

3. Hashing Motivation and Terminology:

a) Sequential search of an array or linked list follows the same search pattern for any given target value being searched for, i.e., scans the array from one end to the other, or until the target is found.

If  $n$  is the number of items being searched, what is the average and worst case big-oh notation for a sequential search?

average case  $O(n)$

worst case  $O(n)$

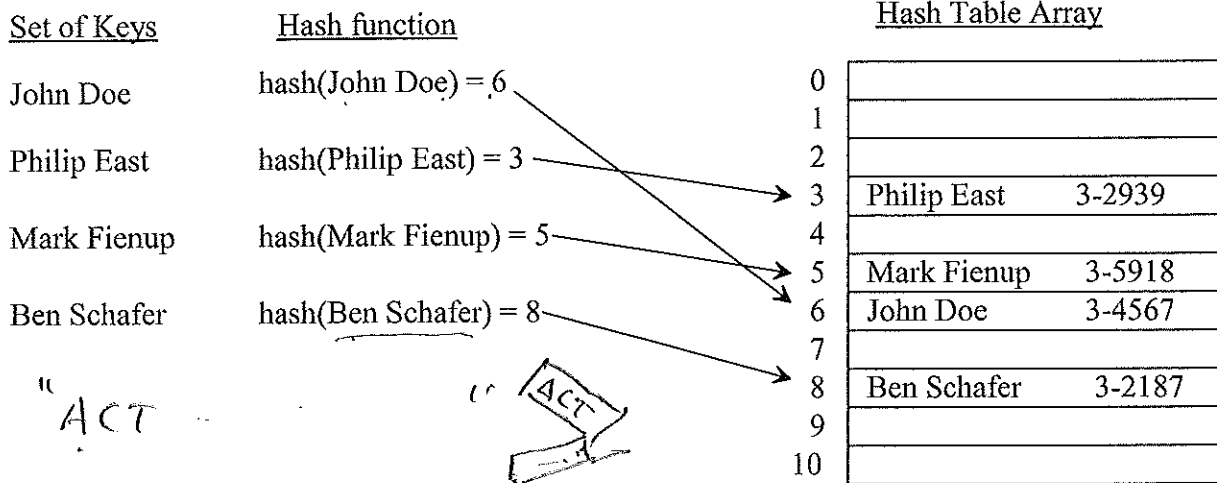
b) Similarly, binary search of a sorted array (or AVL tree) always uses a fixed search strategy for any given target value. For example, binary search always compares the target value with the middle element of the remaining portion of the array needing to be searched.

If  $n$  is the number of items being searched, what is the average and worst case big-oh notation for a search?

average case  $O(\log_2 n)$

worst case  $O(\log_2 n)$

Hashing tries to achieve average constant time (i.e.,  $O(1)$ ) searching by using the target's value to calculate where in the array/Python list (called the *hash table*) it should be located, i.e., each target value gets its own search pattern. The translation of the target value to an array index (called the target's *home address*) is the job of the *hash function*. A *perfect hash function* would take your set of target values and map each to a unique array index.



a) If  $n$  is the number of items being searched and we had a perfect hash function, what is the average and worst case big-oh notation for a search?

average case  $O(1)$

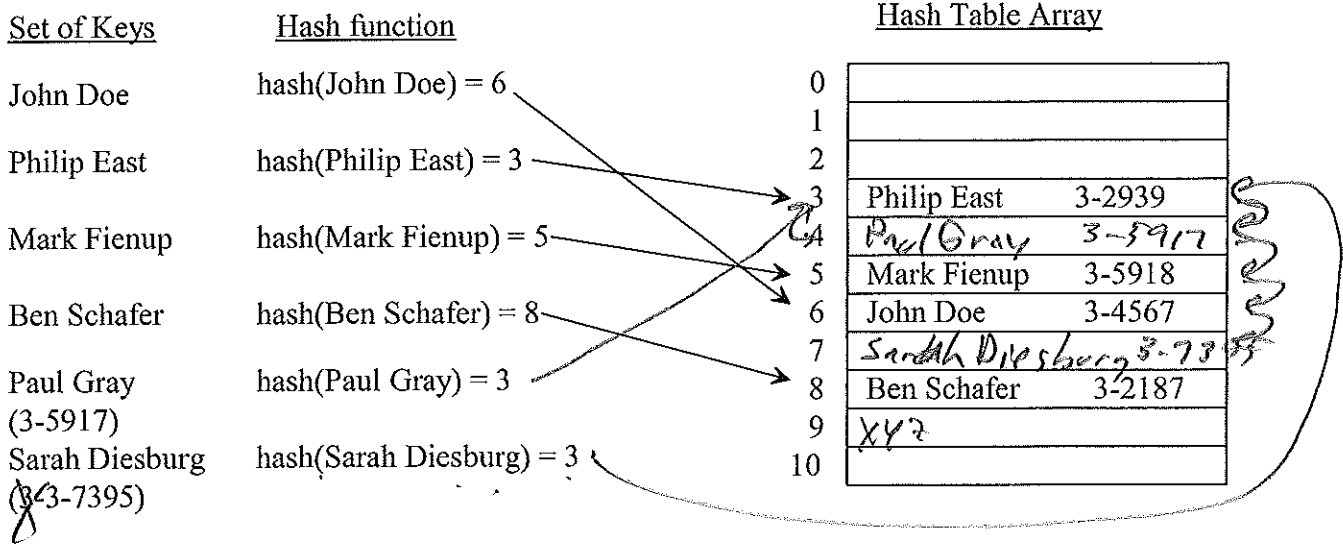
worst case  $O(1)$

4. Unfortunately, perfect hash functions are a rarity, so in general many target values might get mapped to the same hash-table index, called a *collision*.

Collisions are handled by two approaches:

- *open-address* with some *rehashing* strategy: Each hash table home address holds at most one target value. The first target value hashed to a specify home address is stored there. Later targets getting hashed to that home address get rehashed to a different hash table address. A simple rehashing strategy is *linear probing* where the hash table is scanned circularly from the home address until an empty hash table address is found.
- *chaining, closed-address, or external chaining*: all target values hashed to the same home address are stored in a data structure (called a *bucket*) at that index (typically a linked list, but a BST or AVL-tree could also be used). Thus, the hash table is an array of linked list (or whatever data structure is being used for the buckets)

5. Consider the following examples using *open-address* approach with a simple rehashing strategy of *linear probing* where the hash table is scanned circularly from the home address until an empty hash table address is found.

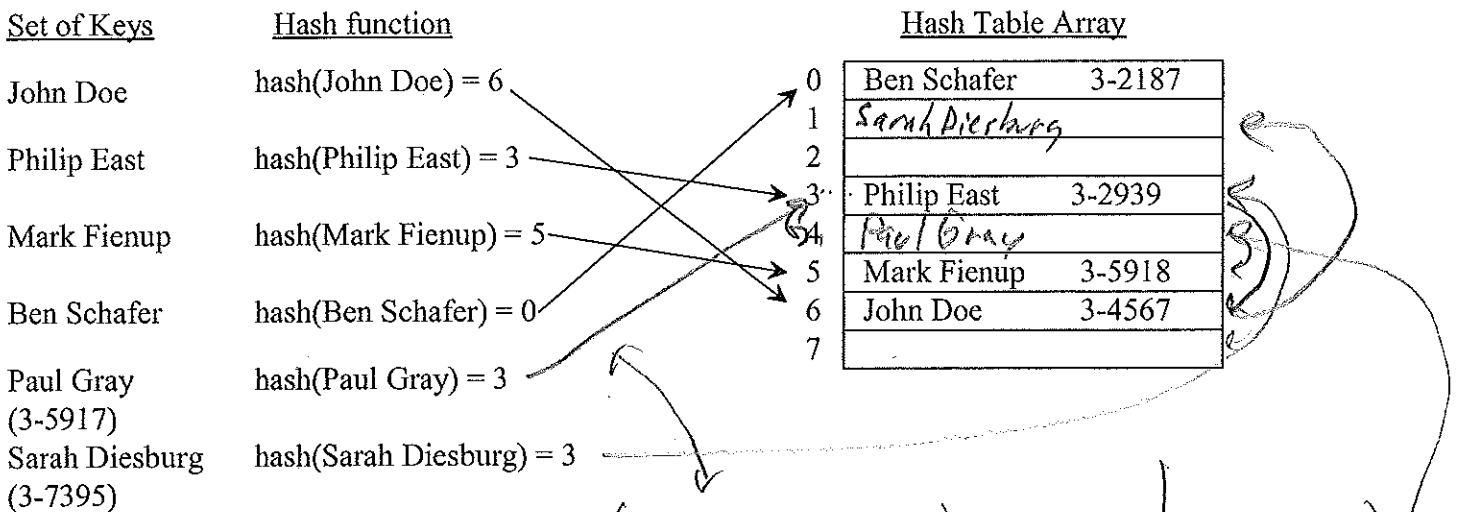


a) Assuming open-address with linear probing where would Paul Gray and then Sarah Diesburg be placed?

Common rehashing strategies include the following.

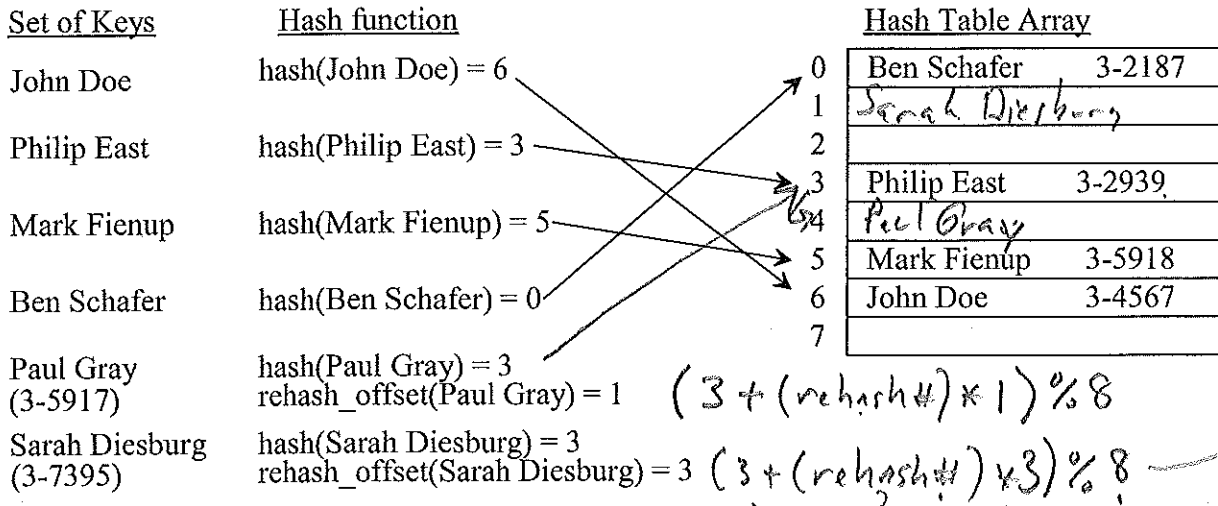
Rehash Strategy	Description
linear probing	Check next spot (counting circularly) for the first available slot, i.e., $(\text{home address} + (\text{rehash attempt \#})) \% (\text{hash table size})$
quadratic probing	Check the square of the attempt-number away for an available slot, i.e., $(\text{home address} + ((\text{rehash attempt \#})^2 + (\text{rehash attempt \#})) / 2) \% (\text{hash table size})$ , where the hash table size is a power of 2
double hashing	Use the target key to determine an offset amount to be used each attempt, i.e., $(\text{home address} + (\text{rehash attempt \#}) * \text{offset}) \% (\text{hash table size})$ , where the hash table size is a power of 2 and the offset hash returns an odd value between 1 and the hash table size

b) Assume quadratic probing, insert "Paul Gray" and "Sarah Diesburg" into the hash table.



$$\begin{aligned}
 & (3 + \frac{1^2 + 1}{2}) \% 8 = 4 \\
 & (3 + \frac{2^2 + 2}{2}) \% 8 = 6 \\
 & (3 + \frac{3^2 + 3}{2}) \% 8 = 1
 \end{aligned}$$

c) Assume double hashing, insert "Paul Gray" and "Sarah Diesburg" into the hash table.



d) For the above double-hashing example, what would be the sequence of hashing and rehashing addresses tried for Sarah Diesburg if the table was full? For the above example, (home address + (rehash attempt #) \* offset) % (hash table size) would be:  $(3 + (\text{rehash attempt \#}) * 3) \% 8$

Rehash Attempt #	0	1	2	3	4	5	6	7	8	9	10
Address	3	6	1	4	7	2	5	0	3		

e) Indicate whether each of the following rehashing strategies suffer from primary or secondary clustering.

- *primary clustering* - keys mapped to a home address follow the same rehash pattern
- *secondary clustering* - rehash patterns from initially different home addresses merge together

Rehash Strategy	Description	Suffers from:	
		primary clustering	secondary clustering
linear probing	Check next spot (counting circularly) for the first available slot, i.e., (home address + (rehash attempt #)) % (hash table size)	Yes	Yes
quadratic probing	Check a square of the attempt-number away for an available slot, i.e., (home address + ((rehash attempt #) <sup>2</sup> + (rehash attempt #)) / 2) % (hash table size), where the hash table size is a power of 2	Yes	No
double hashing	Use the target key to determine an offset amount to be used each attempt, i.e., (home address + (rehash attempt #) * offset) % (hash table size), where the hash table size is a power of 2 and the offset hash returns an odd value between 1 and the hash table size	No	No?

6. Let  $\lambda$  be the load factor (# item/hash table size). The average probes with **linear probing** for insertion or unsuccessful search is:  $(\frac{1}{2})(1 + (\frac{1}{(1-\lambda)^2}))$ . The average for successful search is:  $(\frac{1}{2})(1 + (\frac{1}{(1-\lambda)}))$ .

a) Why is an unsuccessful search worse than a successful search?

With successful search we can stop when we find the target value, but unsuccessful search needs to continue until an empty spot found.

The average probes with **quadratic probing** for insertion or unsuccessful search is:  $\left(\frac{1}{1-\lambda}\right) - \lambda - \log_e(1-\lambda)$

The average probes with quadratic probing for successful search is:  $1 - \left(\frac{\lambda}{2}\right) - \log_e(1-\lambda)$

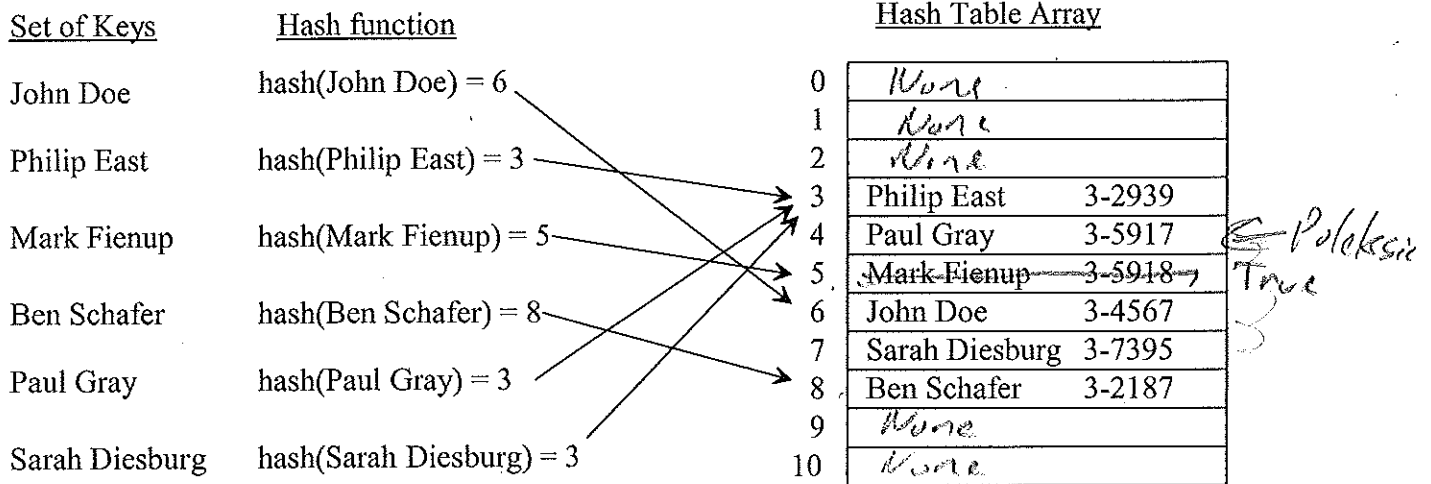
Consider the following table containing the average number probes for various load factors:

Probing Type	Search outcome	Load Factor, $\lambda$				
		0.25	0.5	0.67	0.8	0.99
Linear Probing	unsuccessful	1.39	2.50	5.09	13.00	5000.50
	successful	1.17	1.50	2.02	3.00	50.50
Quadratic Probing	unsuccessful	1.37	2.19	3.47	5.81	103.62
	successful	1.16	1.44	1.77	2.21	5.11

b) Why do you suppose the "general rule of thumb" in hashing tries to keep the load factor between 0.5 and 0.67?

Less than load factor of 0.5, we waste a lot of space since hash table is < half full. Greater than 0.67 load factor we start to get many probes.

7. Allowing deletions from an open-address hash table complicates the implementation. Assuming linear probing we might have the following



a) If "Mark Fienup" is deleted, how will we find Sarah Diesburg? need to somehow know to keep search past the deleted value.

b) How might we fix this problem?

Two markers for unused slots: - empty = None value  
- deleted = True value

$$\lambda = \frac{(\# \text{ items} + \# \text{ deleted})}{\text{hash table size}}$$

searching