{E}[X]}{t}$$

$$\text{Let } Y = \begin{cases} 1 & X \geq t \\ 0 & \text{o.w.} \end{cases} \implies Y \leq \left\lfloor \frac{X}{t} \right\rfloor \leq \frac{X}{t}$$

$$\Pr[X \geq t] = \mathbb{E}[Y] \leq \mathbb{E}\left[\frac{X}{t}\right] = \frac{\mathbb{E}[X]}{t}$$



# Birthday Paradox (pairwise independence)

$n$  balls are thrown into  $m$  bins: **by 2-universal hashing**  
event  $\mathcal{E}$ : each bin receives  $\leq 1$  balls

- Location of  $n$  balls:  $X_1, X_2, \dots, X_n \in [m]$

- Total # of collisions:

$$Y = \sum_{i < j} I[X_i = X_j]$$

- Linearity of expectation:

$$\mathbb{E}[Y] = \sum_{i < j} \Pr[X_i = X_j] \leq \binom{n}{2} \frac{1}{m}$$

**2-universal**

**when**  
 $n \leq \sqrt{2m\epsilon}$

- Markov's inequality:  $\Pr[\neg \mathcal{E}] = \Pr[Y \geq 1] \leq \mathbb{E}[Y] \leq \epsilon$

# Perfect Hashing

$S = \{a, b, c, d, e, f\} \subseteq [N]$  of size  $n$

**2-universal**  $h: [N] \rightarrow [m]$        $\Pr[\textit{imperfect}] = \frac{n(n-1)}{2m}$

Table  $T$ : 

$e$	$b$		$d$		$f$		$c$	$a$	
-----	-----	--	-----	--	-----	--	-----	-----	--

 $m$

For 2-universal family  $\mathcal{H}$  from  $[N]$  to  $[m]$ , if  $m > \binom{n}{2}$ , for any  $S \subseteq [N]$  of size  $n$ , there is an  $h \in \mathcal{H}$  that cause no collisions over  $S$ .

**Query( $x$ ):**

retrieve hash function  $h$ ;

check whether  $T[h(x)] = x$ ;

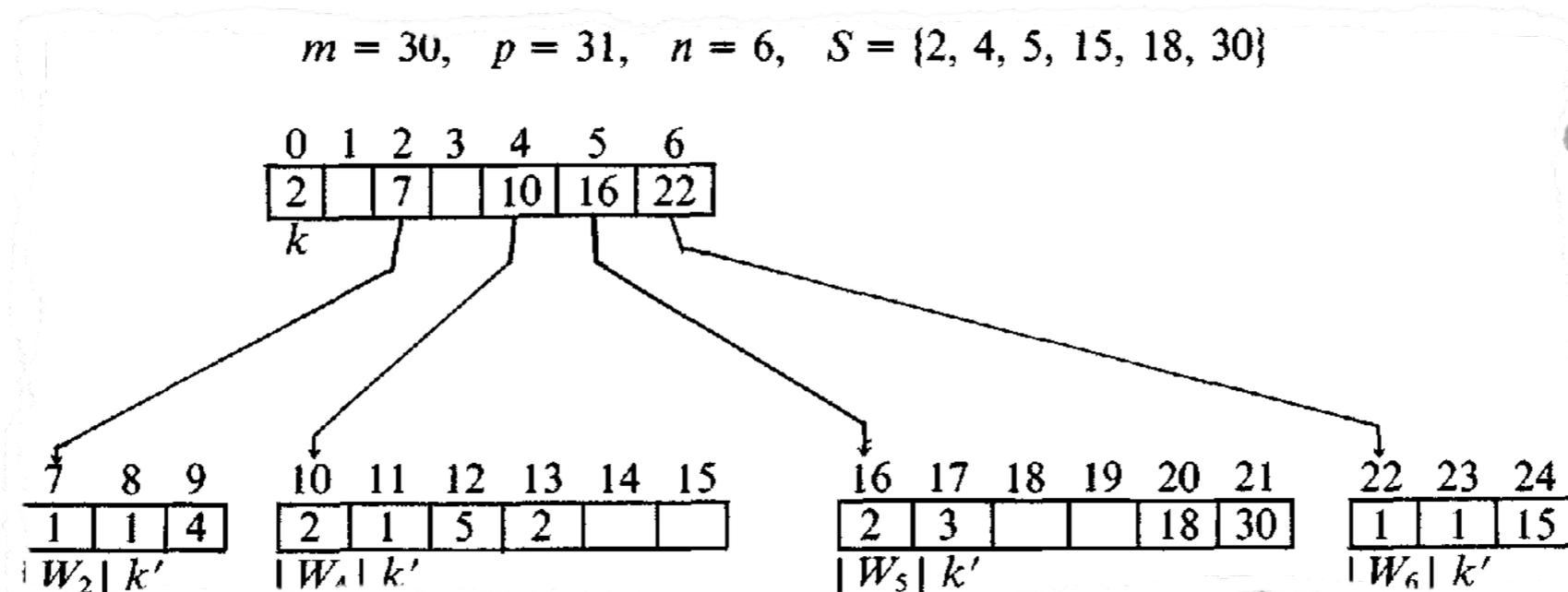
# FKS Perfect Hashing

(Fredman, Komlós, Szemerédi, 1984)

**Data:** a set  $S$  of  $n$  items  $x_1, x_2, \dots, x_n \in U = [N]$

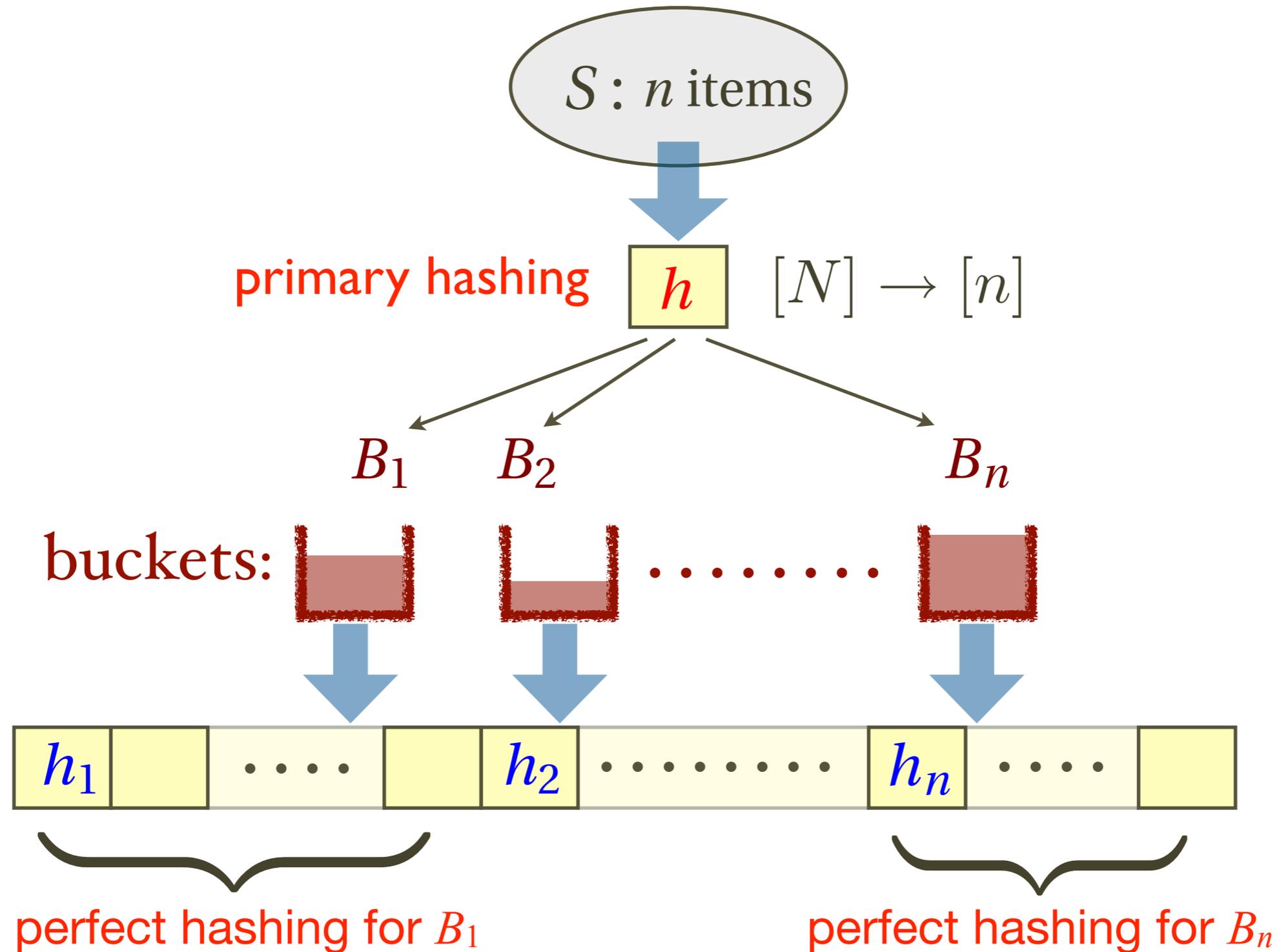
**Query:** an item  $x \in U$

Determine whether  $x \in S$ .



- Space cost:  $O(n)$  words (each of  $O(\log N)$  bits)
- Time cost:  $O(1)$  for each query in the worst case

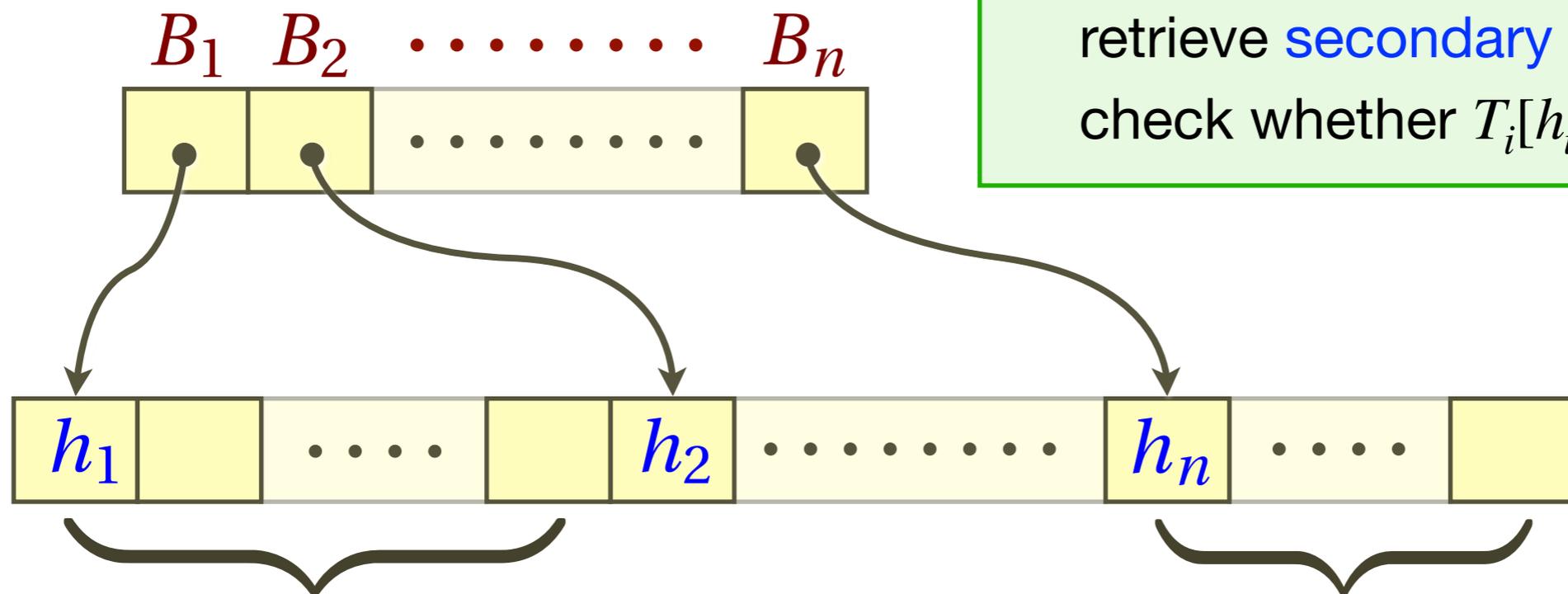
# FKS Perfect Hashing



# FKS Perfect Hashing

Set  $S \subseteq [N]$  of size  $n$

$$h: [N] \rightarrow [n]$$



**Query(x):**

retrieve **primary hash**  $h$ ;

goto **bucket**  $i = h(x)$ ;

retrieve **secondary hash**  $h_i$ ;

check whether  $T_i[h_i(x)] = x$ ;

perfect hashing for  $B_1$   
using space  $|B_1|^2$

perfect hashing for  $B_n$   
using space  $|B_n|^2$

- $\exists h_1, \dots, h_n$  from 2-universal family s.t.  $h_i$  is perfect for  $B_i$  for all  $i$

# Collision Number

$n$  balls are thrown into  $m$  bins by **2-universal hashing**

- Location of  $n$  bins:  $X_1, X_2, \dots, X_n \in [m]$

$$\text{Collision \#}: Y = \sum_{i < j} I[X_i = X_j]$$

- Linearity of expectation:

$$\mathbb{E}[Y] = \sum_{i < j} \Pr[X_i = X_j] \leq \binom{n}{2} \frac{1}{m}$$

**2-universal**

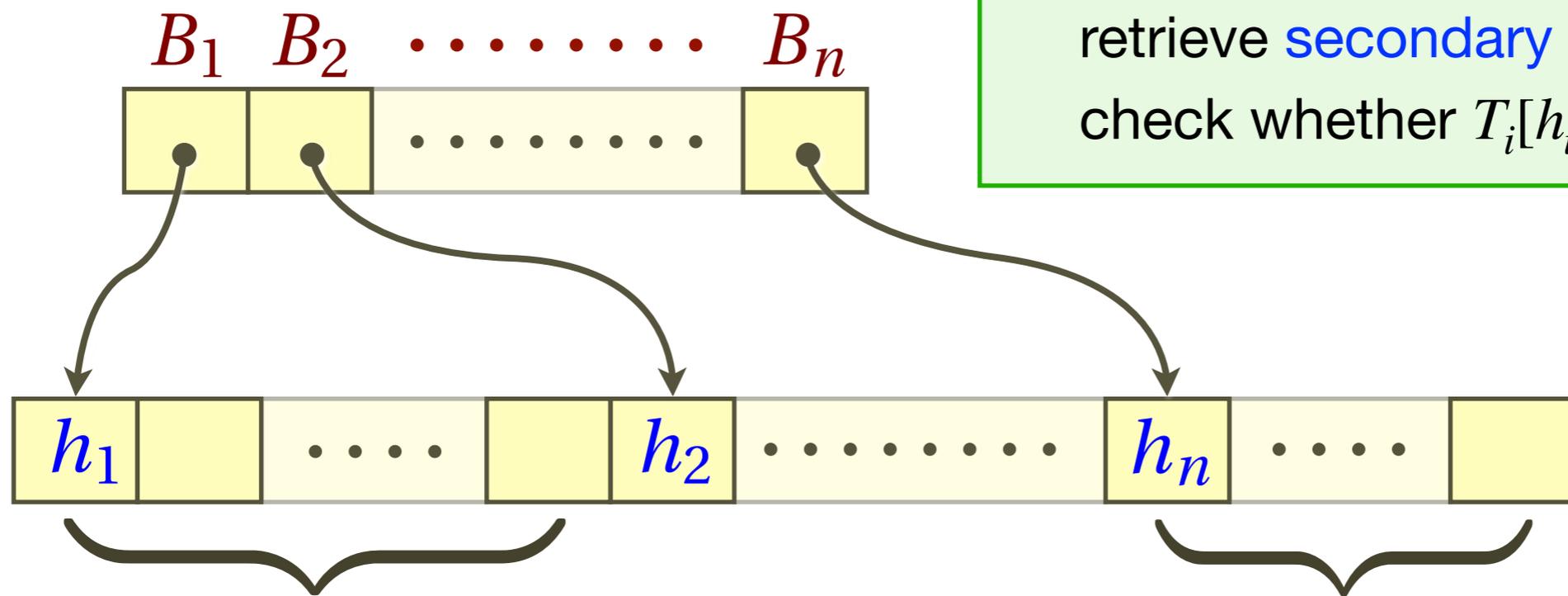
- Size of the  $i$ -th bin:  $|B_i|$

$$Y = \sum_{i=1}^n \binom{|B_i|}{2} = \frac{1}{2} \sum_{i=1}^n |B_i| (|B_i| - 1) \implies \mathbb{E} \left[ \sum_{i=1}^n |B_i|^2 \right] = \frac{n(n-1)}{m} + n$$

# FKS Perfect Hashing

Set  $S \subseteq [N]$  of size  $n$

$$h: [N] \rightarrow [n]$$



**Query( $x$ ):**

retrieve **primary hash**  $h$ ;

goto **bucket**  $i = h(x)$ ;

retrieve **secondary hash**  $h_i$ ;

check whether  $T_i[h_i(x)] = x$ ;

perfect hashing for  $B_1$   
using space  $|B_1|^2$

perfect hashing for  $B_n$   
using space  $|B_n|^2$

- $\exists h$  from a 2-universal family s.t. the total space cost is  $O(n)$

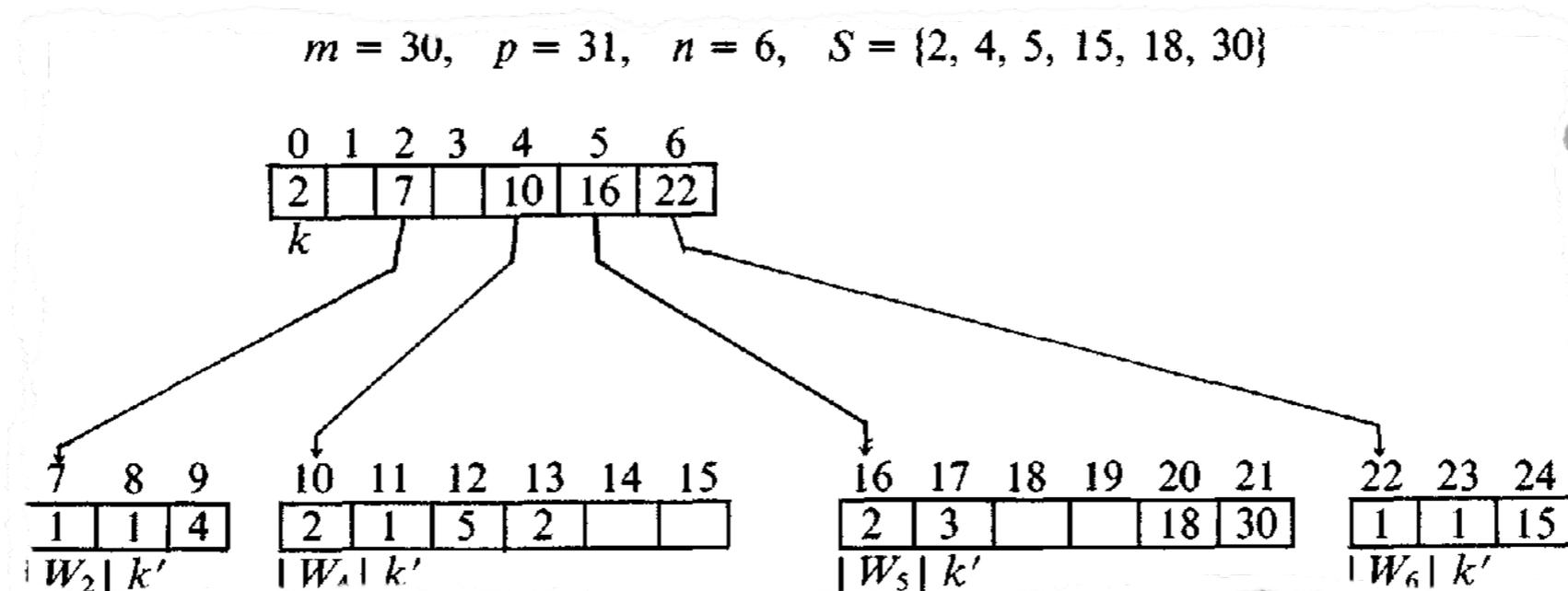
# FKS Perfect Hashing

(Fredman, Komlós, Szemerédi, 1984)

**Data:** a set  $S$  of  $n$  items  $x_1, x_2, \dots, x_n \in U = [N]$

**Query:** an item  $x \in U$

Determine whether  $x \in S$ .



- $O(n \log N)$  space,  $O(1)$  time in the worst case
- Dynamic version: [Dietzfelbinger, Karlin, Mehlhorn, Meyer auf der Heide, Rohnert, Tarjan, 1984]

# Optimal Dynamic Perfect Hashing

(Upper Bound, STOC 2022)

On the Optimal Time/Space Tradeoff for Hash Tables

Michael A. Bender  
Stony Brook University

Martín Farach-Colton  
Rutgers University

John Kuszmaul  
Yale University

William Kuszmaul  
MIT

Mingmou Liu  
NTU

## Abstract

For nearly six decades, the central open question in the study of hash tables has been to determine the optimal achievable tradeoff curve between time and space. State-of-the-art hash tables offer the following guarantee: If keys/values are  $\Theta(\log n)$  bits each, then it is possible to achieve constant-time insertions/deletions/queries while wasting only  $O(\log \log n)$  bits of space per key when compared to the information-theoretic optimum. Even prior to this bound being achieved, the target of  $O(\log \log n)$  wasted bits per key was known to be a natural end goal, and was proven to be optimal for a number of closely related problems (e.g., stable hashing, dynamic retrieval, and dynamically-resized filters).

This paper shows that  $O(\log \log n)$  wasted bits per key is not the end of the line for hashing. In fact, for any  $k \in [\log^* n]$ , it is possible to achieve  $O(k)$ -time insertions/deletions,  $O(1)$ -time queries, and

$$O(\log^{(k)} n) = O\left(\underbrace{\log \log \cdots \log n}_k\right)$$

wasted bits per key (all with high probability in  $n$ ). This means that, each time we increase insertion/deletion time by an *additive constant*, we reduce the wasted bits per key *exponentially*. We further show that this tradeoff curve is the best achievable by any of a large class of hash tables, including any hash table designed using the current framework for making constant-time hash tables succinct.

# Optimal Dynamic Perfect Hashing

## (Lower Bound, FOCS 2023)

Tight Cell-Probe Lower Bounds for Dynamic Succinct Dictionaries

Tianxiao Li <sup>\*</sup>    Jingxun Liang <sup>†</sup>    Huacheng Yu <sup>‡</sup>    Renfei Zhou <sup>§</sup>

### Abstract

A dictionary data structure maintains a set of at most  $n$  keys from the universe  $[U]$  under key insertions and deletions, such that given a query  $x \in [U]$ , it returns if  $x$  is in the set. Some variants also store values associated to the keys such that given a query  $x$ , the value associated to  $x$  is returned when  $x$  is in the set.

This fundamental data structure problem has been studied for six decades since the introduction of hash tables in 1953. A hash table occupies  $O(n \log U)$  bits of space with constant time per operation in expectation. There has been a vast literature on improving its time and space usage. The state-of-the-art dictionary by Bender, Farach-Colton, Kuszmaul, Kuszmaul and Liu [BFCK<sup>+</sup>22] has space consumption close to the *information-theoretic optimum*, using a total of

$$\log \binom{U}{n} + O(n \log^{(k)} n)$$

bits, while supporting all operations in  $O(k)$  time, for any parameter  $k \leq \log^* n$ . The term  $O(\log^{(k)} n) = O(\underbrace{\log \cdots \log}_k n)$  is referred to as the *wasted bits per key*.

In this paper, we prove a matching *cell-probe* lower bound: For  $U = n^{1+\Theta(1)}$ , any dictionary with  $O(\log^{(k)} n)$  wasted bits per key must have expected operational time  $\Omega(k)$ , in the cell-probe model with word-size  $w = \Theta(\log U)$ . Furthermore, if a dictionary stores values of  $\Theta(\log U)$  bits, we show that *regardless of the query time*, it must have  $\Omega(k)$  expected update time. It is worth noting that this is the first cell-probe lower bound on the trade-off between space and update time for general data structures.

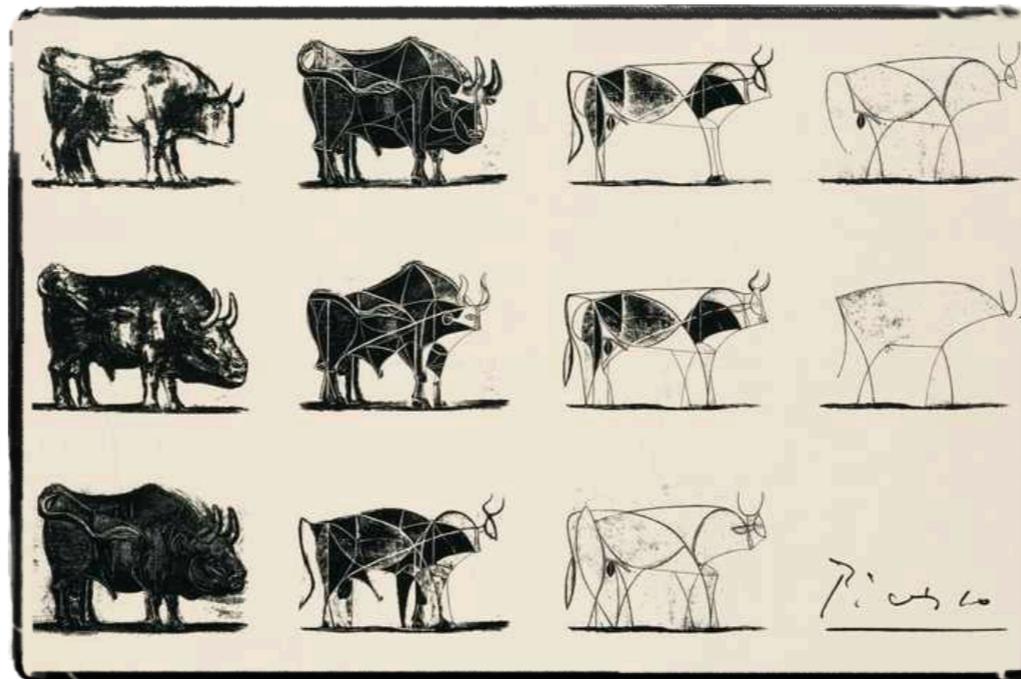
# Data Structure for Set

**Data:** a set  $S$  of  $n$  items  $x_1, x_2, \dots, x_n \in U = [N]$

**Query:** an item  $x \in U$

Determine whether  $x \in S$ .

- **Space cost:** size of data structure (in bits)
  - **entropy** of a set:  $\log \binom{N}{n} = O(n \log N)$  bits (when  $N \gg n$ )
- **Sketch:** lossy representation of  $S$  using  $<$  entropy space



# Approximate Set

**Data:** a set  $S$  of  $n$  items  $x_1, x_2, \dots, x_n \in U = [N]$

**Query:** an item  $x \in U$

Answer whether  $x \in S$  with bounded error.

- uniform hash function  $h : U \rightarrow [m]$  ( $m$  to be fixed)

**Data Structure:** bit array  $A \in \{0,1\}^m$

$A$  is initialized to all 0's;

for each  $x_i \in S$ : set  $A[h(x_i)] = 1$ ;

**Query**  $x$ : answer “yes” iff  $A[h(x)] = 1$

- $x \in S$ : always correct
- $x \notin S$ : **false positive**  $\Pr [A[h(x)] = 1] = 1 - (1 - 1/m)^n \approx 1 - e^{-n/m}$

# Bloom Filters (Bloom 1970)

**Data:** a set  $S$  of  $n$  items  $x_1, x_2, \dots, x_n \in U = [N]$

**Query:** an item  $x \in U$

Answer whether  $x \in S$  with bounded error.

- uniform & independent hash function  $h_1, \dots, h_k : U \rightarrow [m]$   
( $k$  and  $m$  to be fixed)

**Data Structure:** bit array  $A \in \{0,1\}^m$

$A$  is initialized to all 0's;

for each  $x_i \in S$ : set  $A[h_j(x_i)] = 1$  for all  $1 \leq j \leq k$ ;

**Query**  $x$ : “yes” iff  $A[h_j(x)] = 1$  for all  $1 \leq j \leq k$

# Bloom Filters (Bloom 1970)

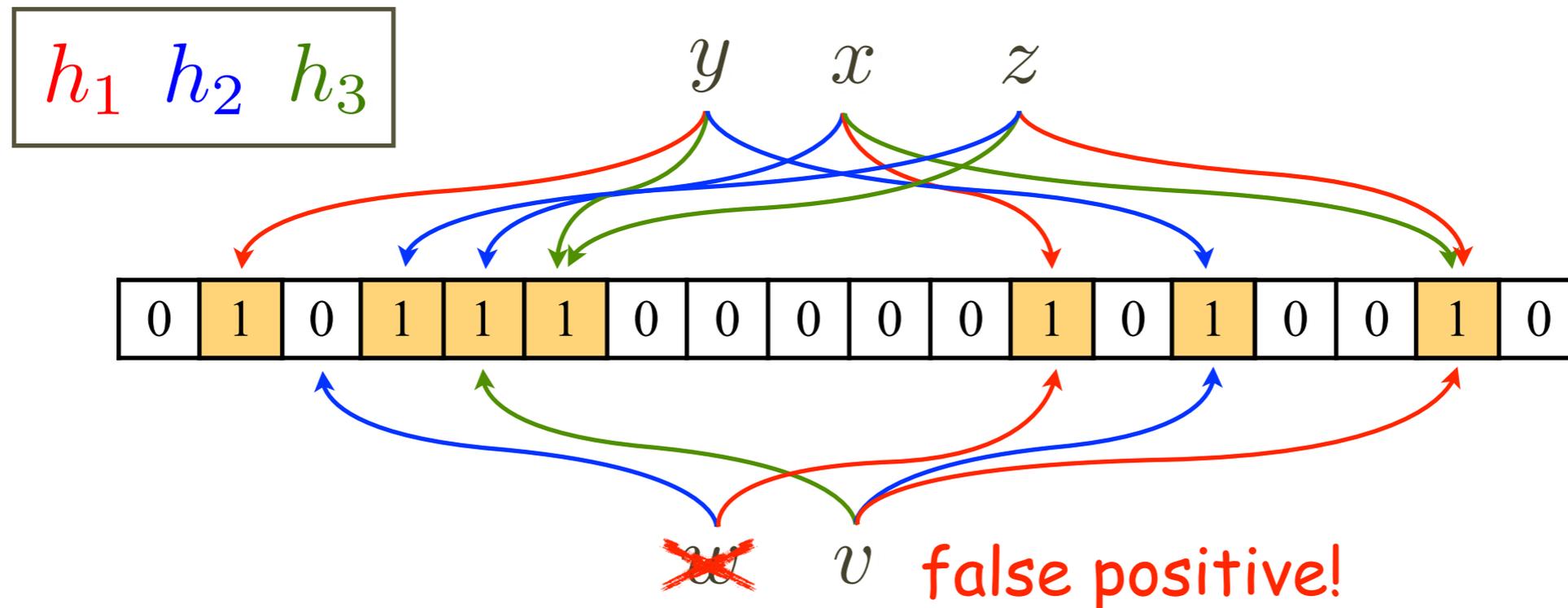
- uniform & independent hash function  $h_1, \dots, h_k : U \rightarrow [m]$

**Data Structure:** bit array  $A \in \{0,1\}^m$

$A$  is initialized to all 0's;

for each  $x_i \in S$ : set  $A[h_j(x_i)] = 1$  for all  $1 \leq j \leq k$ ;

**Query**  $x$ : “yes” iff  $A[h_j(x)] = 1$  for all  $1 \leq j \leq k$



**Data:** set  $S \subseteq U$  of size  $n$       **Query:**  $x \in U$

- uniform & independent hash function  $h_1, \dots, h_k : U \rightarrow [m]$

**Data Structure:** bit array  $A \in \{0,1\}^m$

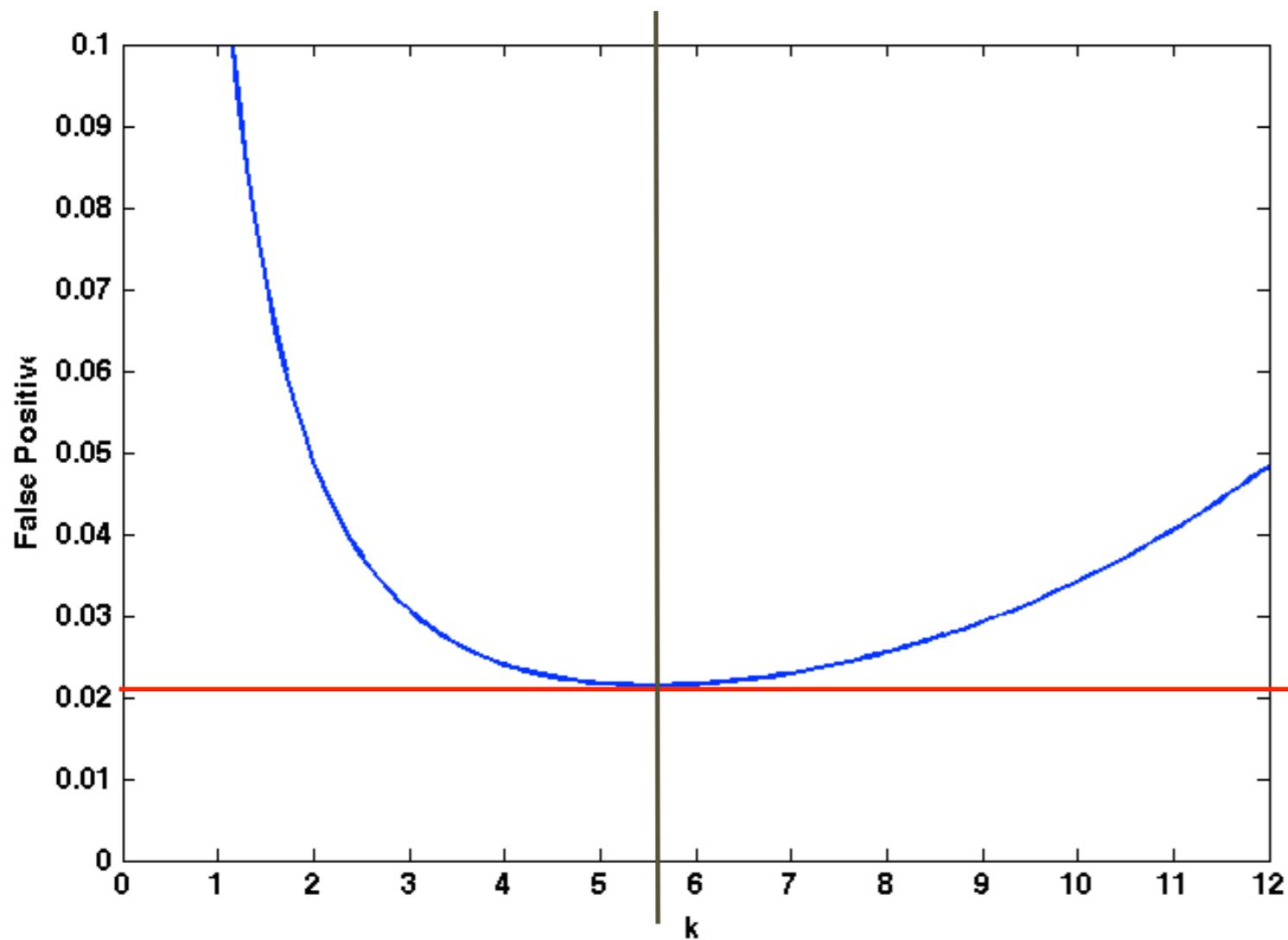
$A$  is initialized to all 0's;

for each  $x_i \in S$ : set  $A[h_j(x_i)] = 1$  for all  $1 \leq j \leq k$ ;

**Query**  $x$ : “yes” iff  $A[h_j(x)] = 1$  for all  $1 \leq j \leq k$

- $x \in S$ : always correct
- $x \notin S$ : **false positive**

$$\begin{aligned} & \Pr \left[ \forall 1 \leq j \leq k : A[h_j(x)] = 1 \right] \\ &= \left( \Pr \left[ A[h_j(x)] = 1 \right] \right)^k = \left( 1 - \Pr \left[ A[h_j(x)] = 0 \right] \right)^k \\ &\leq \left( 1 - (1 - 1/m)^{kn} \right)^k \approx \left( 1 - e^{-kn/m} \right)^k \end{aligned}$$



$y: x \in U$

$\dots, h_k : U \rightarrow [m]$

all  $1 \leq j \leq k;$

$\leq j \leq k$

- $x \notin S$ : false positive

$$\Pr \left[ \forall 1 \leq j \leq k : A[h_j(x)] = 1 \right]$$

choose  $k = c \ln 2$

$$m = cn$$

$$= \left( \Pr \left[ A[h_j(x)] = 1 \right] \right)^k = \left( 1 - \Pr \left[ A[h_j(x)] = 0 \right] \right)^k$$

$$\leq \left( 1 - (1 - 1/m)^{kn} \right)^k \approx \left( 1 - e^{-kn/m} \right)^k = 2^{-c \ln 2} \leq (0.6185)^c$$

# Bloom Filters (Bloom 1970)

**Data:** a set  $S$  of  $n$  items  $x_1, x_2, \dots, x_n \in U = [N]$

**Query:** an item  $x \in U$

Answer whether  $x \in S$  with bounded error.

- uniform & independent hash function  $h_1, \dots, h_k : U \rightarrow [m]$

**Data Structure:** bit array  $A \in \{0,1\}^m$

$A$  is initialized to all 0's;

for each  $x_i \in S$ : set  $A[h(x_i)] = 1$ ;

**Query**  $x$ : answer “yes” iff  $A[h(x)] = 1$

- choose  $k = c \ln 2$  and  $m = cn$ 
  - space cost:  $m = cn$  bits, time cost:  $k = c \ln 2$
  - false positive  $\leq (0.6185)^c$

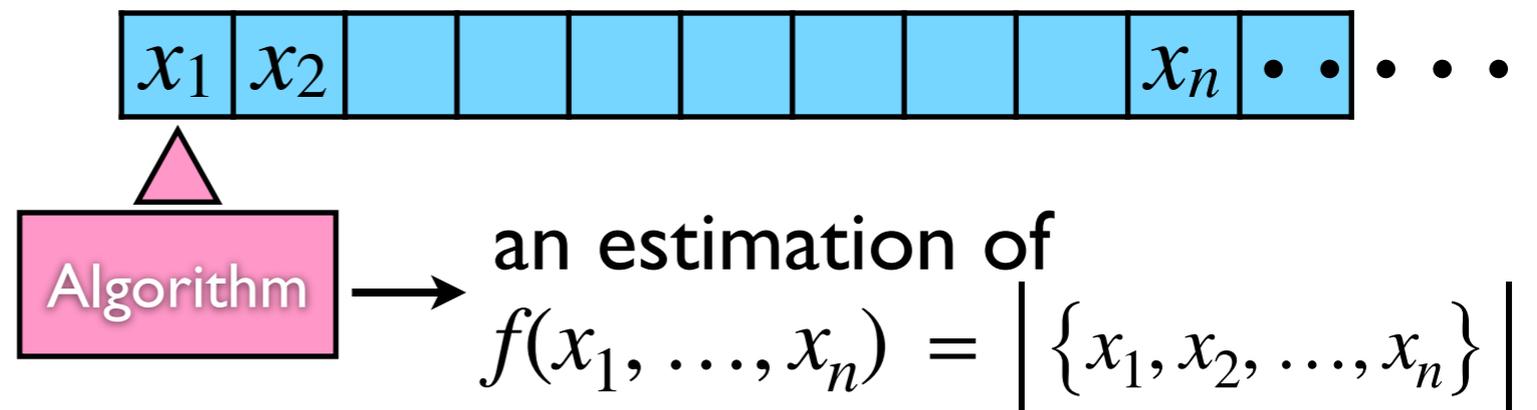
# **Distinct Elements** ***(Frequency Moments)***

# Count Distinct Elements

**Input:** a sequence  $x_1, x_2, \dots, x_n \in U = [N]$

**Output:** an estimation of  $z = \left| \{x_1, x_2, \dots, x_n\} \right|$

- **Data stream** model: input data item comes one at a time



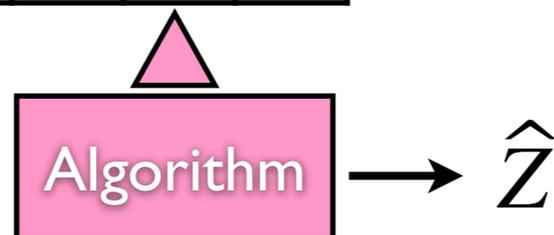
- Naïve algorithm: store all distinct data items using  $\Omega(z \log N)$  bits
- **Sketch:** (lossy) representation of data using space  $\ll z$
- **Lower bound (Alon-Matias-Szegedy):** any deterministic (exact or approx.) algorithm must use  $\Omega(N)$  bits of space in the worst case

# Count Distinct Elements

**Input:** a sequence  $x_1, x_2, \dots, x_n \in U = [N]$

**Output:** an estimation of  $z = \left| \{x_1, x_2, \dots, x_n\} \right|$

- **Data stream** model: input data item comes one at a time



- **$(\epsilon, \delta)$ -estimator:** randomized variable  $\hat{Z}$

$$\Pr \left[ (1 - \epsilon)z \leq \hat{Z} \leq (1 + \epsilon)z \right] \geq 1 - \delta$$

Using only memory equivalent to 5 lines of printed text, you can estimate with a typical accuracy of 5% and in a single pass the total vocabulary of Shakespeare.

**Input:** a sequence  $x_1, x_2, \dots, x_n \in U = [N]$

**Output:** an estimation of  $z = \left| \{x_1, x_2, \dots, x_n\} \right|$

### Simple Uniform Hash Assumption (SUHA):

A uniform function is available, whose preprocessing, representation and evaluation are considered to be easy.

- (idealized) uniform hash function  $h : U \rightarrow [0,1]$ 
  - $x_i = x_j \longrightarrow$  the same hash value  $h(x_i) = h(x_j) \in_r [0,1]$
  - $\{h(x_1), \dots, h(x_n)\}$ :  $z \times$  uniform and independent values in  $[0,1]$
  - partition  $[0,1]$  into  $z + 1$  subintervals (with *identically distributed* lengths)

$$\mathbb{E} \left[ \min_{1 \leq i \leq n} h(x_i) \right] = \mathbb{E}[\text{length of a subinterval}] = \frac{1}{z + 1} \quad (\text{by symmetry})$$

- estimator:  $\hat{Z} = \frac{1}{\min_i h(x_i)} - 1$  ? Variance is too large!

# Markov's Inequality

## Markov's Inequality

For *nonnegative* random variable  $X$ , for any  $t > 0$ ,

$$\Pr[X \geq t] \leq \frac{\mathbb{E}[X]}{t}$$

## Corollary

For random variable  $X$  and *nonnegative-valued* function  $f$ , for any  $t > 0$ ,

$$\Pr[f(X) \geq t] \leq \frac{\mathbb{E}[f(X)]}{t}$$

# Chebyshev's Inequality

## Chebyshev's Inequality

For random variable  $X$ , for any  $t > 0$ ,

$$\Pr \left[ |X - \mathbb{E}[X]| \geq t \right] \leq \frac{\mathbf{Var}[X]}{t^2}$$

- **Variance:**

$$\mathbf{Var}[X] = \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{E}[X^2] - (\mathbb{E}[X])^2$$

Apply Markov's inequality to  $Y = (X - \mathbb{E}[X])^2$ :

$$\Pr \left[ |X - \mathbb{E}[X]| \geq t \right] = \Pr[Y \geq t^2] \leq \frac{\mathbb{E}[Y]}{t^2} \leq \frac{\mathbf{Var}[X]}{t^2}$$

**Input:** a sequence  $x_1, x_2, \dots, x_n \in U = [N]$

**Output:** an estimation of  $z = \left| \{x_1, x_2, \dots, x_n\} \right|$

- (idealized) uniform hash function  $h : U \rightarrow [0,1]$

**Min Sketch:**

let  $Y = \min_{1 \leq i \leq n} h(x_i)$ ;

return  $\hat{Z} = \frac{1}{Y} - 1$ ;

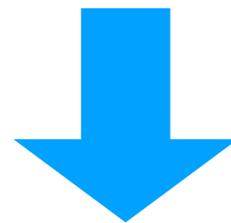
- By symmetry:

$$\mathbb{E}[Y] = \frac{1}{n+1}$$

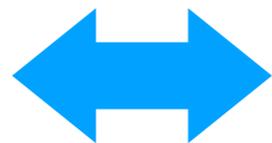
- Goal:

$$\Pr \left[ \hat{Z} < (1 - \epsilon)z \text{ or } \hat{Z} > (1 + \epsilon)z \right] \leq \delta$$

assuming  $\epsilon \leq 1/2$



$$\left| Y - \mathbb{E}[Y] \right| > \frac{\epsilon/2}{z+1}$$



$$\left| Y - \frac{1}{z+1} \right| > \frac{\epsilon/2}{z+1}$$

**Input:** a sequence  $x_1, x_2, \dots, x_n \in U = [N]$

**Output:** an estimation of  $z = \left| \{x_1, x_2, \dots, x_n\} \right|$

- (idealized) uniform hash function  $h : U \rightarrow [0,1]$

**Min Sketch:**

let  $Y = \min_{1 \leq i \leq n} h(x_i)$ ;

return  $\hat{Z} = \frac{1}{Y} - 1$ ;

- Uniform independent hash values:

$$H_1, \dots, H_z \in [0,1]$$



- $Y = \min_{1 \leq i \leq z} H_i$

**geometry  
probability:**

$$\Pr[Y > y] = (1 - y)^z \quad \longrightarrow \quad \text{pdf: } p(y) = z(1 - y)^{z-1}$$

$$\mathbb{E}[Y^2] = \int_0^1 y^2 p(y) dy = \int_0^1 y^2 z(1 - y)^{z-1} dy = \frac{2}{(z+1)(z+2)}$$

$$\mathbf{Var}[Y] = \mathbb{E}[Y^2] - \mathbb{E}[Y]^2 = \frac{z}{(z+1)^2(z+2)} \leq \frac{1}{(z+1)^2}$$

**Input:** a sequence  $x_1, x_2, \dots, x_n \in U = [N]$

**Output:** an estimation of  $z = \left| \{x_1, x_2, \dots, x_n\} \right|$

- (idealized) uniform hash function  $h : U \rightarrow [0,1]$

**Min Sketch:**

let  $Y = \min_{1 \leq i \leq n} h(x_i)$ ;

return  $\hat{Z} = \frac{1}{Y} - 1$ ;

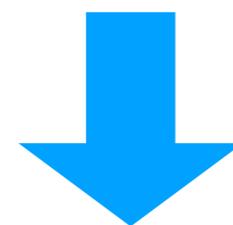
- By symmetry:

$$\mathbb{E}[Y] = \frac{1}{z+1}$$

- Goal:

$$\Pr \left[ \hat{Z} < (1 - \epsilon)z \text{ or } \hat{Z} > (1 + \epsilon)z \right] \leq \delta$$

assuming  $\epsilon \leq 1/2$



$$\text{Var}[Y] \leq \frac{1}{(z+1)^2} \xrightarrow{\text{(Chebyshev)}} \Pr \left[ |Y - \mathbb{E}[Y]| > \frac{\epsilon/2}{z+1} \right] \leq \frac{4}{\epsilon^2}$$

# The Mean Trick (for Variance Reduction)

- Variance and covariance:

$$\mathbf{Var}[X] = \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{E}[X^2] - (\mathbb{E}[X])^2$$

$$\mathbf{Cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]$$

- Useful properties:

$$\mathbf{Var}[X + a] = \mathbf{Var}[X]$$

$$\mathbf{Var}[aX] = a^2 \mathbf{Var}[X]$$

$$\mathbf{Var} \left[ \sum_i X_i \right] = \sum_i \mathbf{Var}[X_i] + \sum_{i \neq j} \mathbf{Cov}(X_i, X_j)$$

- For **pairwise independent** **identically distributed**  $X_i$ 's:

$$\mathbf{Var} \left[ \frac{1}{k} \sum_{i=1}^k X_i \right] = \frac{1}{k^2} \sum_{i=1}^k \mathbf{Var}[X_i] = \frac{1}{k} \mathbf{Var}[X_1]$$

**Input:** a sequence  $x_1, x_2, \dots, x_n \in U = [N]$

**Output:** an estimation of  $z = \left| \{x_1, x_2, \dots, x_n\} \right|$

- uniform & independent hash functions  $h_1, \dots, h_k : U \rightarrow [0,1]$

**Min Sketch:**

for each  $1 \leq j \leq k$ , let  $Y_j = \min_{1 \leq i \leq n} h_j(x_i)$ ;

return  $\hat{Z} = \frac{1}{\bar{Y}} - 1$  where  $\bar{Y} = \frac{1}{k} \sum_{j=1}^k Y_j$ ;

- For every  $1 \leq j \leq k$ :

$$\mathbb{E}[Y_j] = \frac{1}{z+1}$$

linearity of  
expectation



$$\mathbb{E}[\bar{Y}] = \frac{1}{z+1}$$

$$\text{Var}[Y_j] \leq \frac{1}{(z+1)^2}$$

independence



$$\text{Var}[\bar{Y}] \leq \frac{1}{k(z+1)^2}$$

**Input:** a sequence  $x_1, x_2, \dots, x_n \in U = [N]$

**Output:** an estimation of  $z = \left| \{x_1, x_2, \dots, x_n\} \right|$

- uniform & independent hash functions  $h_1, \dots, h_k : U \rightarrow [0,1]$

### Min Sketch:

for each  $1 \leq j \leq k$ , let  $Y_j = \min_{1 \leq i \leq n} h_j(x_i)$ ;

return  $\hat{Z} = \frac{1}{\bar{Y}} - 1$  where  $\bar{Y} = \frac{1}{k} \sum_{j=1}^k Y_j$ ;

$$\mathbb{E} [\bar{Y}] = \frac{1}{z + 1}$$

$$\mathbf{Var} [\bar{Y}] \leq \frac{1}{k(z + 1)^2}$$

- Goal:**  $\Pr \left[ \hat{Z} < (1 - \epsilon)z \text{ or } \hat{Z} > (1 + \epsilon)z \right] \leq \delta$



assuming  $\epsilon \leq 1/2$

$$\Pr \left[ \left| \bar{Y} - \mathbb{E} [\bar{Y}] \right| > \frac{\epsilon/2}{z + 1} \right] \leq \frac{4}{k\epsilon^2} \leq \delta$$

(Chebyshev)

$$\text{Set } k = \left\lceil \frac{4}{\epsilon^2 \delta} \right\rceil$$

**Input:** a sequence  $x_1, x_2, \dots, x_n \in U = [N]$

**Output:** an estimation of  $z = \left| \{x_1, x_2, \dots, x_n\} \right|$

- uniform & independent hash functions  $h_1, \dots, h_k : U \rightarrow [0,1]$

**Min Sketch:** set  $k = \lceil 4/(\epsilon^2\delta) \rceil$

for each  $1 \leq j \leq k$ , let  $Y_j = \min_{1 \leq i \leq n} h_j(x_i)$ ;

return  $\hat{Z} = \frac{1}{\bar{Y}} - 1$  where  $\bar{Y} = \frac{1}{k} \sum_{j=1}^k Y_j$ ;

$$\Pr \left[ (1 - \epsilon)z \leq \hat{Z} \leq (1 + \epsilon)z \right] \geq 1 - \delta$$

- **Space cost:**  $k = O\left(\frac{1}{\epsilon^2\delta}\right)$  *real numbers* in  $[0,1]$
- Storing  $k$  *idealized* hash functions.

# Universal Hashing

## Universal Hash Family (Carter and Wegman 1979):

A family  $\mathcal{H}$  of hash functions in  $U \rightarrow [m]$  is  **$k$ -universal** if for any distinct  $x_1, \dots, x_k \in U$ ,

$$\Pr_{h \in \mathcal{H}} [h(x_1) = \dots = h(x_k)] \leq \frac{1}{m^{k-1}}.$$

Moreover,  $\mathcal{H}$  is **strongly  $k$ -universal** ( $k$ -wise independent) if for any distinct  $x_1, \dots, x_k \in U$  and any  $y_1, \dots, y_k \in [m]$ ,

$$\Pr_{h \in \mathcal{H}} \left[ \bigwedge_{i=1}^k h(x_i) = y_i \right] = \frac{1}{m^k}.$$

# $k$ -Universal Hash Family

hash functions  $h : U \rightarrow [m]$

- **Linear congruential hashing:**

- Represent  $U \subseteq \mathbb{Z}_p$  for sufficiently large prime  $p$

- $h_{a,b}(x) = ((ax + b) \bmod p) \bmod m$

- $\mathcal{H} = \left\{ h_{a,b} \mid a \in \mathbb{Z}_p \setminus \{0\}, b \in \mathbb{Z}_p \right\}$

**Theorem:**

The linear congruential family  $\mathcal{H}$  is 2-wise independent.

- **Degree- $k$  polynomial in finite field with random coefficients**

- Hashing between binary fields:  $GF(2^w) \rightarrow GF(2^l)$

$$h_{a,b}(x) = (a * x + b) \gg (w-1)$$

# Flajolet-Martin Algorithm

**Input:** a sequence  $x_1, x_2, \dots, x_n \in [N] \subseteq [2^w]$

**Output:** an estimation of  $z = \left| \{x_1, x_2, \dots, x_n\} \right|$

- **2-wise independent** hash function  $h : [2^w] \rightarrow [2^w]$
- For  $y \in [2^w]$ , let **zeros**( $y$ ) =  $\max\{i : 2^i \mid y\}$  denote # of trailing 0's

**Flajolet-Martin Algorithm:**

let  $R = \max_{1 \leq i \leq n} \text{zeros}(h(x_i))$ ;

return  $\hat{Z} = 2^R$ ;

$$\Pr \left[ \hat{Z} < \frac{z}{C} \text{ or } \hat{Z} > C \cdot z \right] \leq \frac{3}{C}$$

**Input:** a sequence  $x_1, x_2, \dots, x_n \in [N] \subseteq [2^w]$

**Output:** an estimation of  $z = \left| \{x_1, x_2, \dots, x_n\} \right|$

- **2-wise independent** hash function  $h : [2^w] \rightarrow [2^w]$
- For  $y \in [2^w]$ , let **zeros**( $y$ ) =  $\max\{i : 2^i \mid y\}$  denote # of trailing 0's

**Flajolet-Martin Algorithm:**

let  $R = \max_{1 \leq i \leq n} \text{zeros}(h(x_i))$ ;

return  $\hat{Z} = 2^R$ ;

Let

$$Y_r = \sum_{x \in \{x_1, \dots, x_n\}} I[\text{zeros}(h(x)) \geq r]$$

(linearity of expectation)

$$\mathbb{E}[Y_r] = \sum_{x \in \{x_1, \dots, x_n\}} \Pr[\text{zeros}(h(x)) \geq r] = z2^{-r}$$

(pairwise independence)

$$\text{Var}[Y_r] = \sum_{x \in \{x_1, \dots, x_n\}} \text{Var}[I[\text{zeros}(h(x)) \geq r]] = z2^{-r}(1 - 2^{-r}) \leq z2^{-r}$$

# Pairwise Independent Trials

## Proposition:

If  $X$  is a sum of **pairwise independent** random variables taking values in  $\{0,1\}$ , then  $\mathbf{Var}[X] \leq \mathbb{E}[X]$ .

$$\begin{aligned}\mathbf{Var}[X] &= \sum_i \mathbf{Var}[X_i] = \sum_i (\mathbb{E}[X_i^2] - \mathbb{E}[X_i]^2) = \sum_i (\mathbb{E}[X_i] - \mathbb{E}[X_i]^2) \\ &= \mathbb{E}[X] - \sum_i \mathbb{E}[X_i]^2 \leq \mathbb{E}[X]\end{aligned}$$

## Corollary (Chebyshev's Inequality):

If  $X$  is a sum of **pairwise independent** random variables taking values in  $\{0,1\}$ , for any  $t > 0$ ,

$$\Pr \left[ |X - \mathbb{E}[X]| \geq t \right] \leq \frac{\mathbb{E}[X]}{t^2}$$

**Input:** a sequence  $x_1, x_2, \dots, x_n \in [N] \subseteq [2^w]$

**Output:** an estimation of  $z = \left| \{x_1, x_2, \dots, x_n\} \right|$

- **2-wise independent** hash function  $h : [2^w] \rightarrow [2^w]$
- For  $y \in [2^w]$ , let **zeros**( $y$ ) =  $\max\{i : 2^i \mid y\}$  denote # of trailing 0's

**Flajolet-Martin Algorithm:**

let  $R = \max_{1 \leq i \leq n} \text{zeros}(h(x_i))$ ;

return  $\hat{Z} = 2^R$ ;

Let

$$Y_r = \sum_{x \in \{x_1, \dots, x_n\}} I \left[ \text{zeros}(h(x)) \geq r \right]$$

(linearity of expectation)

$$\mathbb{E}[Y_r] = \sum_{x \in \{x_1, \dots, x_n\}} \Pr \left[ \text{zeros}(h(x)) \geq r \right] = z2^{-r}$$

(pairwise independence)  $\mathbf{Var}[Y_r] \leq \mathbb{E}[Y_r] \leq z2^{-r}$

**Input:** a sequence  $x_1, x_2, \dots, x_n \in [N] \subseteq [2^w]$

**Output:** an estimation of  $z = \left| \{x_1, x_2, \dots, x_n\} \right|$

- **2-wise independent** hash function  $h : [2^w] \rightarrow [2^w]$
- For  $y \in [2^w]$ , let **zeros**( $y$ ) =  $\max\{i : 2^i \mid y\}$  denote # of trailing 0's

**Flajolet-Martin Algorithm:**

let  $R = \max_{1 \leq i \leq n} \text{zeros}(h(x_i))$ ;

return  $\hat{Z} = 2^R$ ;

Let

$$Y_r = \sum_{x \in \{x_1, \dots, x_n\}} I[\text{zeros}(h(x)) \geq r]$$

$$\mathbb{E}[Y_r] = z2^{-r} \quad \mathbf{Var}[Y_r] \leq z2^{-r}$$

(denote  $r^* = \lceil \log_2 Cz \rceil$ )

(observe  $R = \max\{r : Y_r > 0\}$ )

(Markov's inequality)

$$\Pr[\hat{Z} > Cz] \leq \Pr[R \geq r^*]$$

$$\leq \Pr[Y_{r^*} > 0] = \Pr[Y_{r^*} \geq 1]$$

$$\leq \mathbb{E}[Y_{r^*}] = z/2^{r^*} \leq 1/C$$

**Input:** a sequence  $x_1, x_2, \dots, x_n \in [N] \subseteq [2^w]$

**Output:** an estimation of  $z = \left| \{x_1, x_2, \dots, x_n\} \right|$

- **2-wise independent** hash function  $h : [2^w] \rightarrow [2^w]$
- For  $y \in [2^w]$ , let **zeros**( $y$ ) =  $\max\{i : 2^i \mid y\}$  denote # of trailing 0's

**Flajolet-Martin Algorithm:**

let  $R = \max_{1 \leq i \leq n} \text{zeros}(h(x_i))$ ;

return  $\hat{Z} = 2^R$ ;

Let

$$Y_r = \sum_{x \in \{x_1, \dots, x_n\}} I[\text{zeros}(h(x)) \geq r]$$

$$\mathbb{E}[Y_r] = z2^{-r} \quad \mathbf{Var}[Y_r] \leq z2^{-r}$$

(denote  $r^{**} = \lceil \log_2(z/C) \rceil$ )

(observe  $R = \max\{r : Y_r > 0\}$ )

(Chebyshev's inequality)

$$\Pr[\hat{Z} < z/C] \leq \Pr[R < r^{**}]$$

$$\leq \Pr[Y_{r^{**}} = 0]$$

$$\leq \mathbf{Var}[Y_{r^{**}}] / \mathbb{E}[Y_{r^{**}}]^2 \leq 2^{r^{**}} / z$$

$$\leq 2/C$$

**Input:** a sequence  $x_1, x_2, \dots, x_n \in [N] \subseteq [2^w]$

**Output:** an estimation of  $z = \left| \{x_1, x_2, \dots, x_n\} \right|$

- **2-wise independent** hash function  $h : [2^w] \rightarrow [2^w]$
- For  $y \in [2^w]$ , let **zeros**( $y$ ) =  $\max\{i : 2^i \mid y\}$  denote # of trailing 0's

### Flajolet-Martin Algorithm:

let  $R = \max_{1 \leq i \leq n} \text{zeros}(h(x_i))$ ;

return  $\hat{Z} = 2^R$ ;

$$\Pr \left[ \hat{Z} < \frac{z}{C} \text{ or } \hat{Z} > C \cdot z \right] \leq \frac{3}{C}$$

- **Space cost:**  $O(\log \log N)$  bits for maintaining  $R$
- $O(\log N)$  bits for storing 2-wise independent hash function

# BJKST Algorithm

**Input:** a sequence  $x_1, x_2, \dots, x_n \in [N]$

**Output:** an estimation of  $z = \left| \{x_1, x_2, \dots, x_n\} \right|$

- **2-wise independent** hash function  $h : [N] \rightarrow [M] = \{1, \dots, M\}$

## BJKST Algorithm:

let  $Y_1, \dots, Y_k$  be the  $k$  smallest hash values among  
 $\{h(x_1), h(x_2), \dots, h(x_n)\}$ ;

return  $\hat{Z} = \frac{kM}{Y_k}$ ;

(Bar-Yossef, Jayram, Kumar, Sivakumar and Trevisan, 2002)

**Input:** a sequence  $x_1, x_2, \dots, x_n \in [N]$

**Output:** an estimation of  $z = \left| \{x_1, x_2, \dots, x_n\} \right|$

- **2-wise independent** hash function  $h : [N] \rightarrow [M] = \{1, \dots, M\}$

**BJKST Algorithm:**

let  $Y_1, \dots, Y_k$  be the  $k$  smallest hash values among  
 $\{h(x_1), h(x_2), \dots, h(x_n)\}$ ;

return  $\hat{Z} = \frac{kM}{Y_k}$ ;

- **Goal:**  $\Pr \left[ \hat{Z} < (1 - \epsilon)z \text{ or } \hat{Z} > (1 + \epsilon)z \right] \leq \delta$

assuming  $\epsilon \leq 1$



$$\left| Y_k - \frac{kM}{z} \right| > \frac{\epsilon}{2} \cdot \frac{kM}{z}$$

- uniform and **2-wise independent**  $X_1, \dots, X_n \in [N^3]$
- let  $Y_1, \dots, Y_z$  be these elements in non-decreasing order

$$\text{Let } V = \sum_{i=1}^z I \left[ X_i \leq \left(1 - \frac{\epsilon}{2}\right) \frac{kM}{z} \right] \quad W = \sum_{i=1}^z I \left[ X_i \leq \left(1 + \frac{\epsilon}{2}\right) \frac{kM}{z} \right]$$

$$\mathbb{E}[V] = \left(1 - \frac{\epsilon}{2} \pm o(1)\right) k \quad \mathbb{E}[W] = \left(1 + \frac{\epsilon}{2} \pm o(1)\right) k$$

$$Y_k < \left(1 - \frac{\epsilon}{2}\right) \frac{k(M+1)}{z} \implies V \geq k \quad Y_k > \left(1 + \frac{\epsilon}{2}\right) \frac{k(M+1)}{z} \implies W \leq k$$

(Chebyshev's inequality for sum of pairwise independent trials)

$$\Pr[V \geq k] \leq \frac{8}{k\epsilon^2}$$

$$\Pr[W \leq k] \leq \frac{8}{k\epsilon^2}$$

- **Goal:**  $\Pr \left[ \left| Y_k - \frac{kM}{z} \right| > \frac{\epsilon}{2} \cdot \frac{kM}{z} \right] \leq \delta$       Set  $k = \left\lceil \frac{16}{\epsilon^2 \delta} \right\rceil$

**Input:** a sequence  $x_1, x_2, \dots, x_n \in [N]$

**Output:** an estimation of  $z = \left| \{x_1, x_2, \dots, x_n\} \right|$

- **2-wise independent** hash function  $h : [N] \rightarrow [N^3]$

**BJKST Algorithm:** Set  $k = \lceil 16/(\epsilon^2 \delta) \rceil$

let  $Y_1, \dots, Y_k$  be the  $k$  smallest hash values among  
 $\{h(x_1), h(x_2), \dots, h(x_n)\}$ ;

return  $\hat{Z} = \frac{kM}{Y_k}$ ;

$$\Pr \left[ (1 - \epsilon)z \leq \hat{Z} \leq (1 + \epsilon)z \right] \geq 1 - \delta$$

- **Space cost:**  $O(k \log N) = O(\epsilon^{-2} \log N)$  bits when  $\delta = \Omega(1)$

# Frequency Moments

- **Data stream:**  $x_1, x_2, \dots, x_n \in U$
- for each  $x \in U$ , define **frequency** of  $x$  as  $f_x = |\{i : x_i = x\}|$   
 **$k$ -th frequency moments:**  $F_k = \sum_{x \in U} f_x^k$
- **Space complexity** for  $(\epsilon, \delta)$ -estimation: constant  $\epsilon, \delta$ 
  - for  $k \leq 2$ :  $\text{polylog}(N)$  [Alon-Matias-Szegedy '96]
  - for  $k > 2$ :  $\tilde{\Theta}(N^{1-2/k})$  [Indyk-Woodruff '05]
- **Count distinct elements:**  $F_0$ 
  - optimal algorithm [Kane-Nelson-Woodruff '10]:  $O(\epsilon^{-2} + \log N)$  bits

# Frequency Estimation

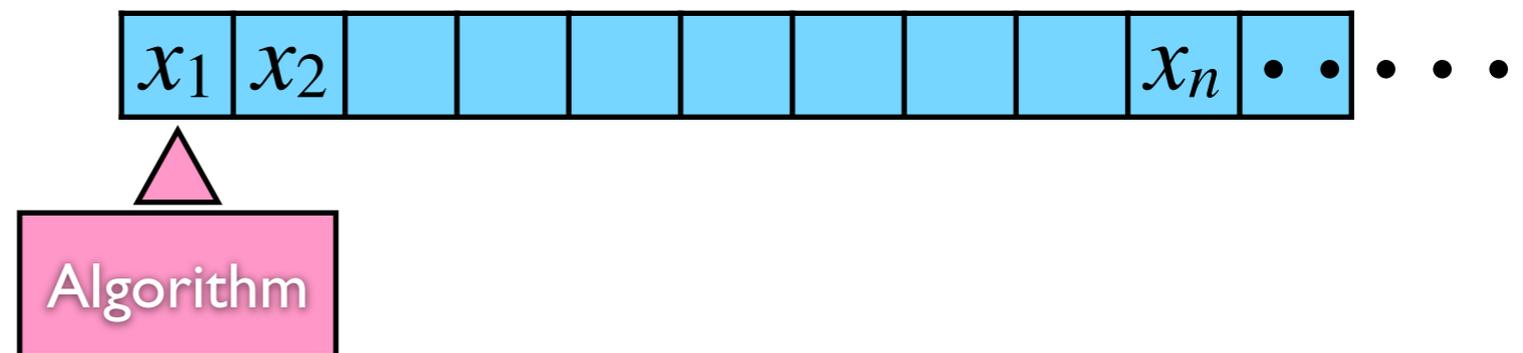
# Frequency Estimation

**Data:** a sequence  $x_1, x_2, \dots, x_n \in U = [N]$

**Query:** an item  $x \in U$

Estimate the *frequency*  $f_x = |\{i : x_i = x\}|$  of  $x$ .

- **Data stream** model: input data item comes one at a time



# Frequency Estimation

**Data:** a sequence  $x_1, x_2, \dots, x_n \in U = [N]$

**Query:** an item  $x \in U$

Estimate the **frequency**  $f_x = |\{i : x_i = x\}|$  of  $x$ .

- **Data stream** model: input data item comes one at a time



$\hat{f}_x$  : estimation of  $f_x$   
within **additive error**

$$\Pr \left[ |\hat{f}_x - f_x| \geq \epsilon n \right] \leq \delta$$

- **Heavy hitters:** items that appears  $> \epsilon n$  times

# Count-Min Sketch

**Data:** a sequence  $x_1, x_2, \dots, x_n \in U = [N]$

**Query:** an item  $x \in U$

Estimate the *frequency*  $f_x = |\{i : x_i = x\}|$  of  $x$ .

- $k$  independent **2-universal** hash functions  $h_1, \dots, h_k : [N] \rightarrow [m]$

**Count-Min Sketch:** CMS[ $k$ ][ $m$ ] (initialized to all 0's)

**Upon** each  $x_i$ : CMS[ $j$ ][ $h_j(x_i)$ ] ++ for all  $1 \leq j \leq k$ ;

**Query**  $x$ : return  $\hat{f}_x = \min_{1 \leq j \leq k} \text{CMS}[j][h_j(x)]$

**Observation:** CMS[ $j$ ][ $h_j(x)$ ]  $\geq f_x$  for all  $1 \leq j \leq k$

$$f_x \leq \hat{f}_x \leq ?$$

**Data:** sequence  $x_1, \dots, x_n \in [N]$     **Query:**  $x \in [N]$

*frequency*  $f_x = |\{i : x_i = x\}|$  of  $x$

- $k$  independent **2-universal** hash functions  $h_1, \dots, h_k : [N] \rightarrow [m]$

**Count-Min Sketch:** CMS[ $k$ ][ $m$ ] (initialized to all 0's)

**Upon** each  $x_i$ : CMS[ $j$ ][ $h_j(x_i)$ ] + + for all  $1 \leq j \leq k$ ;

**Query**  $x$ : return  $\hat{f}_x = \min_{1 \leq j \leq k} \text{CMS}[j][h_j(x)]$

- for any  $x \in [N]$  and any  $1 \leq j \leq k$ :

$$\text{CMS}[j][h_j(x)] = f_x + \sum_{\substack{y \in \{x_1, \dots, x_n\} \setminus \{x\} \\ h_j(y) = h_j(x)}} f_y$$

$$\mathbb{E} \left[ \text{CMS}[j][h_j(x)] \right] = f_x + \sum_{y \in \{x_1, \dots, x_n\} \setminus \{x\}} f_y \Pr[h_j(y) = h_j(x)]$$

**Data:** sequence  $x_1, \dots, x_n \in [N]$     **Query:**  $x \in [N]$

**frequency**  $f_x = |\{i : x_i = x\}|$  of  $x$

- $k$  independent **2-universal** hash functions  $h_1, \dots, h_k : [N] \rightarrow [m]$

**Count-Min Sketch:** CMS[ $k$ ][ $m$ ] (initialized to all 0's)

**Upon** each  $x_i$ : CMS[ $j$ ][ $h_j(x_i)$ ] ++ for all  $1 \leq j \leq k$ ;

**Query**  $x$ : return  $\hat{f}_x = \min_{1 \leq j \leq k} \text{CMS}[j][h_j(x)]$

- for any  $x \in [N]$  and any  $1 \leq j \leq k$ :

$$\begin{aligned} \mathbb{E} \left[ \text{CMS}[j][h_j(x)] \right] &= f_x + \sum_{y \in \{x_1, \dots, x_n\} \setminus \{x\}} f_y \Pr[h_j(y) = h_j(x)] \\ &\leq f_x + \frac{1}{m} \sum_{y \in \{x_1, \dots, x_n\} \setminus \{x\}} f_y \leq f_x + \frac{1}{m} \sum_{y \in \{x_1, \dots, x_n\}} f_y = f_x + \frac{n}{m} \end{aligned}$$

**Data:** sequence  $x_1, \dots, x_n \in [N]$     **Query:**  $x \in [N]$

**frequency**  $f_x = |\{i : x_i = x\}|$  of  $x$

- $k$  independent **2-universal** hash functions  $h_1, \dots, h_k : [N] \rightarrow [m]$

**Count-Min Sketch:** CMS[ $k$ ][ $m$ ] (initialized to all 0's)

**Upon** each  $x_i$ : CMS[ $j$ ][ $h_j(x_i)$ ] ++ for all  $1 \leq j \leq k$ ;

**Query**  $x$ : return  $\hat{f}_x = \min_{1 \leq j \leq k} \text{CMS}[j][h_j(x)]$

$$\forall x, \forall j: \quad \text{CMS}[j][h_j(x)] \geq f_x$$
$$\mathbb{E} \left[ \text{CMS}[j][h_j(x)] \right] \leq f_x + \frac{n}{m}$$

**(Markov's inequality)**  $\Pr \left[ \text{CMS}[j][h_j(x)] - f_x \geq \epsilon n \right] \leq \frac{1}{\epsilon m}$

$$\Pr \left[ |\hat{f}_x - f_x| \geq \epsilon n \right] = \Pr \left[ \forall 1 \leq j \leq k : \text{CMS}[j][h_j(x)] - f_x \geq \epsilon n \right] \leq \left( \frac{1}{\epsilon m} \right)^k$$

**Data:** a sequence  $x_1, x_2, \dots, x_n \in U = [N]$

**Query:** an item  $x \in U$

Estimate the **frequency**  $f_x = |\{i : x_i = x\}|$  of  $x$ .

- $k$  independent **2-universal** hash functions  $h_1, \dots, h_k : [N] \rightarrow [m]$

**Count-Min Sketch:** CMS[ $k$ ][ $m$ ] (initialized to all 0's)

**Upon** each  $x_i$ : CMS[ $j$ ][ $h_j(x_i)$ ] ++ for all  $1 \leq j \leq k$ ;

**Query**  $x$ : return  $\hat{f}_x = \min_{1 \leq j \leq k} \text{CMS}[j][h_j(x)]$

$$\Pr \left[ |\hat{f}_x - f_x| \geq \epsilon n \right] \leq \left( \frac{1}{\epsilon m} \right)^k \leq \delta$$

- Choose  $m = \lceil e/\epsilon \rceil$  and  $k = \lceil \ln(1/\delta) \rceil$ 
  - **space cost:**  $O\left(\frac{1}{\epsilon} \log(1/\delta) \log n\right)$  bits
  - $O(\log(1/\delta) \log N)$  bits for hash functions
  - **time cost for query:**  $O(\log(1/\delta))$