

HASHING

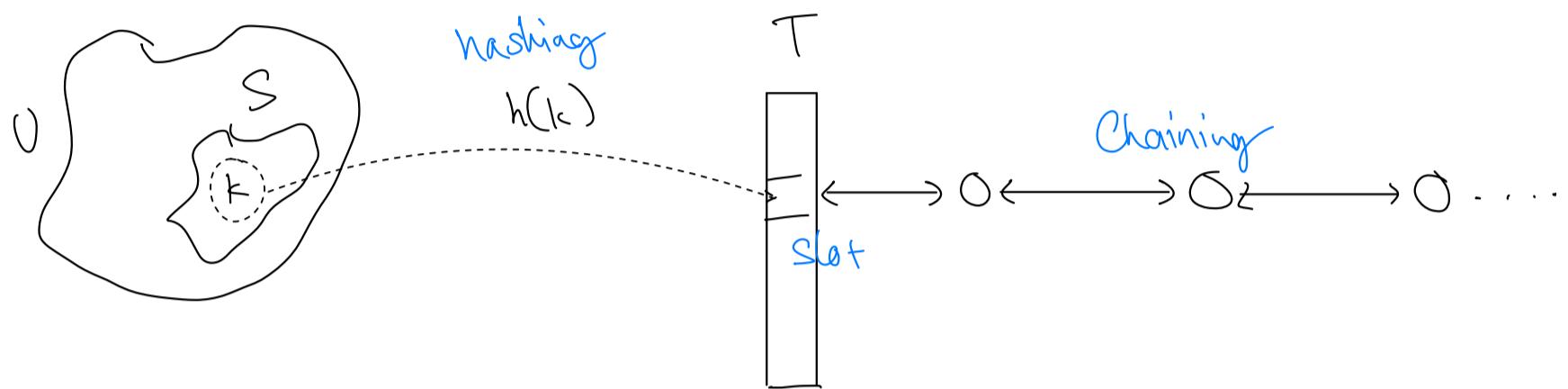
Hashing with Chaining

Universal set: $U = \{ \dots \dots \}^\infty$

Key set $S \subseteq U$, stored in a hash table $T[0 \dots m-1]$

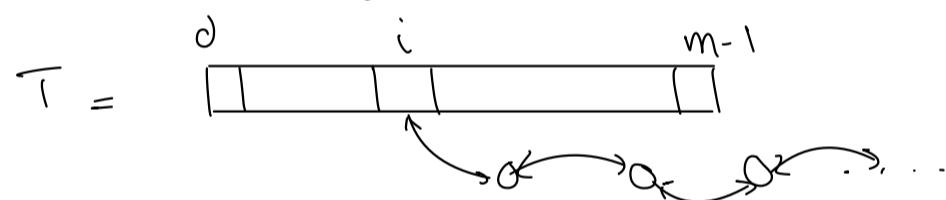
Hash func: $h(k): U \rightarrow \{0, \dots, m-1\}$
where $k \in S$

\Rightarrow key k is hashed into slot $h(k)$



Collision: Many keys k in S hashed into the same slot

\hookrightarrow Solution: (Doubly) Linked-list



Load factor: average number of keys per slot

Complexity:

- . Insert and Delete: $\Theta(1)$
- . Search: $\Theta(n)$ if chain has size $\Theta(n)$

Challenge: Can we design a good hashing strategy that distributes keys evenly across the slots?

\hookrightarrow No. If you hash a set of n keys, assuming $h: U \rightarrow \{0 \dots m-1\}$ and $|U| \geq m \cdot n$, then exist n keys that can be hashed to the same slot

Simple Uniform Hashing

Simple Uniform Hashing Properties

- Any element in U is equally likely to hash into any m slots in T
- $\Leftrightarrow P(h(k) = q) = \frac{1}{m} \quad \forall \text{ key } k \in S \subseteq U, \forall \text{ slot } q$

Lemma

Hashing $S \subseteq U$, $|S|=n$, onto hashtable $T[0..m-1]$ using function

$h: U \rightarrow \{0..m-1\}$ with Simple Uniform Hashing

Let L_k be the number of keys in slot $h(k)$ in T

$$E[L_k] < \alpha + 1 \quad : \alpha = n \cdot \frac{1}{m}$$

Proof:

Let X be number of keys in slot $h(k)$ that is not k .

- Case 1: k is not in T

$$\begin{aligned} \Leftrightarrow X = L_k \\ \text{Since } L_k \sim \text{Bin}(n, \frac{1}{m}) \end{aligned} \quad \Rightarrow X \sim \text{Bin}(n, \frac{1}{m})$$

$$\Rightarrow E[X] = n \cdot \frac{1}{m} = \alpha$$

$$< \alpha + 1$$

- Case 2: k is in T

$$\begin{aligned} \Leftrightarrow L_k = X + 1 \\ \text{Since } L_k \sim \text{Bin}(n, \frac{1}{m}) \end{aligned} \quad \Rightarrow X \sim \text{Bin}(n-1, \frac{1}{m})$$

$$\Rightarrow E[X] = (n-1) \cdot \frac{1}{m} = \alpha - 1$$

$$< \alpha + 1$$

\Rightarrow Proven

Corollary:

The expected time to search for a key under Simple Uniform Hashing is: $\Theta(1 + \alpha)$ $\alpha = \frac{n}{m}$

Problem: In practice, we can't implement Simple Uniform Hashing

\hookrightarrow Solution:

Universal Hashing, provably good expected search time for hashing with chaining

Universal Hashing

Randomized Algorithm : an extra call to Random (a, b) , which will return a number in range $[a, b]$ with (almost) equal probability

$$\Rightarrow E[\text{running-time}] = \frac{1}{n} \sum_{\text{n-outcomes}} \text{running-time}$$

n = outcomes generated by Random (a, b)

Definition (Universal Hashing)

Let $\mathcal{H} = \{h_1, h_2, \dots, h_r\}$ be a set of hash functions size r .

If we can construct set \mathcal{H} s.t. :

$$|\{h_i : h_i(x) = h_i(y)\}| \leq \gamma_m \quad \forall \text{ distinct } x, y$$

collision

$$\text{Then: } P[h_i(x) = h_i(y)] \leq \gamma_m \cdot \frac{1}{r} = \frac{1}{m}$$

In words: If we carefully craft a set of \mathcal{H} hash functions s.t. the number of functions resulted in collision is lesser than γ_m , then the probability of collision is $\frac{1}{m}$ for any $h \in \mathcal{H}$

Lemma:

Let h be random hash function from universal family \mathcal{H} be used to hash n keys into table $T[0 \dots m-1]$

Let L_k be number of keys in slot $h(k)$ in T , $\forall k \in U$.

$$E[L_k] < \alpha + 1$$

Proof: (that $E[L_k] < \alpha + 1 \quad \forall h \in \mathcal{H}$)

$$\left\{ \begin{array}{l} X_{kl} \text{ be i.r.v. s.t. } \\ Y = \sum_{l \neq k} X_{kl} \end{array} \right\} \left\{ \begin{array}{ll} X_{kl} = 1 & \text{if } h(k) = h(l) \\ X_{kl} = 0 & \text{otherwise } (h(k) \neq h(l)) \end{array} \right\} \begin{array}{l} X_{kl} = 1 \quad \text{if } h(k) = h(l) \quad (\text{collision}) \\ X_{kl} = 0 \quad \text{otherwise } (h(k) \neq h(l)) \end{array}$$

• Case 1: k not in T

$$\Rightarrow Y = L_k$$

$$\Rightarrow E[L_k] = E[Y] = E\left[\sum_{l \neq k} X_{kl}\right]$$

$$= \sum_{l \neq k} E[X_{kl}] \quad (\text{linearity of expectation})$$

$$= \sum_{l \neq k} \left(0 \cdot P(X_{kl}=0) + 1 \cdot P(X_{kl}=1) \right)$$

$$= \sum_{l \neq k} P(X_{kl}=1) \quad (\text{collision probability})$$

$$\leq \sum_{l \neq k} \frac{1}{m} = n \cdot \frac{1}{m}$$

$$= \alpha$$

• Case 2: k is in T

$$\Rightarrow L_k = Y + 1$$

$$\Rightarrow E[L_k] = E[Y] + 1$$

$$\text{Since } E[Y] \leq (n-1) \frac{1}{m} \quad (\text{similar reasoning to Case 1})$$

$$< \alpha$$

$$\Rightarrow E[L_k] < \alpha + 1$$

\Rightarrow Prove

Universal Family $H_{p,m}$

How do we actually construct such set f satisfies the properties of Universal Hashing?

Example 1: Using prime number

• Choose a prime p s.t. $|U| = p$ (or $p > |U|$)

• Define hash function $h_{a,b}$ as:

$$h_{a,b}(k) = ((a \cdot k + b) \bmod p) \bmod m \quad (1)$$

value
between
 $[1, p]$

$$\text{where } \begin{cases} a, b \in \mathbb{N} \\ 1 \leq a < p \\ 0 \leq b < p \end{cases}$$

\Rightarrow Set $H_{p,m}$ is collection of all hash function defined by equation (1)

$$|H_{p,m}| = p \cdot (p-1) \rightarrow \text{very large}$$

Proof: $H_{p,m}$ is Universal Family

For any distinct keys k, l .

$$\text{let } H^* = \{ h_{a,b} : h_{a,b}(k) = h_{a,b}(l) \}$$

Show:

$$|H^*| \leq \frac{p(p-1)}{m} \quad r = p(p-1) = |H_{p,m}|$$

$\forall 1 \leq a < p, 0 \leq b < p$

Fix a pair (a, b)

$$\begin{cases} r = (ak + b) \% p \\ s = (al + b) \% p \end{cases} \Rightarrow \begin{cases} h_{a,b}(k) = r \% m \\ h_{a,b}(l) = s \% m \end{cases}$$

(a) Show bijection between pairs (a, b) and (r, s)

By definition, for distinct keys k and l :

\forall pair (r, s) , \exists a pair (a, b) that maps to it through the function

$$\begin{cases} r = (ak + b) \% p \\ s = (al + b) \% p \end{cases}$$

(b) Using the bijection to show "Number of pairs (r, s) for which $h_{a,b}(k) = h_{a,b}(l)$ is at most $\frac{p(p-1)}{m}$ "

By definition, we know:

$$\begin{cases} h_{a,b}(k) = r \% m \\ h_{a,b}(l) = s \% m \end{cases}$$

$$\Rightarrow P[h_{a,b}(k) = h_{a,b}(l)] = P[r \% m = s \% m]$$

$$\text{linearity} \Rightarrow \sum_{H^*} P[h_{a,b}(k) = h_{a,b}(l)] = \sum_{\substack{H \\ \text{not } (r,s)}} P[r \% m = s \% m] = p$$

$$x \% m = i \% m$$

$$(x < p) \quad (0 \leq i < m)$$

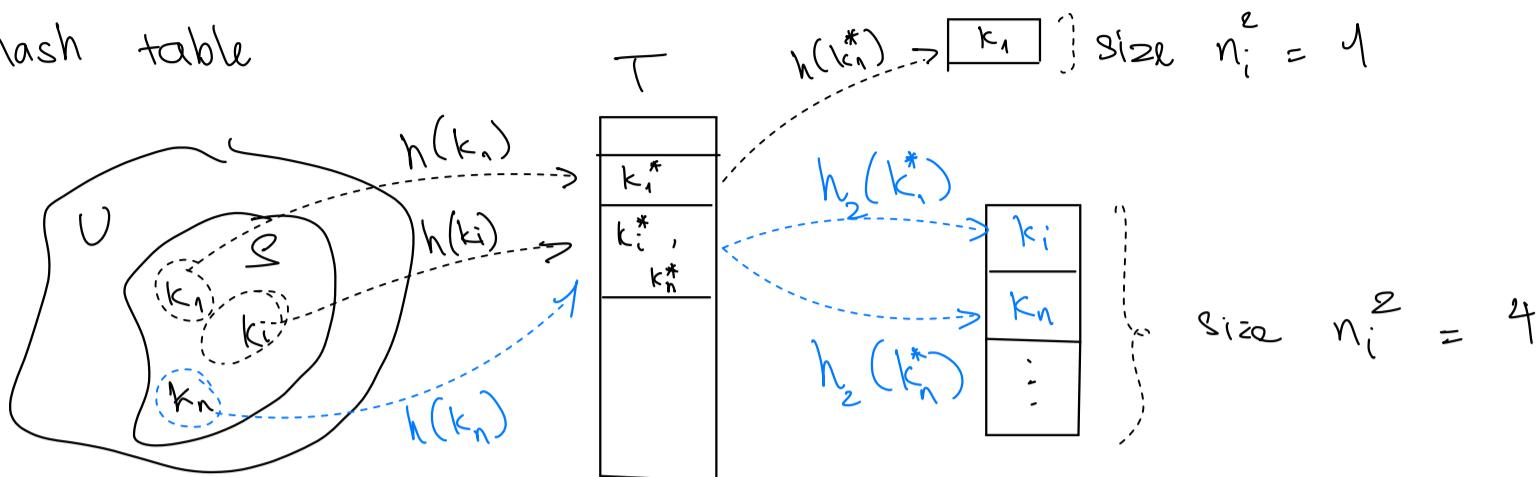
$$\Rightarrow x \in X : |X| \leq \frac{p}{m}$$

Summary: Universal Hashing

- , Formulates hashing as a randomized algorithm
 - a) Makes one initial call to $\text{Random}(1, r)$ to pick a random hash function h from the universal family
 - b) Hashes the set S using hash function h
- Achieves constant expected time for search (using $m = \Theta(n)$)
 - . Insert and Delete take $O(1)$

Perfect Hashing

Main idea: instead of chaining when collision happens, use another hash table



Pseudo:

Step 1: Pick random hash function h from $H_{p,n}$ to map keys in S to primary table $T[0..n-1]$

Step 2: Let n_i be number of keys hashed into i^{th} slot.

If $\sum n_i^2 > 4n$, repeat steps 1-2 with new random hash function until the condition holds.

Step 3: For $i := 0 \rightarrow n-1$:

- Pick random hash function h_i from universal class H_{p,n_i^2}
- Use h_i hash n_i keys already in slot i^{th} of the primary table, into a secondary table of size n_i^2
- If there is collision, repeat Step 3 for i with new random hash function h_i until the keys are hashed into the i^{th} secondary hash table with no collisions.

Claim 1: Expected number of collisions

If n keys are hashed onto table of size m :

$$E[\text{total collisions}] = C(n, 2) \cdot \frac{1}{m}$$

Proof:

Let X be r.v. the number of collisions
 X_{ij} be i.r.v. event that keys i^{th} and j^{th} collide

Then:

$$E[X] = E\left[\sum_{j=i+1}^r \sum_{i=1}^r X_{ij}\right] = \sum_{j=i+1}^r \sum_{i=1}^r E[X_{ij}]$$

$$\leq C(n, 2) \cdot \frac{1}{m}$$

Claim 2: Expected total # of keys squared (primary table)

Let n_i be the keys to be hashed by random function h_i

$$E \left[\sum_{i=0}^{n-1} n_i^2 \right] < 2n \quad \text{and} \quad P \left[\sum_{i=0}^{n-1} n_i^2 > 4n \right] < \frac{1}{2}$$

Proof:

Using this fact: $a^2 = a + 2 \cdot C(a, 2)$ for $a \geq 0, a \in \mathbb{N}$

$$(C(a, 2) = \frac{a \cdot (a-1)}{2})$$

$$\begin{aligned} E \left[\sum_{i=0}^{n-1} n_i^2 \right] &= E \left[\sum_{i=0}^{n-1} (n_i + 2 \cdot C(n_i, 2)) \right] && (\text{above fact}) \\ &= E \left[\sum_{i=0}^{n-1} n_i \right] + 2 E \left[\sum_{i=0}^{n-1} C(n_i, 2) \right] && (\text{linearity}) \\ &= n + 2 \cdot E[\text{total collisions}] \\ &\leq n + 2 \cdot C(n, 2) \cdot \frac{1}{n} && \xrightarrow{\# \text{ collisions in } i^{\text{th}} \text{ slot}} \\ &\leq n + 2 \cdot \frac{n(n-1)}{2} \cdot \frac{1}{n} \\ &\leq n + n - 1 \\ &\leq 2n \end{aligned}$$

$$\begin{aligned} P \left[\sum_{i=0}^{n-1} n_i^2 > 4n \right] &\leq \frac{1}{4n} E \left[\sum_{i=0}^{n-1} n_i^2 \right] && (\text{Markov inequality}) \\ &\leq \frac{1}{4n} 2n && (\text{from above}) \\ &= \frac{1}{2} \end{aligned}$$

\Rightarrow Recall: In step 2, to hash n_i keys into slot i^{th} of primary table, we need this condition to hold $\sum_{i=0}^{n-1} n_i^2 \leq 4n$.

From Claim 2, $P \left[\sum_{i=0}^{n-1} n_i^2 > 4n \right] < \frac{1}{2}$

\Rightarrow Number of iterations taken to run steps 1-2 before we get set of keys n_i that satisfy the condition is 2

Claim 3: Probability collision among n_i keys (secondary table)

If n_i keys are hashed onto secondary table of size n_i^2 , then:

$$P[\text{any collision among } n_i \text{ keys}] < \frac{1}{2}$$

Proof

Let X be number of collisions

By Claim 1, we are hashing n_i keys into n_i^2 slots =

$$E[X] \leq C(n_i, 2) \cdot \frac{1}{n_i^2} < \frac{1}{2}$$

By Markov inequality

$$P[X \geq 1] < E[X] < \frac{1}{2}$$

⇒ Proven

⇒ Result: In step 3, we repeat the iteration with different hash function h'_i picked from H_{p, n_i^2} until there is no collision in the secondary table.

Claim 3 allow us to say that "Expected number of iterations before we picks a "good" function h_i is 2."

Open Addressing (Closed Hashing)

Another hashing technique that uses "probing" to find an empty slot to store the colliding key

- , pros: efficient when load factor is low
- , space efficiency (does not use 2-level hashing scheme like Perfect Hashing)

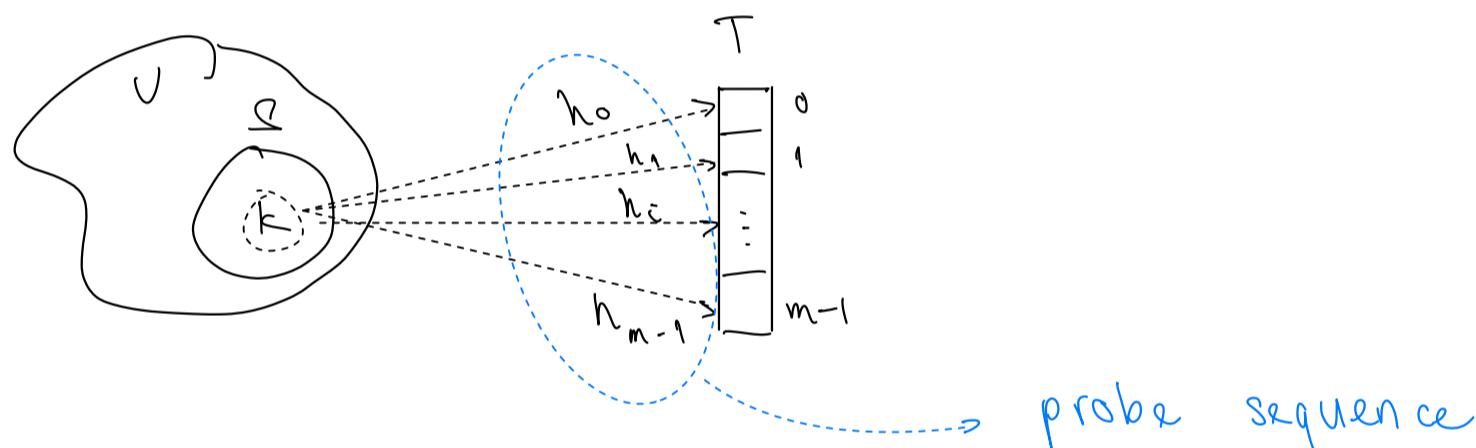
Main idea:

Sequence of hash functions $[h_0, h_1, h_2, \dots, h_{m-1}]$

Probe sequence $[h_0(k), h_1(k), h_2(k), \dots, h_{m-1}(k)]$

↳ is a permutation of $\{0, 1, 2, \dots, m-1\}$

⇒ Different hash functions in the sequence maps key k to different slot in table T



Pseudo : (`SEARCH_KEY` , `INSERT_KEY`)

`SEARCH_KEY (k)`

```

for i = 0 → m-1:
    if T[h_i(k)] == k:
        return h_i(k)
    else if T[h_i(k)] == φ:
        return "absent"
return "full"

```

`INSERT_KEY (k)`

```

for i = 0 → m-1:
    if T[h_i(k)] == k:
        return h_i(k)
    else if T[h_i(k)] == φ:
        T[h_i(k)] ← k
        return "full"

```

Uniform Assumption:

Just like Chaining, we assume the probing sequence is equally likely to be any of the permutations created from set $\{0, 1, 2, \dots, m-1\}$

\Rightarrow Consequences of Uniform Assumption.

Suppose there are currently n elements in the hash table T size m

- Uniformity of hash value $h_i(k)$:

Fix key k and index i , $h_i(k)$ is equally likely to be hashed in any of the slot in table T

\Rightarrow Probability of a slot being occupied is $\frac{n}{m}$:

$$P[h_i(k) \neq \emptyset] = \frac{n}{m} \quad \textcircled{1}$$

- Independence of sub-probe sequence

Ignore the first hash function $h_0(k)$, the remaining probe sequence $[h_1(k), h_2(k), \dots, h_{m-1}(k)]$ is equally likely to be any permutation created from set $\{0, 1, 2, \dots, m-1\} \setminus \{h_0(k)\}$

\Rightarrow If a slot is occupied, recursively perform the same operation on the sub-probe sequence $\{0, 1, 2, \dots, m-1\} \setminus \{h_0(k)\}$ $\textcircled{2}$

$\textcircled{1} \text{ } \& \text{ } \textcircled{2} \Rightarrow$ The recurrence relation for expected number of probes
(# of attempt to find empty slot)

Recurrence for expected number of probes

or "Expected number of attempts to find an empty slot"

$$E[T(m, n)] = 1 + \frac{n}{m} \cdot E[T(m-1, n-1)]$$

Proof: $E[T(m, n)] \leq \frac{1}{1-\alpha} \quad \alpha = \frac{n}{m}$

Base case: $n=0$

$E[T(m, 0)] = 1$ Because guaranteed to find an empty slot in the beginning

Induction: Assume $E[T(m-1, n-1)] = 1 + \frac{n-1}{m-1} \cdot E[T(m-2, n-2)]$

$$E[T(m, n)] = 1 + \frac{n}{m} E[T(m-1, n-1)]$$

$$\leq 1 + \frac{n}{m} \cdot \frac{m-1}{m-1} / \frac{m-n}{m-n}$$

$$< 1 + \frac{n}{m} \cdot \frac{m}{m-n}$$

$$= \frac{m}{m-n} \quad (\text{true})$$

Rewrite $\frac{m}{m-n} = \frac{1}{1-\alpha}$ where $\alpha = \frac{n}{m}$ is load factor

\Rightarrow Expected time for an unsuccessful search is $O(1) \left(\frac{1}{1-\alpha}\right)$

In practice, we can't generate the ideal (uniform) probe sequences, so we must rely on simpler probing scheme to resolve collision

Linear Probing

The simplest probing scheme

$$h_i(k) := (h(k) + i) \bmod m \quad \forall 0 \leq i \leq m-1$$

- Pros:

- Better cache performance (since it visits consecutive entries in the hash table)
- If load factor is strictly $< 1 \Rightarrow$ the expected length of any probe sequence provably constant

- Cons:

- Number of probes grow quickly as load factor approach 1 (since occupied cells in the table tend to stick together)

Example:

1. Initial hash: key k is hashed into slot $h(k)$

2. If slot $h(k)$ is occupied, trigger probing:

$$h_0(k) := (h(k) + 0) \bmod m$$

3. Continue with $i := 1 \rightarrow m-1$ until find empty slot

Note: Wrap around, if $h_i(k)$ reaches end of table, it wraps around to the beginning of the table due to modulo

Binary Probing

Let hash table size $m = 2^l$ ($2, 4, 8, 16, \dots$)

Unlike linear probing, binary probing use bitwise XOR operation to determine next slot

$$h_i(k) := h(k) \oplus i \quad \forall 0 \leq i \leq m-1$$

- Pros:

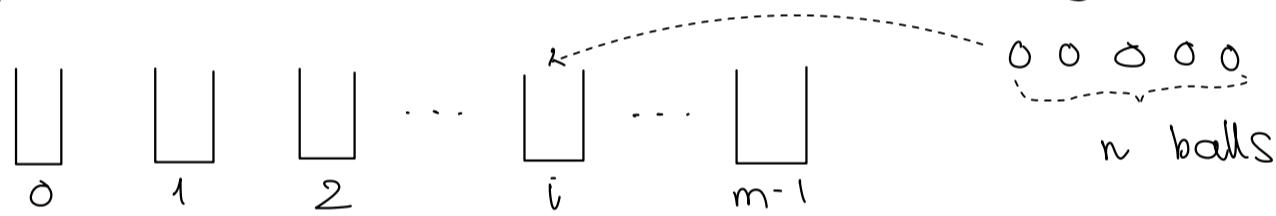
- Less cache misses (since the pattern of binary probing values tend to check current cache before moving to the next cache line)

Example: Consider a hash table of size 16
with cache line of 8.

- Initial hash: $h(k) = 5$
- First probe: $h_1(k) = 5 \oplus 1 = 4$
- Second probe: $h_2(k) = 5 \oplus 2 = 7$
- Third probe: $h_3(k) = 5 \oplus 3 = 6$
- \vdots $h_4(k) = 5 \oplus 4 = 1$
- \vdots $h_5(k) = 5 \oplus 5 = 0$
- \vdots $h_6(k) = 5 \oplus 6 = 3$
- Seventh probe: $h_7(k) = 5 \oplus 7 = 2$ checked all of cache line
- Eighth probe: $h_8(k) = 5 \oplus 8 = 13$ start checking outside cache line

Balls and Bins as Hashing

Throwing n balls into m bins \sim hashing n keys into m -size table



Each ball is thrown independently and uniformly: $P(\text{ball } i \rightarrow \text{bin } j) = \frac{1}{m}$

Let $\{X_j\}$ be r.v. of number of balls in bin- j $\Rightarrow X_j \sim \text{Bin}(n, \frac{1}{m})$
 $\lambda = E[X_j] = \frac{n}{m}$ be the load factor

Rule of thumb for analysis

- If $n \gg m$ (large λ), use Chernoff bounds
- If $n \approx m$ ($\lambda = \Theta(1)$), use Poisson approximation

Case 2: $n \approx m$

Recall that for r.v. $X \sim \text{Bin}(n, p)$ with large n and small p , we can say that $X \sim \text{Pois}(\lambda)$ with PMF:

$$P(X=k) = \frac{e^{-\lambda} \lambda^k}{k!}, \quad k=0, 1, 2, \dots \quad \lambda = n.p = \frac{n}{m}$$

Recall: basic properties of Poisson distribution

$$E[X] = \lambda$$

$$\text{Var}[X] = \lambda$$

$$X + Y \sim \text{Pois}(\lambda + \lambda') \quad \text{for } X \sim \text{Pois}(\lambda) \text{ and } Y \sim \text{Pois}(\lambda')$$

Sol:

$$\text{Since } \begin{cases} X_j \sim \text{Bin}(n, \frac{1}{m}) \\ n \approx m \end{cases}$$

$$\Rightarrow X_j \sim \text{Pois}\left(\frac{n}{m}\right)$$

Facts

- Probability of n balls in same bin: $\left(\frac{1}{m}\right)^n$
- Probability that bin j is empty: $\left(1 - \frac{1}{m}\right)^n \approx e^{-\frac{n}{m}}$
- Let Y be number of empty bins, $E[Y] = ?$
Denote Y_j be i.r.v of if bin- j is empty
 $\Rightarrow Y_j = \begin{cases} 1 & \text{if bin-}j \text{ is empty} \\ 0 & \text{otherwise} \end{cases}$
- Then: $E[Y] = \sum_{j=0}^{m-1} E[Y_j]$ (linearity of expectations)
 $= \sum_{j=0}^{m-1} E[X_j = 0]$
 $= m \left(1 - \frac{1}{m}\right)^n$
 $\approx m \cdot e^{-n/m}$
- Probability j^{th} bin contains 1 ball
 - Exact computation:
 $P[X_j = 1] = C(n, 1) \cdot \left(\frac{1}{m}\right)^1 \cdot \left(1 - \frac{1}{m}\right)^{n-1}$
 $= \frac{n}{m} \left(1 - \frac{1}{m}\right)^n \cdot \left(1 - \frac{1}{m}\right)^{-1}$
 - Poisson approx:
 $P[X_j = 1] = e^{-\lambda} \cdot \lambda = e^{-n/m} \cdot \frac{n}{m}$
- Probability j^{th} bin contains exactly r balls
 - Exact computation:
 $P[X_j = r] = C(n, r) \cdot \left(\frac{1}{m}\right)^r \cdot \left(1 - \frac{1}{m}\right)^{n-r}$
 $= \frac{1}{r!} \left(\frac{n}{m} \cdot \frac{n-1}{m} \cdot \frac{n-2}{m} \cdots \frac{n-r+1}{m}\right) \cdot \left(1 - \frac{1}{m}\right)^n \left(1 - \frac{1}{m}\right)^{-r}$
 $= \lambda^r \cdot e^{-\lambda} \cdot 1^{-r}$
 - Poisson approx:
 $P[X_j = r] = \frac{e^{-\lambda} \cdot \lambda^r}{r!}$

Probability at least one bin has a collision

$$P[\text{at least 1 bin has more than 1 balls}]$$

$$= 1 - P[\text{every bin-}j \text{ has } X_j \leq 1]$$

Fix a ball k^{th} :

$$P[\text{ball } k \rightarrow \text{non-empty bin}] = \frac{k-1}{m}$$

$$P[\text{ball } k \rightarrow \text{empty bin}] = 1 - \frac{k-1}{m}$$

Therefore:

$$\begin{aligned} P[\text{every bin-}j \text{ has } X_j \leq 1] &= \prod_{i=1}^{n-1} \left(1 - \frac{i-1}{m}\right) \\ &\approx \prod_{i=1}^{n-1} e^{-i/m} \end{aligned}$$

$$\Rightarrow P[\text{at least 1 bin has collision}] = 1 - \prod_{i=1}^{n-1} e^{-i/m}$$

Coupon Collector's problem: How many balls we have to throw to assure that every bin contains ≥ 1 balls

Let $\begin{cases} Y \text{ be number of balls to throw until no empty bins left} \\ Y_i \text{ be number of balls to throw until hits an empty bin, given that } i-1 \text{ bins have been filled} \end{cases}$

$$\Rightarrow Y_0 = 1 \quad \text{and} \quad Y = \sum_{i=0}^{m-1} Y_i$$

$$= 1 + Y_1 + Y_2 + \dots + Y_{m-1}$$

$$\Rightarrow \begin{cases} P[\text{ball } i \rightarrow \text{non-empty bin}] = \frac{i-1}{m} \end{cases}$$

$$\begin{cases} P[\text{ball } i \rightarrow \text{empty bin}] = 1 - \frac{i-1}{m} \end{cases}$$

$$\Rightarrow P[Y_i = k] = \left(\frac{i-1}{m}\right)^{k-1} \left(1 - \frac{i-1}{m}\right)$$

$$\Rightarrow Y_i \sim \text{Geometric } (p = 1 - \frac{i-1}{m})$$

$$\Rightarrow E[Y_i] = \frac{1}{p} = \frac{m}{m+i+1}$$

By linearity of expectation:

$$\begin{aligned} E[Y] &= \sum_{i=0}^{m-1} E[Y_i] = \sum_{i=0}^{m-1} \frac{m}{m+i+1} \\ &= m \sum_{i=0}^{m-1} \frac{1}{i+1} \approx m (\ln(m) + O(1)) \end{aligned}$$

Maximum load theorem

If we throw n balls into m bins independently and uniformly, where $m \approx n$, then:

$$\text{Max load of a bin} \leq \frac{4 \lg n}{\lg \lg n} \quad \text{with probability} \leq 1 - \frac{1}{n}$$