

## 12 HASH TABLES

DYNAMIC SET IS A GENERAL TERM FOR ANY DATA STRUCTURE REPRESENTING A SET WHICH MAY GROW, SHRINK, OR OTHERWISE CHANGE OVER TIME.

EXAMPLES INCLUDE STACKS, QUEUES, LISTS, HEAPS, PRIORITY QUEUES, TREES, I.E. EVERYTHING WE ARE STUDYING.

$$S = \{ \dots \underbrace{(\cdot, \cdot, \cdot, \text{key})}_{\text{SATELLITE DATA}} \dots \}$$

A DICTIONARY IS A DYNAMIC SET WHICH SUPPORTS THE OPERATIONS

$\text{Insert}(S, x)$ : INSERT ELEMENT (POINTED TO BY)  $x$  INTO  $S$ .

$\text{Delete}(S, x)$ : DELETE ELEMENT (POINTED TO BY)  $x$  FROM  $S$ .

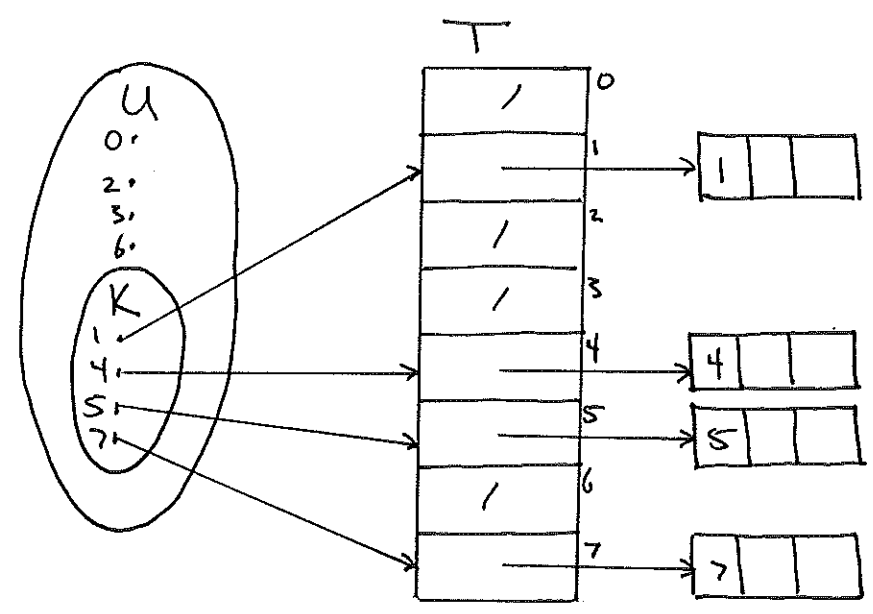
$\text{Search}(S, k)$ : RETURN A POINTER  $x$  TO AN ELEMENT IN  $S$  SUCH THAT  $\text{key}[x] = k$ , OR NIL IF NO SUCH ELEMENT EXISTS.

SUPPOSE AN APPLICATION NEEDS A DICTIONARY IN WHICH EACH ELEMENT HAS A KEY DRAWN FROM A UNIVERSE

$$U = \{0, 1, \dots, m-1\}$$

OF POSSIBLE KEYS. WE ASSUME NO TWO ELEMENTS HAVE THE SAME KEY, BUT NOT ALL KEYS ARE NECESSARILY IN USE.

ONE WAY TO PROCEED IS TO HAVE AN ARRAY  $T[0 \dots (m-1)]$  WHOSE INDICES ARE JUST THE KEYS.



AN ARRAY USED IN THIS WAY IS CALLED A DIRECT ADDRESS TABLE.

THE DICTIONARY OPERATIONS ARE TRIVIAL TO IMPLEMENT, AND RUN IN  $O(1)$  TIME (WORST CASE), FOR INSTANCE

Search (T, k)

1.) return  $T[k]$

THIS WORKS FINE IF  $|U|$  IS NOT TOO LARGE. IF  $|U|$  IS LARGE, AND  $|K|$  IS SMALL RELATIVE TO  $|U|$ , THEN THIS ARRANGEMENT WASTES A LOT OF SPACE.

ONE ALTERNATIVE IS TO USE A LINKED LIST INSTEAD OF AN ARRAY. THE SEARCH OPERATION WOULD THEN TAKE  $\Theta(n)$  TIME (AVERAGE CASE), WHICH IS CONSIDERED SLOW.

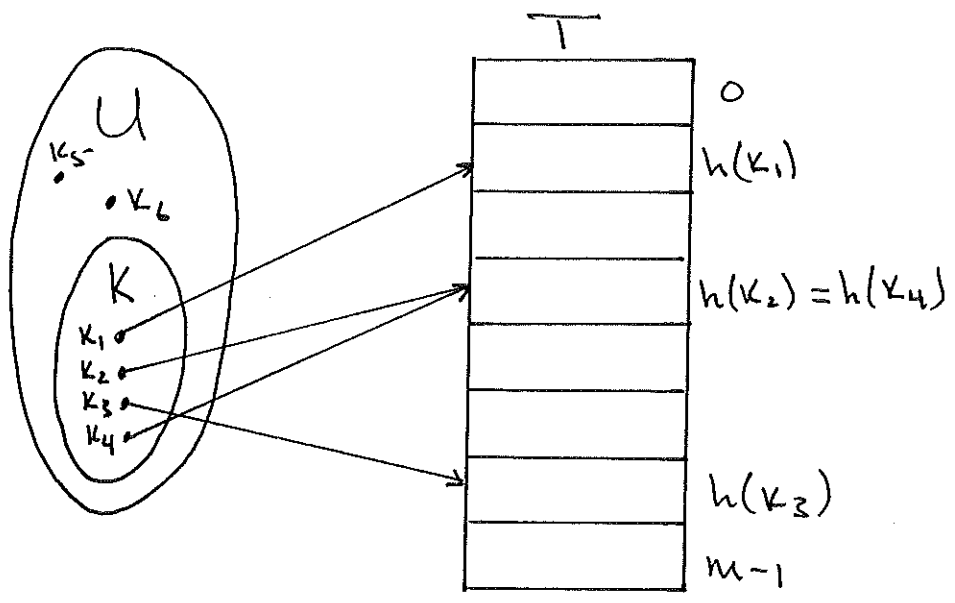
A THIRD ALTERNATIVE IS TO USE A HASH TABLE, WHERE STORAGE REQUIREMENTS ARE  $\Theta(|K|)$  AS IN A LINKED LIST, AND SEARCH TAKES  $O(1)$  TIME (NOW IN AVERAGE CASE).

A HASH TABLE IS AN ARRAY  $T[0 \dots (m-1)]$  IN WHICH AN ELEMENT WITH KEY  $k$  IS STORED, NOT IN  $T[k]$ , BUT IN  $T[h(k)]$ .

HERE  $h: U \rightarrow \{0, 1, \dots, m-1\}$  IS CALLED A HASH FUNCTION. WE SAY KEY  $K$  HASHES TO SLOT  $h(k)$ , AND THAT  $h(k)$  IS THE HASH VALUE OF  $K$ .

IT MAY BE HOWEVER THAT TWO (OR MORE) KEYS HASH TO THE SAME SLOT:  $h(k_1) = h(k_2)$ . THIS IS CALLED A COLLISION.

IN OTHER WORDS  $h$  IS NOT INJECTIVE, SINCE THE WHOLE POINT IS TO SAVE SPACE  $m < |U|$ .

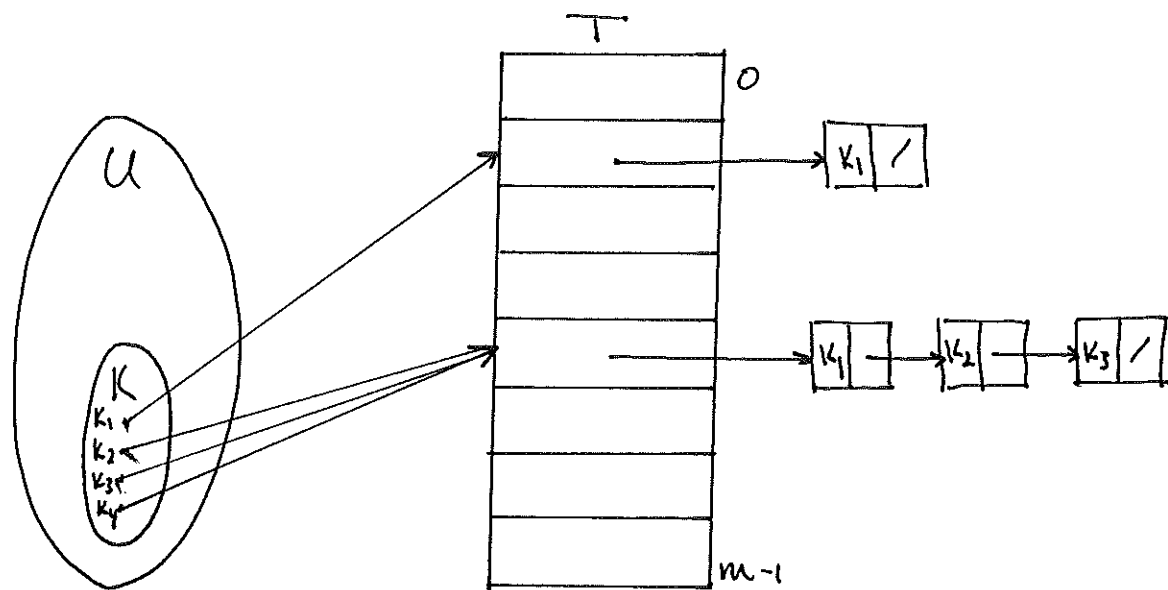


FORTUNATELY THERE ARE EFFECTIVE METHODS FOR RESOLVING COLLISIONS.

EVEN SO, THE FUNCTION  $h$  SHOULD BE CHOSEN IN SUCH A WAY AS TO MINIMIZE THE PROBABILITY THAT TWO RANDOM KEYS COLLIDE. IN OTHER WORDS,  $h$  "DISTRIBUTES" THE PROBABILITY OF A COLLISION EVENLY ACROSS THE ARRAY INDICES  $\{0, \dots, m-1\}$ .

$h$  IS THEREFORE DESIGNED TO LOOK "RANDOM" IN SOME SENSE, HENCE THE NAME HASH FUNCTION.

CHAINING IS A COLLISION RESOLUTION METHOD IN WHICH ALL ELEMENTS WHOSE KEYS HASH TO THE SAME SLOT ARE PLACED IN A LINKED LIST.



THE DICTIONARY OPERATIONS ARE THEN IMPLEMENTED BY CALLING LIST OPERATIONS.

Insert  $(T, x)$

insert  $x$  INTO ITEMS OF LIST  $T[h(\text{key}[x])]$

Delete  $(T, x)$

delete  $x$  FROM LIST  $T[h(\text{key}[x])]$

Search  $(T, k)$

SEARCH FOR AN ELEMENT  $x$  SUCH THAT  $\text{key}[x] = k$  IN LIST  $T[h(k)]$ .

Insert TAKES  $O(1)$  TIME (WORST CASE).

THE TIME TO RUN SEARCH DEPENDS ON THE NUMBER OF ELEMENTS IN THE LIST  $T[h(k)]$ , WHICH DEPENDS ON  $k$ .

Delete TAKES  $O(1)$  TIME IF THE LISTS ARE DOUBLY LINKED.

WE WILL ANALYSE THE RUN TIME OF SEARCH UNDER CHAINING.

GIVEN A HASH TABLE  $T[0 \dots (m-1)]$  WITH  $m$  SLOTS WHICH STORES  $n = |K|$  ELEMENTS, DEFINE THE LOAD FACTOR  $\alpha$  TO BE

$$\alpha = \frac{n}{m} .$$

Thus  $\alpha$  is the average number of elements stored in any list.

As noted earlier, the performance of search depends on the choice of hash function  $h: U \rightarrow \{0, \dots, m-1\}$ .

We make the following assumption about  $h$ .

- Any randomly selected key is equally likely to hash into any of the  $m$  slots. i.e.

$$\text{Prob}(h(k) = i) = \frac{1}{m} \quad \text{for } 0 \leq i < m .$$

We call this the assumption of simple uniform hashing.

In the language of probability theory,  $U$  is the sample space and  $h: U \rightarrow \{0, \dots, m-1\}$  is a discrete random variable which is uniformly distributed, i.e. the probability density function of  $h$  is the constant  $\frac{1}{m}$ .

NOTE WE DO NOT ASSUME THAT ANY TWO KEYS ARE EQUALLY LIKELY TO BE SELECTED. SOME KEYS MAY BE MORE LIKELY TO BE USED THAN OTHERS.

THUS TO CONFORM TO THE ASSUMPTION OF SIMPLE UNIFORM HASHING,  $h$  MUST BE CHOSEN WITH THE DISTRIBUTION OF KEYS IN MIND.

TO ANALYZE SEARCH IN THE AVERAGE CASE WE RECALL

DEFN

IF  $f(i)$  IS A FUNCTION DEFINED ON A SEQUENCE OF RANDOM EVENTS  $i$ , THE EXPECTED VALUE OR AVERAGE VALUE OF  $f$  IS

$$E[f] = \sum_i f(i) \cdot (\text{PROB. OF EVENT } i)$$

EX. TWO FAIR DICE ARE THROWN AND THE SUM  $i$  IS RECORDED. LET  $f(i) = i$ . THEN

$i =$	2	3	4	5	6	7	
	11	12	13	14	15	16	8
	21	22	23	24	25	26	9
	31	32	33	34	35	36	10
	41	42	43	44	45	46	11
	51	52	53	54	55	56	12
	61	62	63	64	65	66	

$$E[L] = 2 \cdot \frac{1}{36} + 3 \cdot \frac{2}{36} + \dots + 7 \cdot \frac{6}{36} + \dots + 12 \cdot \frac{1}{36} = 7$$

WE CONSIDER TWO CASES FOR THE RUN TIME OF SEARCH.

CASE 1: THE SEARCH IS UNSUCCESSFUL.

LET  $f(i)$  BE THE NUMBER OF KEYS STORED IN SLOT  $i$ , ( $0 \leq i \leq m-1$ ). THEN

$$\sum_{i=0}^{m-1} f(i) = n,$$

THE NUMBER OF KEYS IN USE. THE AVERAGE LIST LENGTH IS THEN

$$E[L] = \sum_{i=0}^{m-1} f(i) \cdot P(h(k) = i) = \sum_{i=0}^{m-1} \frac{f(i)}{m} = \frac{n}{m} = \alpha.$$

THE RUN TIME OF AN UNSUCCESSFUL SEARCH IS THEREFORE PROPORTIONAL TO  $\alpha$ .

CASE 2: THE SEARCH IS SUCCESSFUL.

AGAIN LET  $f(i)$  BE THE NUMBER OF KEYS STORED IN SLOT  $i$ , ( $0 \leq i \leq m-1$ ). ON AVERAGE THE KEY SEARCHED FOR WILL BE FOUND HALFWAY DOWN THE LIST.

THE AVERAGE NUMBER OF KEYS EXAMINED IN SLOT  $i$  IS THEREFORE  $f(i)/2$ .

THE NUMBER OF KEYS EXAMINED, AVERAGED OVER ALL SLOTS IS

$$E\left[\frac{f}{2}\right] = \sum_{i=0}^{m-1} \frac{f(i)}{2} \cdot \frac{1}{m} = \frac{n}{2m} = \frac{\alpha}{2}.$$

THE RUN TIME OF A SUCCESSFUL SEARCH IS ALSO PROPORTIONAL TO  $\alpha$ .

IT TAKES TIME  $\Theta(1)$  TO COMPUTE  $h(k)$ . THEREFORE THE AVERAGE TIME TO RUN SEARCH IS IN EITHER CASE

$$\Theta(1) + \Theta(\alpha) = \Theta(1 + \alpha).$$

IF IT HAPPENS THAT  $n = \Theta(m)$ , THEN  $\alpha = \frac{n}{m} = \Theta(1)$ , AND SEARCH RUNS IN CONSTANT TIME, ON AVERAGE.

A GOOD HASH FUNCTION WILL SATISFY, AT LEAST APPROXIMATELY, THE ASSUMPTION OF SIMPLE UNIFORM HASHING. NOTE THAT

$$\text{PROB. } (h(k) = i) = \sum_{\{k: h(k)=i\}} P(k)$$

WHERE  $P(k)$  IS THE PROBABILITY THAT KEY  $k$  IS SELECTED. SIMPLE UNIFORM HASHING CAN THEN BE EXPRESSED AS

$$\sum_{\{k: h(k)=i\}} P(k) = \frac{1}{m} \quad (0 \leq i \leq m-1).$$

USUALLY THE PROBABILITY FUNCTION  $P$  IS UNKNOWN, THOUGH SOME QUALITATIVE INFORMATION MAY BE AVAILABLE.

IF IT HAPPENS THAT THE KEYS ARE UNIFORMLY DISTRIBUTED OVER  $U$ , I.E.

$$P(k) = \frac{1}{|U|} \quad (\forall k \in U)$$

THEN ANY FUNCTION WHICH PLACES APPROXIMATELY  $|U|/m$  KEYS IN EACH SLOT WILL SATISFY (APPROXIMATELY) THE SIMPLE UNIFORM HASHING CONDITION.