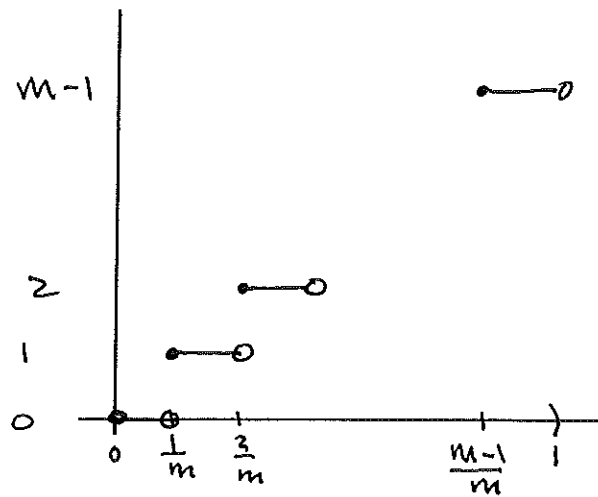


EX.

IF THE KEYS ARE UNIFORMLY DISTRIBUTED OVER  $U = \mathbb{Z}$ , THEN  $h(k) = k \bmod m$  WORKS.

EX.

IF THE KEYS ARE UNIFORMLY DISTRIBUTED OVER  $0 \leq k < 1$ , i.e.  $U = [0, 1)$  THEN  $h(k) = \lfloor km \rfloor$  SATISFIES SIMPLE UNIFORM HASHING.



OBSERVE THAT

$$h: \left[ \frac{i}{m}, \frac{i+1}{m} \right) \rightarrow \{i\} \quad (0 \leq i \leq m-1).$$


IF THE PROBABILITY FUNCTION  $P(k)$  IS NOT KNOWN, WE MUST TAKE A HEURISTIC APPROACH TO CHOOSING  $h$ .

WE WILL ASSUME FROM NOW ON THAT  $U = M = \{0, 1, \dots\}$ .

EX. (DIVISION METHOD)

$$h(k) = k \bmod m.$$

WHEN THE DISTRIBUTION OF KEYS IS UNKNOWN, THE CHOICE OF  $m$  CAN BE CRITICAL. FOR INSTANCE, WE SHOULD NOT CHOOSE  $m = 2^p$  FOR THEN  $h(k)$  IS JUST THE  $p$ . LOWEST ORDER BITS OF  $k$ .

e.g.  $m = 2^3$ ,  $k = 10110(101) = 10110000 + 101$   
 $\therefore h(k) = 101$  

IT WOULD BE BETTER TO MAKE  $h(k)$  DEPEND ON ALL BITS IN  $k$ .

GOOD VALUES FOR  $m$  ARE PRIMES NOT CLOSE TO EXACT POWERS OF 2.

EX. (MULTIPLICATION METHOD)

FIX A NUMBER  $A$  WITH  $0 < A < 1$ .

DEFINE

$$h(k) = \lfloor m (kA - \lfloor kA \rfloor) \rfloor$$

FRACTIONAL PART OF  $kA$

THE ADVANTAGE HERE IS THAT THE VALUE OF  $m$  IS NOT CRITICAL.

IF  $m = 2^p$  THEN THIS FUNCTION IS EASY TO IMPLEMENT AT THE ASSEMBLY LANGUAGE LEVEL.

FOR INSTANCE, SUPPOSE THE WORD SIZE ON OUR COMPUTER IS  $w = 8$ , AND  $k$  FITS IN A SINGLE WORD.

SAY  $m = 2^5$  AND  $A = .8203125$ .  
THEN

$$A = \quad \quad \quad .11010010$$

$$K = \quad \quad \quad 10110101$$

$$KA = \quad \quad \quad 10010100 \cdot \boxed{0111}010$$

$$KA - \lfloor KA \rfloor = \quad \quad \quad .0111011$$

$$m(KA - \lfloor KA \rfloor) = \quad \quad \quad 0111.011$$

$$h(k) = \lfloor m(KA - \lfloor KA \rfloor) \rfloor = \quad \quad \quad 0111$$

NOTICE THAT IN THIS EXAMPLE  $h(k) = 0111$  BEARS NO OBVIOUS RELATIONSHIP TO  $k = 10110101$  (IN TERMS OF BIT PATTERN.)

SOME VALUES OF  $A$  ARE THOUGHT TO BE BETTER THAN OTHERS.

OPEN ADDRESSING IS AN ALTERNATIVE TO CHAINING FOR RESOLUTION OF COLLISIONS.

IN THIS METHOD ALL ELEMENTS ARE STORED IN THE HASH TABLE ITSELF, SO IT IS POSSIBLE FOR THE TABLE TO FILL UP, OR "OVERFLOW".

TO PERFORM INSERTION USING OPEN ADDRESSING WE EXAMINE OR PROBE POSITIONS IN THE HASH TABLE UNTIL WE FIND AN EMPTY SLOT.

THE HASH FUNCTION IS OF THE FORM

$$h : U \times \{0, \dots, m-1\} \rightarrow \{0, 1, \dots, m-1\}.$$

WE REQUIRE THAT FOR EACH  $k \in U$ , THE PROBE SEQUENCE

$$(h(k,0), h(k,1), h(k,2), \dots, h(k,m-1))$$

BE A PERMUTATION OF  $(0, 1, 2, \dots, m-1)$ , SO THAT EVERY SLOT IS EVENTUALLY CONSIDERED AS A POSITION FOR A NEW TABLE ENTRY.

EX.

$$h(k, 2) = 0$$

$$h(k, 5) = 1$$

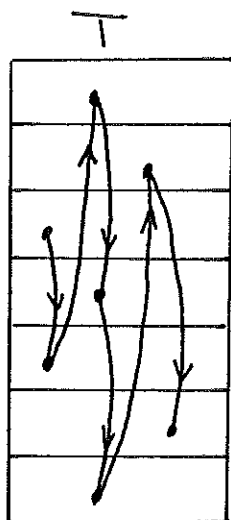
$$h(k, 0) = 2$$

$$h(k, 3) = 3$$

$$h(k, 1) = 4$$

$$h(k, 6) = 5$$

$$h(k, 4) = 6$$



$$m = 7$$

IN THE FOLLOWING WE ASSUME THAT TABLE ENTRIES ARE KEYS WITH NO SATELLITE DATA.

Insert (T, k)

- 1.)  $i \leftarrow 0$
- 2.) WHILE  $i < m$
- 3.)      $j \leftarrow h(k, i)$
- 4.)     IF  $T[j] = \text{NIL}$
- 5.)          $T[j] \leftarrow k$
- 6.)         Return  $j$
- 7.)     ELSE
- 8.)          $i \leftarrow i + 1$
- 9.) error "overflow"

NOTE THAT Insert RETURNS A VALUE, NAMELY THE SLOT  $j$  IN WHICH KEY  $k$  WAS EVENTUALLY PLACED.

Search (T, k)

- 1.)  $i \leftarrow 0$
- 2.) Repeat
- 3.)  $j \leftarrow h(k, i)$
- 4.) If  $T[j] = k$
- 5.) Return  $j$
- 6.)  $i \leftarrow i + 1$
- 7.) until  $i = m$  or  $T[j] = \text{NIL}$ .
- 8.) Return NIL.

DELETIONS CAN CAUSE PROBLEMS IN AN OPEN ADDRESS TABLE, SINCE THE VALUE NIL IS USED TO CANCEL A SEARCH.

ONE SOLUTION IS TO DEFINE A SPECIAL VALUE CALLED DELETED, WHICH DOESN'T STOP A SEARCH, BUT WHICH DOES STOP THE INSERT PROCESS. STEP (4) OF INSERT BECOMES: IF  $T[j] = \text{NIL}$  OR  $T[j] = \text{DELETED}$ ,

Delete (T, i)

- 1.)  $T[i] \leftarrow \text{DELETED}$

THERE ARE A NUMBER OF WAYS TO DESIGN THE HASH FUNCTION

$$h: U \times \{0 \dots m-1\} \rightarrow \{0 \dots m-1\}$$

FOR AN OPEN ADDRESS TABLE.

Ex. (Linear Probing)

$$h(k, i) = (h'(k) + ci) \bmod m$$

WHERE  $h' : U \rightarrow \{0 \dots m-1\}$  IS AN ORDINARY HASH FUNCTION AND  $c$  IS A POSITIVE INTEGER, RELATIVELY PRIME TO  $m$ .

PROBE SEQUENCE:

$$(h'(k), h'(k) + c, h'(k) + 2c, \dots, h'(k) + (m-1)c)$$

Ex. (Quadratic Probing)

$$h(k, i) = (h'(k) + c_1 i + c_2 i^2) \bmod m$$

BOTH OF THESE METHODS SUFFER FROM CLUSTERING: LONG RUNS OF OCCUPIED SLOTS, INCREASING THE AVERAGE SEARCH TIME.

Ex. (Double Hashing)

$$h(k, i) = (h_1(k) + ih_2(k)) \bmod m$$

WHERE  $h_1, h_2 : U \rightarrow \{0 \dots m-1\}$  ARE AUXILIARY HASH FUNCTIONS -

THE PROBE SEQUENCE IS (MOD  $m$ )

$$(h_1(k), h_1(k) + h_2(k), h_1(k) + 2h_2(k), \dots, h_1(k) + (m-1)h_2(k))$$

OBSERVE: THE INITIAL PROBE AND THE INTERVAL BETWEEN SUCCESSIVE PROBES MAY VARY WITH  $k$ . THUS NO TWO PROBE SEQUENCES LOOK ALIKE.

NOTE ALSO THAT  $h_2(k)$  MUST BE RELATIVELY PRIME TO  $m$  FOR THE WHOLE HASH TABLE TO BE SEARCHED.

e.g.  $m=6$ ,  $h_2(k)=3$  YIELDS THE PROBE SEQUENCE

$$(h_1(k), h_1(k)+3, h_1(k)+6, h_1(k)+9, h_1(k)+12, h_1(k)+15)$$

BUT (MOD 6) THERE ARE ONLY TWO TWO TERMS IN THIS SEQUENCE, NOT 6.

ONE WAY TO AVOID THIS IS TO CHOOSE  $m$  TO BE PRIME. ANOTHER WAY IS TO CHOOSE  $m=2^p$ , AND DESIGN  $h_2$  SO THAT ITS VALUES ARE ALWAYS ODD.