

1. Consider the following sequential search (linear search) code:

Textbook's Listing 5.1	Faster sequential search code
<pre>def sequentialSearch(alist, item): """ Sequential search of unordered list """ pos = 0 found = False while pos < len(alist) and not found: if alist[pos] == item: found = True else: pos = pos+1 return found</pre>	<pre>def linearSearch(aList, target): """Returns the index of target in aList or -1 if target is not in aList""" for position in range(len(aList)): if target == alist[position]: return position return -1</pre>

a) What is the *basic operation* of a search? *comparison of items' target with list item.*

b) For the following aList value, which target value causes linearSearch to loop the fewest ("best case") number of times? *10*

	0	1	2	3	4	5	6	7	8	9	10	
aList:	10	15	28	42	60	69	75	88	90	93	97	<i>O(1)</i> <i>best case</i>

c) For the above aList value, which target value causes linearSearch to loop the most ("worst case") number of times? *unsuccessful O(n)*

d) For a *successful search* (i.e., target value in aList), what is the "average" number of loops? *O(n/2)*
O(n)

Textbook's Listing 5.2	Faster sequential search code
<pre>def orderedSequentialSearch(alist, item): """ Sequential search of order list """ pos = 0 found = False stop = False while pos < len(alist) and not found and not stop: if alist[pos] == item: found = True else: if alist[pos] > item: stop = True else: pos = pos+1 return found</pre>	<pre>def linearSearchOfSortedList(target, alist): """Returns the index position of target in sorted alist or -1 if target is not in alist""" breakOut = False for position in range(len(aList)): if target <= alist[position]: breakOut = True break if not breakOut: return -1 elif target == alist[position]: return position else: return -1</pre>

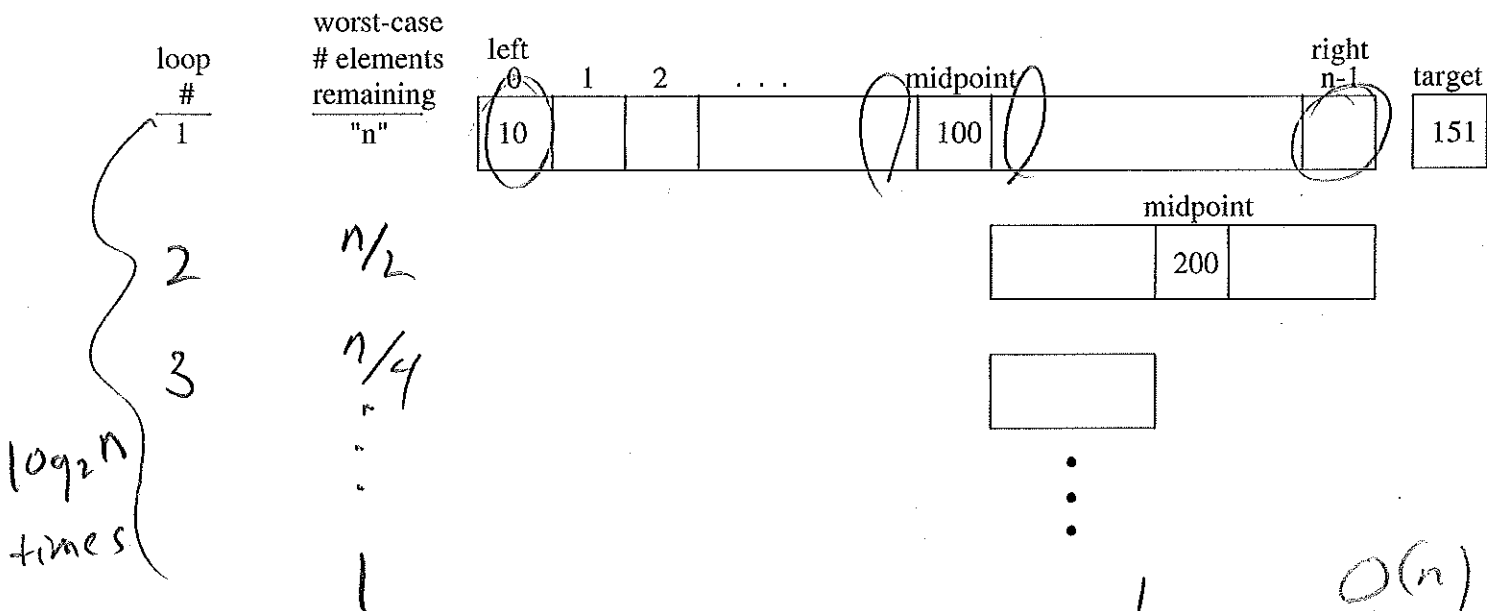
e) The above version of linear search assumes that aList is sorted in ascending order. When would this version perform better than the original linearSearch at the top of the page?

Stops early if run across list item that's bigger than target item.

2. Consider the following binary search code:

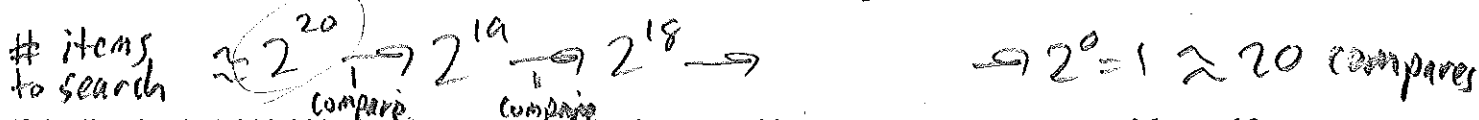
Textbook's Listing 5.3	Faster binary search code
<pre>def binarySearch(alist, item): first = 0 last = len(alist)-1 found = False while first<=last and not found: midpoint = (first + last)//2 if alist[midpoint] == item: found = True else: if item < alist[midpoint]: last = midpoint-1 else: first = midpoint+1 return found</pre>	<pre>def binarySearch(target, lyst): """Returns the position of the target item if found, or -1 otherwise.""" left = 0 right = len(lyst) - 1 while left <= right: midpoint = (left + right) // 2 if target == lyst[midpoint]: return midpoint elif target < lyst[midpoint]: right = midpoint - 1 else: left = midpoint + 1 return -1</pre>

a) "Trace" binary search to determine the worst-case basic total number of comparisons?



- b) What is the worst-case big-oh for binary search? $O(\log_2 n)$ exact time
- c) What is the best-case big-oh for binary search? $O(1)$
- d) What is the average-case (expected) big-oh for binary search? $O(\log_2 n)$

e) If the list size is 1,000,000, then what is the maximum number of comparisons of list items on a successful search?



f) If the list size is 1,000,000, then how many comparisons would you expect on an unsuccessful search?

≈ 20

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 $O(n/2) = O(n)$

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 $O(\log n)$

worst case $O(\log n)$

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.

<u>Set of Keys</u>	<u>Hash function</u>	<u>Hash Table Array</u>
John Doe	hash(John Doe) = 6	0
Philip East	hash(Philip East) = 3	1
Mark Fienup	hash(Mark Fienup) = 5	2
Ben Schafer	hash(Ben Schafer) = 8	3
		4
		5
		6
		7
		8
		9
		10

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 $O(1)$

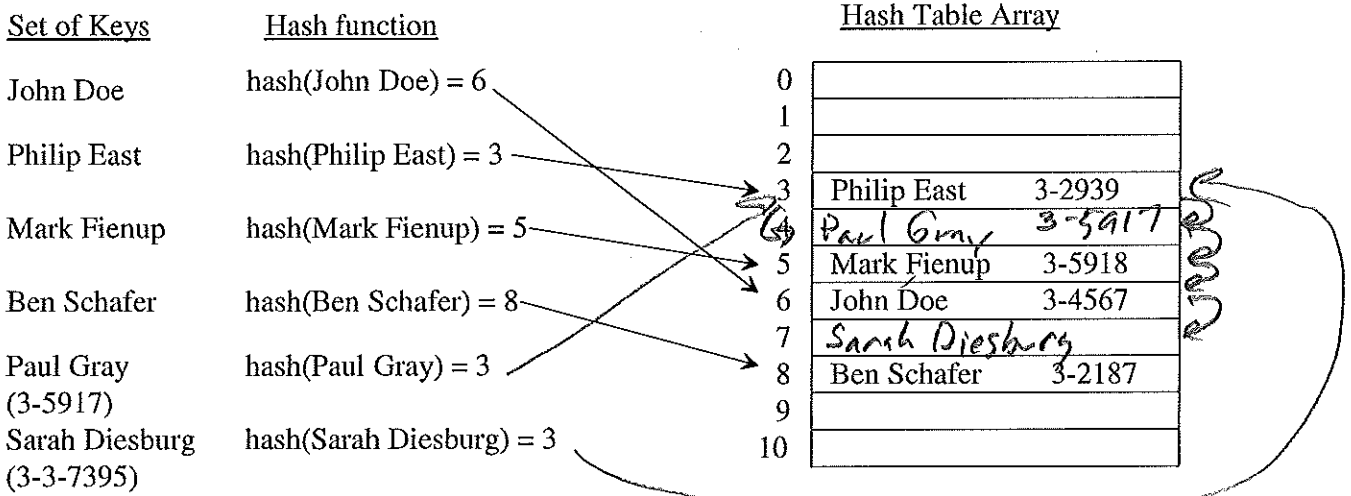
worst case $O(1)$

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 specific 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.



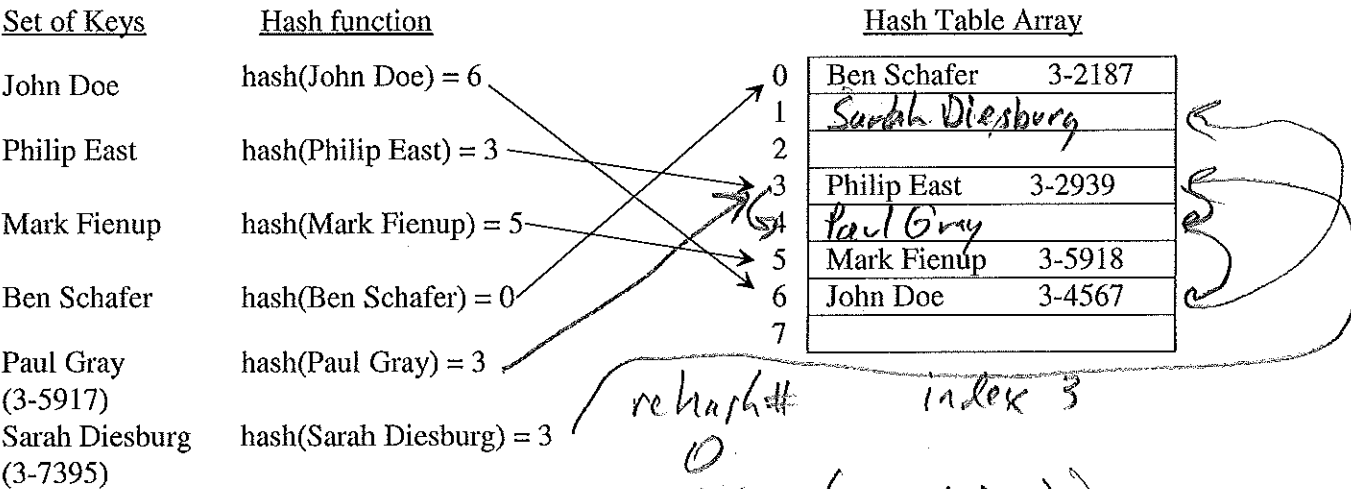
a) Assuming open-address with linear probing where would Paul Gray and then Sarah Diesburg be placed?

Common rehashing strategies include the following.

$11 \% 11 = 0$
 $10 \% 11 = 10$
 $8 \% 11 = 8$
 $7 \% 11 = 7$

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.



rehash# index 3

0

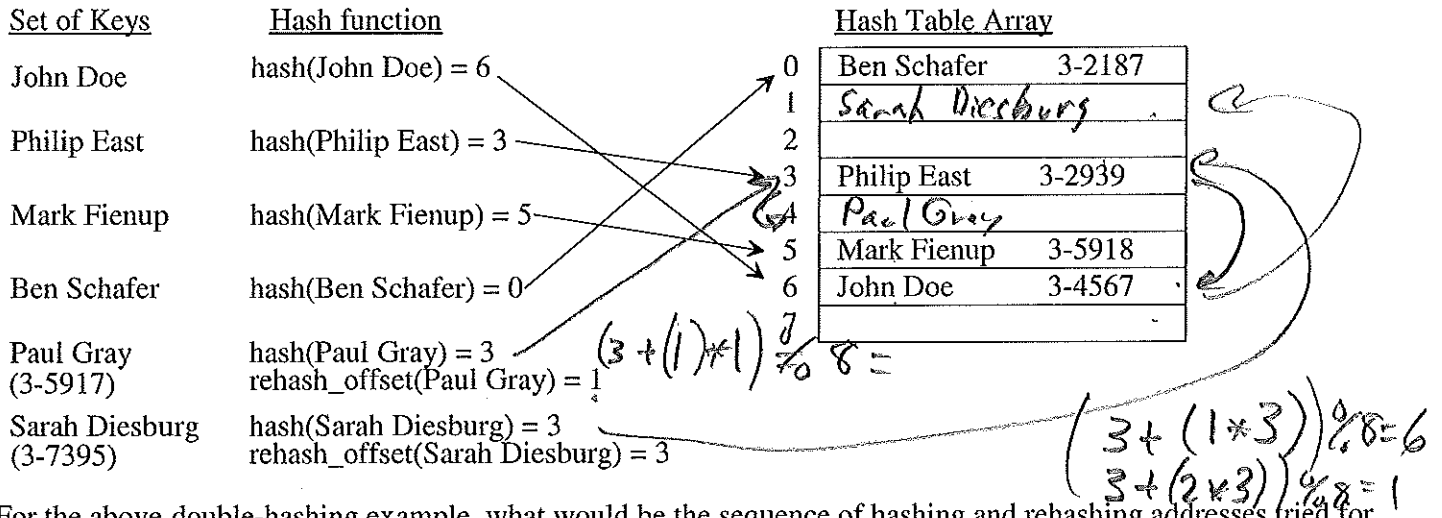
1 $(3 + \frac{(1^2 + 1)}{2}) \% 8 = 4$

2 $(3 + \frac{(2^2 + 2)}{2}) \% 8 = 6$

3 $(3 + \frac{(3^2 + 3)}{2}) \% 8 = 1$

4 $(3 + \frac{(4^2 + 4)}{2}) \% 8 = 5$

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 #) ² + (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 λ 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? *can spot when find target*

Consider $\lambda = 0.5$

unsuccessful: $(\frac{1}{2})(1 + (\frac{1}{(1-0.5)^2}))$
 $= \frac{1}{2}(1 + \frac{1}{0.25})$
 $= \frac{1}{2}(1 + 4) = 2.5$

successful: $(\frac{1}{2})(1 + (\frac{1}{(1-0.5)}))$
 $= (\frac{1}{2})(1 + \frac{1}{0.5})$
 $= \frac{1}{2}(1 + 2)$
 $= 1.5$

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

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

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

Probing Type	Search outcome	Load Factor, λ				
		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

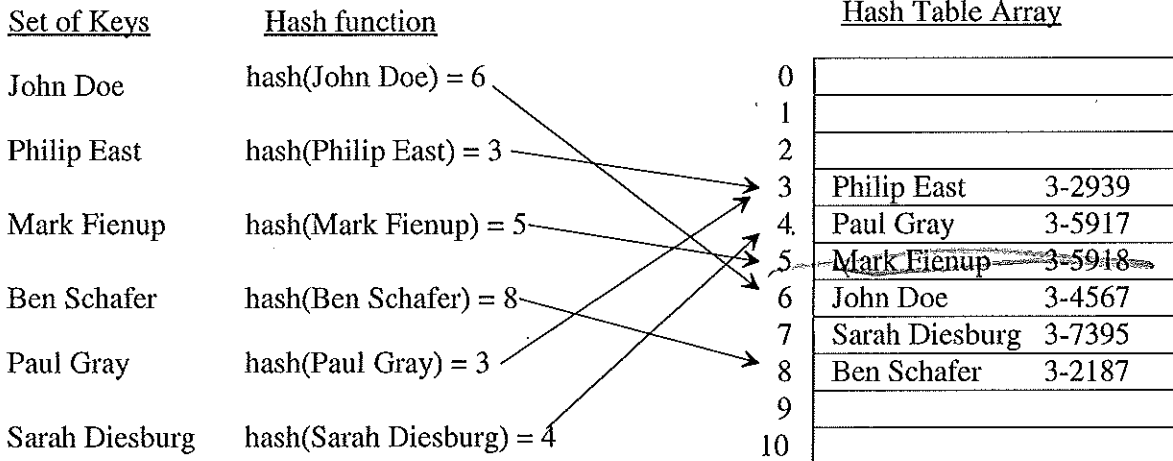
b) Why do you suppose the "general rule of thumb" in hashing tries to keep the load factor between 0.5 and 0.67?

waste memory

"baby bears just right"

starts to get slow

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?

b) How might we fix this problem? Two special "values": empty \equiv None

(Periodically, rehash items to empty hash table to eliminate "deleted spot" entries.)

deleted spot \equiv True

Successful search: load factor = (#items + deleted) / hash table

Unsuccessful search: =

insertion load factor = (#items / hash table)