

HASHING

1

Hashing or hash addressing is searching technique which is essentially independent of the number n .

Assumptions

- (a) There is file F of n records with a set K of keys which uniquely determine records in F .
- (b) F is maintained in memory by a table T of m memory locations & L is the set of memory addresses of the location in T .

Notations

$F \rightarrow$ File

$n \rightarrow$ records.

$K \rightarrow$ Keys. (it determines records in F)

$T \rightarrow$ table in memory in which F is maintained

$m \rightarrow$ memory locations

$L \rightarrow$ set of memory addresses of locations in T .

EXAMPLE

Suppose a company with 68 employees assign a 4-digit employee number to each employee which is used as the primary key in the company employee file.

$F \rightarrow$ company employee file

$n \rightarrow 68$ (records)

$K \rightarrow$ 4 digit employee number

$L \rightarrow$ Employee number as the address of record in memory.

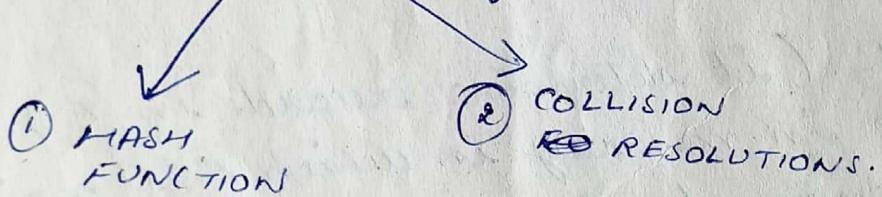
→ Modification of key (k) is done, so that great deal of space is not wasted. This modification takes the form of a function H from the set K of keys into the set L of memory addresses.

$$H : K \rightarrow L$$

H → hash functⁿ or hashing function

H (hash functⁿ) does not yield distinct values i.e. it is possible that two different keys k_1 & k_2 will yield the same address & this situation is called collision.

Division of Hashing in two parts



① HASH FUNCTION

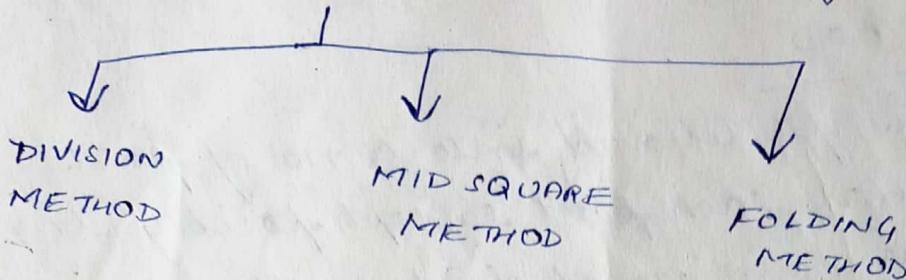
Selection of hash function as follows:

- H should be very easy & quick to compute
- Distribution of hash address should be uniform throughout L so that minimum collision.

→ General technique of hashing is to chop a key k into pieces & combine the pieces in some form to form hash address $H(k)$.

(2)

TYPES OF HASH FUNCTION



(1) DIVISION METHOD / MODULUS DIVISION METHOD

→ Choose a number 'm' larger than the no. n keys in K.

(m is generally prime number or a no. without small divisor)

$$④ H(K) = K \pmod m \quad \left[\begin{array}{l} \text{when range } 0 \text{ to } m-1 \\ \text{OR} \quad K \% m \end{array} \right]$$

$$⑤ H(K) = K \pmod m + 1 \quad \left[\begin{array}{l} \text{when range } 1 \text{ to } m \\ \text{where} \end{array} \right]$$

$K \pmod m$ → denotes the remainder when K is divided by m.

(2) MID SQUARE METHOD

→ The key K is squared, then the hash function H is defined by

$$H(K) = l$$

→ where l is obtained by deleting digits from both ends of K^2 . # same position of K^2 must be deleted for all of the keys.

③ FOLDING METHOD

- The key K is partitioned into a no. of parts K_1, \dots, K_r where each part except possibly the last, has the same no. of digit as the required address.
- Then the parts are added together, ignoring the last ~~last~~ carry i.e.

$$H(K) = K_1 + K_2 + \dots + K_r$$

OR

Alternatively, one may reverse the second part before adding.

EXAMPLE :

A company with 68 employees is assigned a unique 4-digit employee number. Suppose L consists of 100 two-digit addresses: 00, 01, 02, ..., 99.

Apply hash function

3205, 7148, 2345

④ DIVISION METHOD: (MODULUS-DIVISION METHOD)

$m = 97$ (\because it is prime number close to 99)

$$H(3205) = 3205 \% 97 = 4 \quad (3205 / 97 \text{ give remainder } 4)$$

$$H(7148) = 7148 \% 97 = 67$$

$$H(2345) = 2345 \% 97 = 17$$

[In this case memory add^r starting 00, so formula ① is used]

(3)

5. MIDSQUARE METHOD

$$K = 3205 \quad 7148 \quad 2345$$

$$K^2 \quad \underline{10272} \overset{45}{,} \underline{025} \quad \underline{51093} \overset{45}{,} \underline{904} \quad \underline{5499} \overset{45}{,} \underline{025}$$

$$H(K) : \quad 72 \quad 93 \quad 90$$

C) FOLDING METHOD

Chopping the key K into two parts

$$H(3205) = 32 + 05 = 37$$

$$H(7148) = 71 + 48 = 19$$

$$H(2345) = 23 + 45 = 68$$

leading digit 1 in $H(7148)$ is ignored

Alternatively, reverse the second part before adding,
thus producing the foll. hash addresses

$$H(3205) = 32 + 50 = 82$$

$$H(7148) = 71 + 84 = 55$$

$$H(2345) = 23 + 54 = 77$$

Example 2

- Consider a company with 68 employee. Each has been assigned a 4 digit employee number which is used as the primary key in the company's employee file. Suppose VL (set of memory addresses of the locations) consist of 100 two digit addresses: 00, 01, 02, ..., 99. Then apply the division method using a prime no. closest to 99, mid square method & folding method to find out the 2-digit hash add⁸ for each of the foll employ no. 9614, 5882, 1825

Division method

$$9614 \rightarrow 9614 \% 97 = 11$$

$$5882 \rightarrow 5882 \% 97 = 62$$

$$1825 \rightarrow 1825 \% 97 = 79$$

Mid square

$$K : 9614 \quad 5882 \quad 1825$$

$$K^2 : (9614)^2 \quad (5882)^2 \quad (1825)^2$$

$$\begin{array}{r} 924 \ 28 \ 996 \\ \hline 28 \end{array} \quad \begin{array}{r} 3459 \ 7924 \\ \hline 97 \end{array} \quad \begin{array}{r} 333 \ 0625 \\ \hline 06 \end{array}$$

Folding method

$$K: 9614 \rightarrow 96 + 14 = 00$$

$$: 5882 \rightarrow 58 + 82 = 30$$

$$1825 \rightarrow 18 + 25 = 33$$

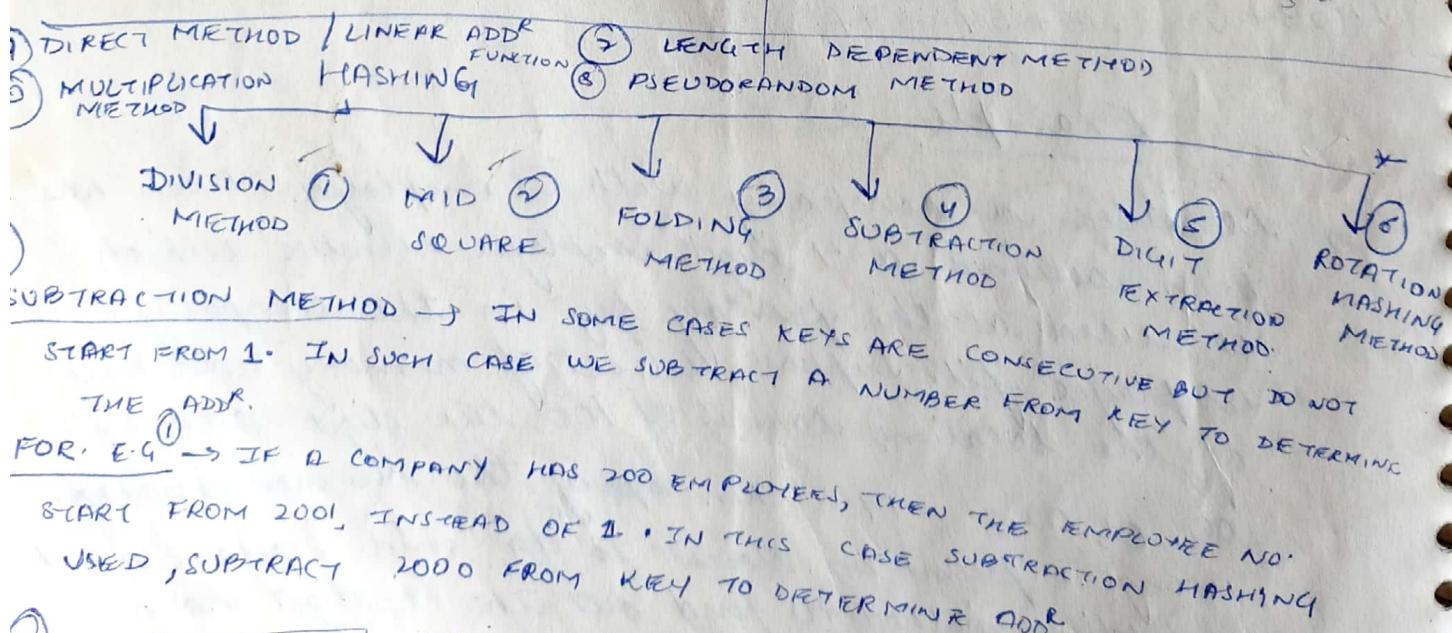
Reverse

$$96 + 41 = 37$$

$$58 + 28 = 76$$

$$18 + 52 = 60$$

$$\begin{array}{r} 59 \\ 32 \\ \hline 30 \end{array}$$

$$\begin{array}{r} 58 \\ 82 \\ \hline 30 \end{array}$$


101	JANE
102	X
103	Y
104	Z
105	A
106	B

→ FIND ADD OF 'Y'

→ SOLN → USING SUBTRACTION METHOD JUST SUBTRACT 100 FROM THE KEY 103. → ADD 15 →

- (5) → DIGIT EXTRACTION METHOD → (a) SELECTED KEYS ARE EXTRACTED FROM THE KEY & MADE USE AS ITS ADD^R USING A METHOD CALLED DIGIT EXTRACTION.
- (b) IN THIS CASE, WE SELECT SPECIFIC DIGITS FROM KEY K & USE IT AS AN ADD^R.
- FOR E.G. (1) SUPPOSE WE WANT A HASH A 6 DIGIT EMPLOYEE NUMBER, 123457, TO A THREE DIGIT ADD^R, WE COULD SELECT FIRST, THIRD, & FOURTH DIGITS FROM LEFT & USE THEM AS AN ADD^R SO ADD^R WILL BE 124

Ques:

(2) SUPPOSE ROLL NO. OF A STUDENT WILLIAMS IS 160252. HASH THE NUMBER TO A THREE DIGIT ADD^R USING DIGIT EXTRACTION METHOD.

Ans: SELECT FIRST, THIRD, & FOURTH DIGIT & USE IT AS THE ADD^R

$$\underline{1} \underline{6} \underline{0} \underline{2} \underline{2} \underline{5} \underline{2} \rightarrow 102$$

$$\therefore \text{ADD}^R \text{ IS } \underline{\underline{102}}$$

(6) ROTATION HASHING METHOD →

→ THIS METHOD IS GENERALLY NOT USED BY ITSELF, BUT IS USED IN COMBINATION WITH OTHER HASHING METHODS. THIS METHOD IS ESPECIALLY USEFUL WHEN KEYS ARE ASSIGNED SEQUENTIALLY

→ E.G.

<u>ORIGINAL KEY</u>	<u>ROTATION</u>	<u>KEY AFTER ROTATION</u>
500201	500201	150020
500202	500202	250020
500203	500203	350020
500204	500204	450020
500205	500205	550020

(7) LENGTH DEPENDENT METHOD

- IN THIS METHOD LENGTH OF THE KEY IS USED ALONG WITH SOME PORTION OF THE KEY.
- FOR E.G. KEY IS 12345 THE 123 IS TAKEN FROM KEY AND ADDED WITH LENGTH OF THE KEY i.e. 5
- ∴ ID IS 128

⑧ PSEUDORANDOM METHOD

→ A random number is generated using the key. Using the modulo-division method a pointer is generated to the key. Formula to generate the pointer is

$$P = (a * \text{key} + c) \% \text{Tablesize}$$

where $a \rightarrow$ coefficient & c is constant.

⑨ DIRECT METHOD OR LINEAR PROBING

- The key is the address without any algorithmic manipulation.
- Data structure must \ni contain an element for every possible key.

E.G.

⑩ Multiplication Method

→ Steps

- ① Multiply the key by a constant A , $0 < A < 1$
- ② Extract the fractional part of the product
- ③ multiply this value by m .

$$\therefore H(k) = \text{FLOOR}(m * (KA - \text{FLOOR}(KA)))$$

E.G. Ques. \rightarrow 80 location
 $A = 0.618033 \rightarrow$ good choice
 $\rightarrow 2345, 5378, 7321$

$$\begin{aligned}H(2345) &= \text{FLOOR}(79 * (2345 * 0.618033 - \text{FLOOR}(2345 * 0.618033))) \\&= \text{FLOOR}(79 * (1426.071885 - \text{FLOOR}(1426.071885))) \\&= \text{FLOOR}(79 * (1426.071885 - 1426)) \\&= \text{FLOOR}(5.678915) \\&= 5\end{aligned}$$

$$H(5378) = 61 \quad H(7321) = 48, H(6911) = 17$$

COLLISION

- A collision is the event that occurs when hashing algorithm produce an address for an insertion key & that address is already occupied.
- Address produced by hashing algorithm is known as home address.
- Prime area → memory that contains all of the home address.
- When two keys collide at a home address, to resolve the collision by placing one of the key & its data in another locations. Each calculation of an address & test for successive place is known as probe.

CLUSTERING

- Some hashing algorithms tend to cause data to group within the list. This tendency of data to build up unevenly across a hashed list is known as clustering.
- TYPES OF CLUSTERING
 - ↓
 - PRIMARY CLUSTERING
 - ↓
 - It occurs when data become clustered around a home address
 - It is easy to identify
 - Random probing & quadratic probing
 - ↓
 - SECONDARY CLUSTERING
 - ↓
 - When data grouped along a collision path throughout a list
 - It is not easy to identify
 - It can be solved by second hashing function i.e. double hashing.

Collision Resolution Technique

→ To avoid collisions.

(a) Spread out records →

Finding a hashing algorithm that distributes the record fairly randomly among available addresses.

(b) Use Extra Memory →

→ If there are more memory addresses to distribute the record

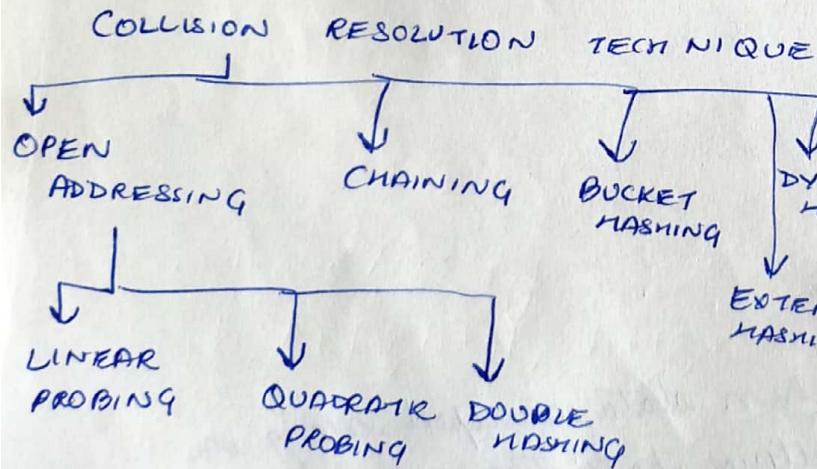
→ Then it is easier to find a hashing algorithm.

Adv: → Records are spread out evenly.

Disadv: → waste spaces.

(c) Use Bucket →

→ put more than one record at single address.



(d) Linear probing

→ Expected no. of probes using linear probing for insertions & unsuccessful search is

$$\frac{1}{2} \left(1 + \frac{1}{(1-\lambda)^2} \right)$$

Successful search

$$\frac{1}{2} \left(1 + \frac{1}{(1-\lambda)} \right)$$

where $\lambda \rightarrow$ load factor
load factor of a hashed list is the no. of elements in the list divided by no. of physical elements allocated for the list expressed as percentage

$$\lambda = K/n * 100$$

where $n \rightarrow$ table size

$K \rightarrow$ no. of element in table

(d) Quadratic probe

Adv:

- (a) simple to implement
- (b) data tend to remain near the home address

Disadv:

- (a) produce primary clustering.
- (b) complicated search algo when data is deleted

Quadratic Probing

- increment is collision probe no.
- $h+1^2, h+2^2, \dots, h+i^2$ squared

Disadvantage

- Time required to square the probe
- It is not possible to generate new address for every element in the list.

Double hashing

E.9

0	845
1	444
2	444
3	
4	845
5	902
6	981
7	345
8	125
9	286
10	369
11	947
12	792
13	

INSERT 652 into hash table

result in overflow at 2 we rehash

$$h_2(652) = 652 \bmod 11 \xrightarrow{\text{NO. OF ELEMENTS}} = 3$$

$$rh(2, 652) = (2+3) \bmod 13 = 5 \text{ already full}$$

then rehash

$$rh(5, 652) = (5+3) \bmod 13 = 8$$

with space so store here.

Adv:

- improvement over linear probing
- i.e. New add^{re} is computed using another hash function instead of sequentially searching the hash table for an empty space.

Disadv-

- suffers from the displacement problem.

Example

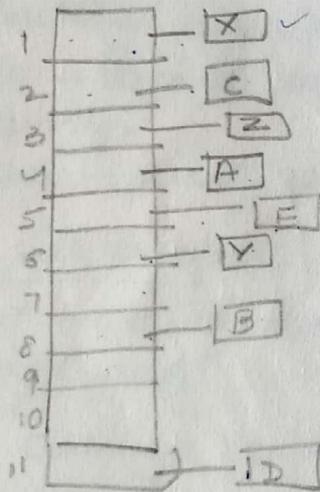
Suppose table T has 11 memory locations
 $T[1], T[2], \dots, T[11]$ suppose F consist of 8 records
 A, B, C, D, E, X, Y and Z with foll hash addr.

RECORD: A, B, C, D, E, X, Y, Z
 H(k) : 4, 8, 2, 11, 4, 11, 5, 1

find S (Successful search probe) & U (Unsuccessful search probe) for

- (a) linear probing (b) quadratic probing (c) double hashing.

(a) LINEAR PROBING

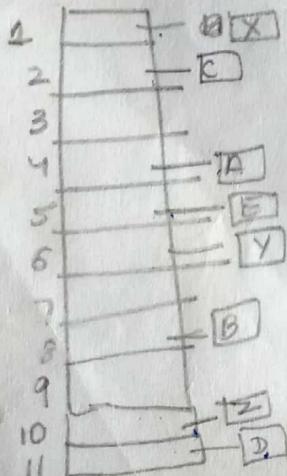


$$S = \frac{1+1+1+1+2+2+2+3}{8} = \frac{13}{8} = 1.6$$

$$U = \frac{7+6+5+4+3+2+1+2+1+1+8}{11} = \frac{40}{11} = 3.6$$

(b) Quadratic Probing

$$h, h+1, h+4, h+9, h+16, \dots, h+c^2$$



$$X \rightarrow$$

$$E \rightarrow h+1 = 5$$

$$X = h+1 = 1$$

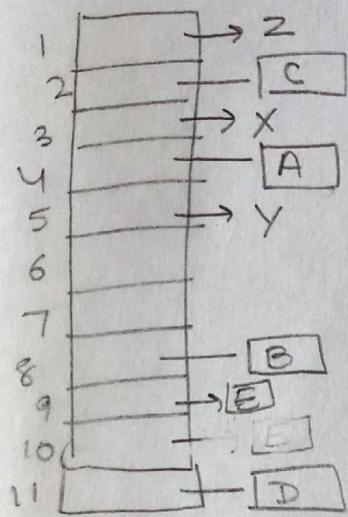
$$Y = h,$$

$$Z = h, h+1, h+4, h+9, h+16, h+25, h+36$$

$$S = \frac{1+1+1+1+2+2+2+4}{8}$$

$$U = \frac{(3+2+1+4+3+2+1+2+1+5+4)}{11} = \frac{19}{11} = 1.727$$

③ Double hashing



INSERT E

$$\rightarrow h' = 4 \bmod 8 = 5$$

$$\text{or } h(19, 5) = (19) \bmod 11 = 9$$

$\rightarrow X \rightarrow$

$$h' = 11 \bmod 8 = 3$$

$$h(14, 3) = 14 \bmod 11 = 3$$

$$S = \frac{1+1+1+1+2+2+1+1}{8} = \frac{10}{8} = 1.25$$

$$U = \frac{6+5+4+3+2+1+1+3+2+1+7}{11} = \frac{35}{11} = 3.18$$