# Hash Tables

Subset of Universe $\mathcal{U}$ $\longrightarrow$ store in an array of size $m$
$= \{0 .. 2^w - 1\}$

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

Ideally $h(x) \neq h(y)$ for all $x, y$ in input

Fiction: $O(1)$ time?

$$\boxed{h(x) = \cancel{x \bmod m}}$$

Knuth:
$$\left( h(x) = \cancel{\lfloor mx \phi \rfloor \bmod m} \right)$$

Deterministic hash function
guarantees predictable collisions.

OTOH, perfect randomness is also useless

---

Fix a set $\mathcal{H}$ of hash functions in advance ("family")
when we create a hash table, pick $h \in \mathcal{H}$ at random
Use $h$ for the life time of the table.

Choose parameters
called "salt"

---

Properties we want:

- ~~Uniform:~~ $\Pr_{h \in \mathcal{H}}[h(x) = i] = 1/m$

  $h_0(x) = 0$ for all $x$
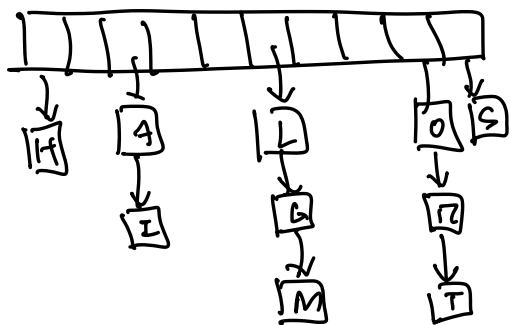  $h_1(x) = 1$ for all $x$
  $\vdots$
  $h_{m-1}(x) = m-1$ for all $x$

  $\{h_0, h_1 \ldots h_{m-1}\}$ is uniform

- Near Universality: $\Pr_{h \in \mathcal{H}}[h(x) = h(y)] \leq \frac{O(1)}{m}$ for all $x \neq y$



Chained hash table

Resolve collisions by
storing a list at every $T[i]$

Expected time to look up $x$ is
$$\leq O(1 + E[\ell(x)])$$

$$= O\left(1 + E[\#\, y \text{ s.t. } h(x)=h(y)]\right)$$

$$= O\left(1 + \sum_{y} Pr[h(x)=h(y)]\right)$$

$$= \boxed{O\left(1 + n/m\right)}$$

$$\boxed{\text{load factor } \alpha = \frac{n}{m}}$$

---

[Carter Wegman 1969]

① Multiplicative    choose prime $p > |\mathcal{U}|$

$[p] = \{0 \dots p-1\}$      $[p]^+ = \{1 \dots p-1\}$

Choose salt $a \in [p]^+$ uniformly at random

$$\boxed{h_a(x) = (ax \bmod p) \bmod m}$$

Near-universal     $Pr[h(x)=h(y)] \leq \frac{2}{m}$

② Multiply-add

Choose $a \in [p]^+$   $b \in [p]$

$$\boxed{h_{a,b}(x) = (ax+b \bmod p) \bmod m}$$
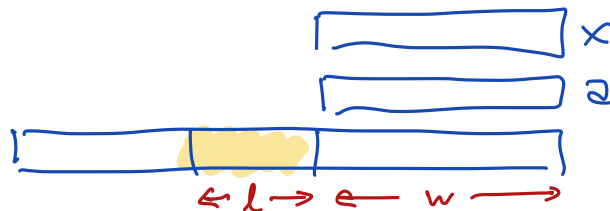
universal    uniform    2-uniform

③ Binary multiplication        $\mathcal{U} = \{0 \dots 2^w - 1\}$

Salt: $a \in [2^w]$        $m = 2^\ell$

$$\boxed{h_a(x) = \left\lfloor \frac{(a \cdot x) \bmod 2^w}{2^{w-\ell}} \right\rfloor}$$



$\leftarrow \ell \rightarrow \leftarrow w \longrightarrow$

$$\boxed{((a) * (x)) >> (\text{WORDSIZE} - \text{HASHBITS})}$$

④ Tabulation hashing $|n| = 2^w$  $m = 2^\ell$
$$= 2^{w/2} \times 2^{w/2}$$

Define two random arrays

$$A[0..2^{w/2}-1] \qquad B[0..2^{w/2}-1]$$

filled with random $\ell$-bit labels

$$\boxed{h_{A,B}(x,y) = A[x] \oplus B[y]}$$

<span style="color:red">universal    2-uniform     3-uniform     not 4-uniform</span>

⑤    Let M be a random matrix



$$\cdot \begin{bmatrix} \\ \\ \end{bmatrix} = \begin{bmatrix} \\ \end{bmatrix} \quad \text{mod } 2$$

<span style="color:red">near-universal</span>

---

For each $x$ we have $E\{\ell(x)\} = O(1)$

We want $E(\max_x \ell(x)) = O(1)$    <span style="color:red"><u>TOO BAD.</u></span>

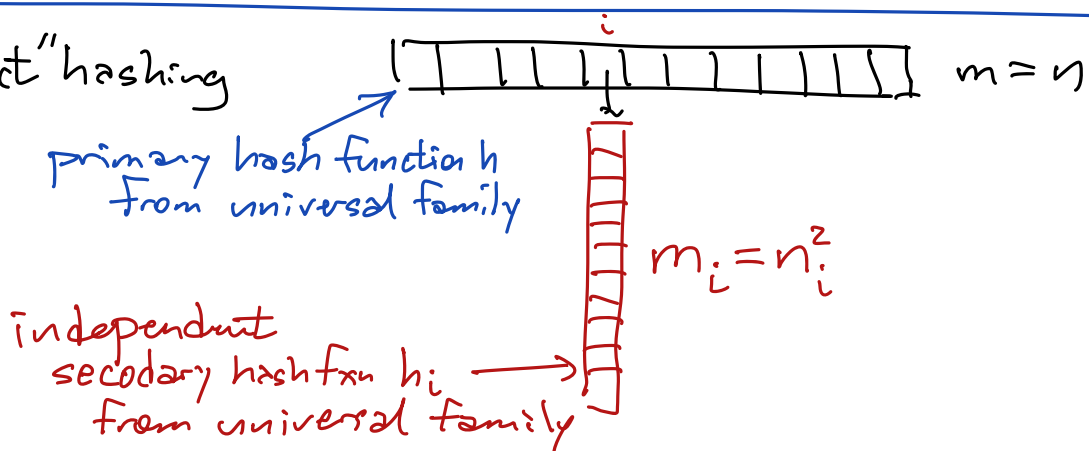<span style="color:red">Ideal random hashing     $m = n$</span>

<span style="color:red">$\Rightarrow \max_x \ell(x) = \Theta\left(\dfrac{\log n}{\log \log n}\right)$   whp</span>

<span style="color:blue"><u>$m = n^2$</u>    $E[\#\text{collisions}] \le \binom{n}{2}\dfrac{1}{m} < \dfrac{1}{2}$</span>

<span style="color:blue">$Pr(\underline{\text{any}} \text{ collisions}) < \frac{1}{2}$</span>

---

"Perfect" hashing  $m = n$

<span style="color:blue">primary hash function h
from universal family</span>

<span style="color:red">$m_i = n_i^2$</span>

<span style="color:red">independent
secondary hash fxn $h_i$ ⟶
from universal family</span>

Lookup(x):
  $i \leftarrow h(x)$
  $j \leftarrow h_i(x)$
  return $H[i][j]$

$$E[\text{Space}] = n + \sum_{i=1}^{n} E[n_i^2]$$

$$E[n_i^2] = E\left[\sum_{x,y=1}^{n} [h(x)=i][h(y)=i]\right]$$

$$= E\left[\sum_{x=1}^{n} [h(x)=i]^2 + 2\sum_{x<y} [h(x)=i=h(y)]\right]$$

$$= 1 \quad + 2 \cdot E\left[\sum_{x<y} [h(x)=h(y)=i]\right]$$

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

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

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

$$\leq n + 2\sum_{x<y} \frac{1}{n} \quad = n + 2\binom{n}{2}\frac{1}{n}$$

$$\leq 3n$$