# Hashing

hash function $h: \mathcal{U} \to [m]$    ← $\{0, 1, \dots m-1\}$

    drawn at random from a set $\mathcal{H}$ of functions

$h(s, x)$

    ↑ ↑ item to be hashed

    salt — fixed when table created, random

$\mathcal{H}$ is <u>universal</u> if $\Pr\limits_{h \in \mathcal{H}}[h(x) = h(y)] \leq \frac{1}{m}$

                for all $x \neq y$



Chained hashing

$$E\left[\substack{\text{time to search} \\ \text{for item } \underline{not} \text{ in } T}\right] = O\left(\frac{n}{m} + 1\right)$$

    #items in table

    size of main array

$\frac{n}{m}$ = <u>load factor</u>

If $m = \Theta(n)$, $E[T(x)] = \Theta(1)$ for all $x$

But we'd like $E[\max\limits_{x} T(x)] = O(1)$

In fact, if $n = m$, <u>ideal</u> <u>random</u> hashing

    $E[\max \text{ chain length}] = \Theta\left(\frac{\log n}{\log \log n}\right)$

---

"Perfect" hashing     Komlos Szemeredi _____

    Replace linked lists with secondary

                       hash tables

$m = n$

Secondary Table of size $m_i = n_i^2$

$n_i = \#\{x \in T \mid h(x) = i\}$

**Assuming universal hashing**

Lemma: If $m = n^2$, $E[\#\text{collisions}] < 1$

$$E[\#\text{collisions}] = \sum_{x \neq y} Pr[h(x) = h(y)]$$

$$\leq \binom{n}{2} \cdot \frac{1}{m} = \frac{n(n-1)}{2} \cdot \frac{1}{n^2}$$

$$< \frac{1}{2}$$

Markov $\Rightarrow Pr[\geq 1 \text{ collision}] < \frac{1}{2}$

So after $\leq 2$ tries, <u>no</u> collisions!

Lookup(x):
$\quad i \leftarrow h(x)$ $\longleftarrow$ primary hash function
$\quad j \leftarrow h_i(x)$ $\longleftarrow$ secondary associated with $T[i]$
$\quad$ if $T_i[j] = x$
$\quad\quad$ return True
$\quad$ else
$\quad\quad$ return FALSE

$E[\text{total space}]$

$$= E\left[\sum_i n_i^2\right] = \sum_i E[n_i^2]$$

$$= \sum_i E\left[\left(\sum_x [h(x)=i]\right)^2\right]$$

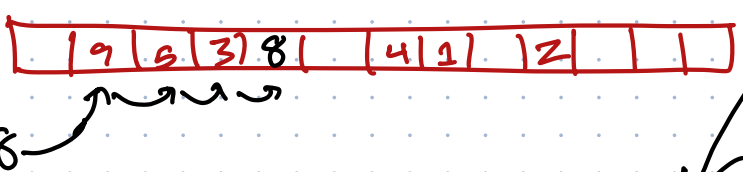$$= \sum_i E\left[\left(\sum_x [h(x)=i]\right) \cdot \left(\sum_y [h(y)=i]\right)\right]$$

$$= \sum_i E\left[\sum_x [h(x)=i] + 2\sum_{x<y} [h(x)=h(y)=i]\right]$$

$$= n + E\left[2 \sum_{x<y} [h(x)=h(y)]\right]$$

$$\leq n + 2\binom{n}{2}\frac{1}{n} = 2n-1 = \boxed{O(n)}$$

---

Bad cache behavior — bad locality

Better: open adressing — linear probing

| | 9 | 6 | 3 | 8 | | 4 | 1 | | 2 | | | |

Everything is in main table

8

Lookup(x):

$i \leftarrow h(x)$
while ($T[i] \neq$ Null and $T[i] \neq x$)
  $i \leftarrow i+1 \bmod m$

Binary probing
  Lookup(x):

  $i \leftarrow h(x)$
  for $j \leftarrow 0$ to $m-1$
    if $T[i \oplus j] =$ Null
      False
    else if $T[i \oplus j] = x$
      True

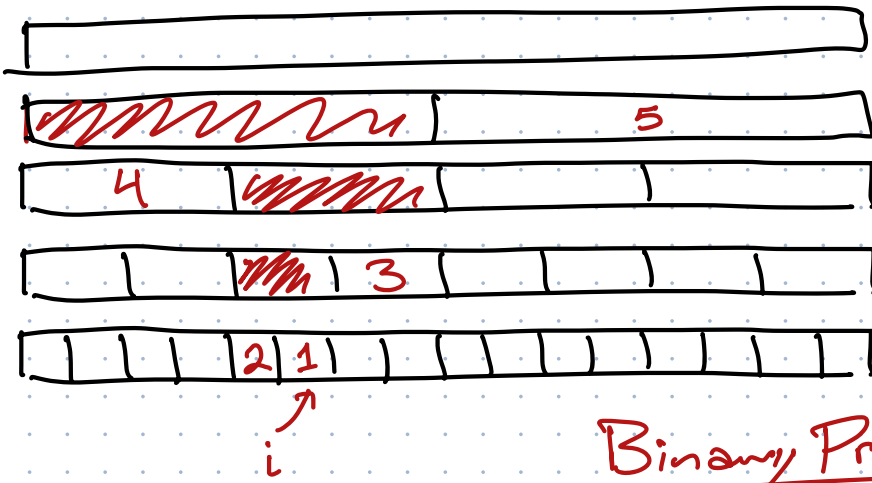IF $H$ is 3-uniform    $E[T(x)] = O(\log n)$
        5-uniform      $E[T(x)] = O(1)$

random matrix
tabulated

$h(x) = a + bx + cx^2 + dx^3 + ex^4$
              $\bmod m$

Twisted tabulation
$h(x,y) = A[x] \oplus B[y] \oplus C[x+y]$

$$m = 4n$$

## Binary Probing:
probe $T[i]$
for $\ell \leftarrow 0$ to $\lg n - 1$
  probe sibling of level-$\ell$
  block containing $T[i]$

size $2^\ell \longrightarrow$

$$E[\text{Time}] \leq \sum_\ell 2^\ell \cdot \Pr\left(\text{level-}\ell \text{ block containing } T[i] \text{ is } \underline{full}\right)$$

$$\Updownarrow$$

$$\Pr\left[\geq 2^\ell \text{ items hashed into block at level } \ell\right]$$

$$\|$$

$$\Pr\left[N_\ell \geq \textcircled{4} E[N_\ell]\right]$$

$\uparrow$ #items hashing into level-$\ell$ block

$m = 4n \Rightarrow$

$E[\text{# hashing into } B_\ell] = \boxed{2^\ell/4}$

Assuming uniformity

Chebyshev's inequality
$X$ is sum of pairwise
independent
0/1 vars.

$$\Pr\left[X \geq \textcircled{(1+\delta)}\mu\right] \leq \frac{1}{\delta\mu}$$

$$\frac{1}{3 \cdot \frac{2^\ell}{4}} = \frac{4}{3 \cdot 2^\ell}$$

$$\Rightarrow E[\text{time}] \leq \sum_\ell \frac{4 \cdot 2^\ell}{3 \cdot 2^\ell} = O(\log n)$$

4-way independence

$$\Pr\left[X \geq (1+\delta)\mu\right] = O\left(\frac{1}{(\delta\mu)^2}\right)$$

$$\Rightarrow E[\text{time}] \leq \sum_\ell 2^\ell O\left(\frac{1}{4^\ell}\right) = \sum_\ell O(2^{-\ell})$$
$$= O(1)$$