

CSCB63 – Design and Analysis of Data Structures

Anya Tafliovich¹

¹based on notes by Anna Bretscher and Albert Lai

dictionaries (again)

Recall that a dictionary is an ADT that supports the following operations on a set of elements with well-ordered key-values $k-v$:

1. `insert(k , v)`: insert new key-value pair $k-v$
2. `delete(k)`: delete the node with key k
3. `search(k)`: find the node with key k (or value associated with key k)

Q. If we know the keys are integers from 1 to K , what is a fast and simple way to represent a dictionary?

A. Use array of size K , store array $[i]$ at index $i-1$

This data structure is called direct addressing.

Q. What is the asymptotic worst-case time for each of the important operations?

A. $\Theta(1)$

direct addressing

Q. What may be a problem with direct addressing?

A. Space not enough

Example 1: Reading a text file

Suppose we want to keep track of the frequencies of each letter in a text file.

Q: Why is this a good application of direct addressing?

A. Example 2: Reading a data file of 32-bit integers

Suppose we want to keep track of the frequencies of each number.

Q: Is this a good or bad application of direct addressing?

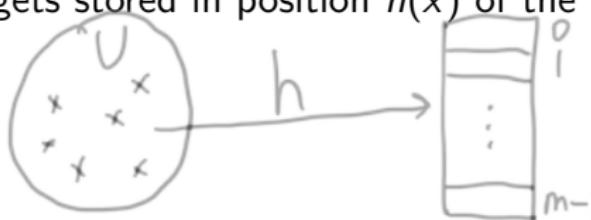
A. array size 2^{32} , too large

hashing: idea

- the range of keys is large
- but many keys are not “used”
- don’t need to allocate space for all possible keys

A hash table:

- if keys come from a universe (set) U
- allocate a table (an array) of size m (where $m < |U|$)
- use a hash function $h : U \rightarrow \{0, \dots, m - 1\}$ to decide where to store the element with key x
- x gets stored in position $h(x)$ of the hash table



hashing: problem

If $m < |U|$, then there must be $k_1, k_2 \in U$ such that $k_1 \neq k_2$ and yet $h(k_1) = h(k_2)$.

This is called a collision.

How we deal with collisions is called collision resolution. When we study hashing, we mostly study collision resolution.

collision resolution: idea

Say we have a small address book and one of the letters fills up, for example, "N"s. Where do you add the next "N" entry?

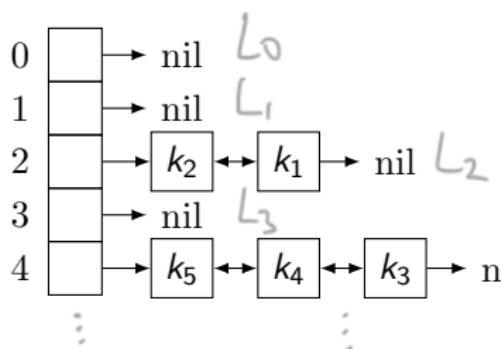
- flip to the next page
- have an overflow page at the very end
- write a little note explaining where to find rest of the "N" names

Two general collision resolution approaches:

1. Closed Addressing: Keys are always stored in the bucket they hash to — use additional data structure to store the keys in the same bucket.
2. Open Addressing: Give a general rule of where to look next (directions to another bucket).

closed addressing: chaining

Idea: store a doubly linked list at each entry in the hash table



$\text{insert}(k, v)$:

① compute $h(k)$
 ↓
 index i

② prepend (k, v) to L_i

An element with key k_1 and an element with key k_2 can both be stored at position $h(k_1) = h(k_2)$.

$\Theta(1)$ complexity

This is called chaining.

chaining: complexity

- Assume we can compute the hash function h in constant time.
- $\text{insert}(k, v)$ takes:
- $\text{delete}(k)$ takes:
 -
 -
- $\text{search}(k)$ takes:

chaining: worst case

Q. What happens if $|U| > m * n$?

A.

Q. What is the worst case?

A.

simple uniform hashing

We assume hash function h has the simple uniform hashing property:

- any element is equally likely to hash into any of m buckets
 - independently of where any other element has hashed to, and
- h distributes elements of U evenly across m buckets
- formally:
 - sample space: set of elements with key-values from U
 - for any probability distribution on U
 - - $\Pr(h(k) = i) = \frac{1}{m}$ for all $1 \leq i \leq m, k \in U$ and
 - $\sum_{k \in U_i} \Pr(k) = \frac{1}{m}$ where $U_i = \{k \in U \mid h(k) = i\}$


$$\Pr(h(k) = i) = \frac{1}{m} \text{ for all } 1 \leq i \leq m, k \in U \text{ and}$$
$$\sum_{k \in U_i} \Pr(k) = \frac{1}{m} \text{ where } U_i = \{k \in U \mid h(k) = i\}$$

load factor

Q. If the table has n elements, how many would you expect in any one entry of the table ?

A.

- We call this ratio n/m the load factor, denoted by α .
- This simple uniform hashing assumption may or may not be accurate depending on U , h and the probability distribution for $k \in U$.

average case analysis

Calculating the average-case run time:

- Let T_k be a random variable which counts the number of elements checked when searching for key k .
- Let L_i be the length of the list at entry i in the hash table.
- Either we are searching for an item in the table or not in the table.

search(k):

- ① compute $h(k) \Rightarrow$ returns i
- ② search through L_i

then $T_k = \# \text{ of comparisons} = \# \text{ of elements of } L_i \text{ checked}$

average case analysis: unsuccessful search

$$E(T) = \sum_{k \in U} P(k) \cdot T_k$$

$$= \sum_{i=1}^n \left[\sum_{k \in U_i} P(k) \right] \cdot T_k$$

$$= \sum_{i=1}^n \frac{1}{m} \cdot |L_i| \rightarrow \# \text{ of elements in } L_i$$

$$= \frac{1}{m} \sum_{i=1}^n |L_i| = \frac{n}{m} = \alpha \text{ "load factor"}$$

average case analysis: successful search

- Suppose we are searching for any of the n elements in the hash table, with equal probability, $1/n$.
- The number of elements examined before we reach the element x we are looking for is determined by the number of elements inserted **after** x . $\Leftrightarrow \# \text{ of elements in front of } k$
- Expected number of elements examined is:

Let:

- k_1, k_2, \dots, k_n : keys inserted, in order
- X_{ij} indicator variable of event that $h(k_i) = h(k_j)$
- then $E[X_{ij}] = P(X_{ij}=1) = P(h(k_i)=h(k_j))$

average case analysis: successful search

$$= \sum_{s=1}^m P(h(k_i) = s \wedge h(k_j) = s)$$

$$= \sum_{s=1}^m [P(h(k_i) = s) \cdot P(h(k_j) = s)] \text{ by independence of } k_i, k_j$$

$$= \sum_{s=1}^m \frac{1}{m} \cdot \frac{1}{m} = \frac{1}{m}$$

average case analysis: successful search

$$E(T) = \frac{1}{n} \cdot \sum_{i=1}^n \left[1 + E\left(\sum_{j=i+1}^n X_{ij}\right) \right]$$

probability of k_i looking
at k_i looking at all
other

$$= \frac{1}{n} \sum_{i=1}^n \left[1 + \sum_{j=i+1}^n E(X_{ij}) \right]$$

$$= \frac{1}{n} \sum_{i=1}^n \left[1 + \sum_{j=i+1}^n \frac{1}{m} \right]$$

$$= \frac{1}{n} \sum_{i=1}^n \left(1 + \frac{n-i}{m} \right) = 1 + \frac{\alpha}{2} - \frac{\alpha}{2n}$$

average case of search — closed addressing

- So the average-case running time of search under simple uniform hashing with chaining is $\Theta(1 + \alpha)$.
- If the number of slots is proportional to number of elements in the table, then n is $\mathcal{O}(m)$ and so search takes constant time on average.

open addressing

- Each entry in the hash table stores a fixed number c of elements.
- This has the immediate implication that we only use it when $n \leq cm$.
- We will keep c at 1 for today's class.

To insert a new element if we get a collision:

- Find a new location to store the new element.
- We need to know where we put it: for future retrieval.
- Search a well-defined sequence of other locations in the hash table, until we find one that's not full.

This sequence is called a probe sequence.

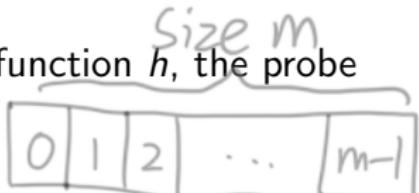
probe sequences

Many methods for generating a probe sequence. For example:

- linear probing: try $A[(h(k) + i) \bmod m]$, $i = 0, 1, 2, \dots$
- quadratic probing: try $A[(h(k) + c_1i + c_2i^2) \bmod m]$
- double hashing: try $A[(h(k) + i \cdot h'(k)) \bmod m]$
where h' is another hash function

linear probing

For a hash table of size m , key k and hash function h , the probe sequence is calculated as:



$$s_i = (h(k) + i) \bmod m \quad \text{for } i = 0, 1, 2, \dots$$

- $s_0 = h(k)$ is called the home location for the item
- the problem:
 - division method (obtain remainders)
- when we hash to a location within a group of filled locations
 -
 - - initial location: $S_0 = h(k) \bmod m$
 - if S_0 filled $S_1 = (h(k)+1) \bmod m$
 - if S_1 filled $S_2 = (h(k)+2) \bmod m$
 - ⋮

non-linear probing

Idea: the probe sequence does not involve steps of fixed size.

Example: Quadratic probing is where the probe sequence is calculated as:

$$s_i = (h(k) + c_1 i + c_2 i^2) \bmod m \quad \text{for } i = 0, 1, 2, \dots$$

But: Usually $s_i = (h(k) + i^2) \bmod m \quad i=0,1,2,\dots$

initial location $s_0 = h(k) \bmod m$

if s_0 filled $s_1 = (h(k) + 1^2) \bmod m$

if s_1 filled $s_2 = (h(k) + 2^2) \bmod m$

⋮

Method: draw table for key k , $h(k)$, $h'(k)$:

k	$h(k)$	$h'(k)$
element 1		
element 2		

- In double hashing we use a different hash function $h_2(k)$ to calculate the step size.

initial location: $s_0 = h(k) + 0 \cdot h'(k)$

- The probe sequence is:

if s_0 filled: $s_1 = h(k) + 1 \cdot h'(k) = h(k)$

if s_1 filled: $s_2 = h(k) + 2 \cdot h'(k)$

$$s_i = (h(k) + i \cdot h'(k)) \bmod m \quad \text{for } i = 0, 1, 2, \dots$$

- Note that $h'(k)$ should not be 0 for any k .
- Also, we want to choose h' so that, if $h(k_1) = h(k_2)$ for two keys k_1, k_2 , it won't be the case that $h'(k_1) = h'(k_2)$.
- That is, the two hash functions don't cause collisions on the same pairs of keys.

open addressing: complexity

- consider the complexity of $\text{search}(k)$
- worst case scenario?
-

Suppose:

- the hash table has m locations
- the hash table contains n elements and $n < m$
- we search for a random key k in the table, with probability $\frac{1}{n}$

Consider a random probe sequence for k :

- probe sequence is equally likely to be any permutation of $\langle 0, 1, \dots, m - 1 \rangle$

open addressing: unsuccessful search

Let T be the number of probes performed in an **unsuccessful search**.

$$\begin{aligned} \text{Then } E(T) &= \sum_{i=0}^{\infty} i \cdot P(T=i) \\ &= \sum_{i=0}^{\infty} i \cdot [P(T \geq i) - P(T \geq i+1)] \\ &= \sum_{i=1}^{\infty} P(T \geq i) \end{aligned}$$

$T \geq i$: need to have at least $i+1$ probes to find empty slot

open addressing: unsuccessful search

Let A_i denote the event that the i -th probe occurs and it is to an occupied slot.

$P(A_i) = \text{prob. that } S_i \text{ is not empty at location } i$

Then, $T \geq i$ iff A_1, A_2, \dots, A_{i-1} all occur. *for uniform hashing function h*

$$\Pr(T \geq i) = P(A_1 \cap A_2 \cap \dots \cap A_{i-1})$$

$$= P(A_1) \cdot P(A_2 | A_1) \dots P(A_{i-1} | A_1 \cap A_2 \cap \dots \cap A_{i-2})$$

$$= \frac{n-1+1}{m-1+1} \cdot \frac{n-2+1}{m-2+1} \dots \frac{n-(i-1)+1}{m-(i-1)+1}$$

$$\leq \frac{n}{m} \cdot \frac{n}{m} \dots \frac{n}{m} = \left(\frac{n}{m}\right)^{i-1} = \alpha^{i-1}$$

open addressing: unsuccessful search

$$Pr(A_j | A_1 \cap \dots \cap A_{j-1}) = ?$$

e.g.

0	1	2	3	4	5
---	---	---	---	---	---

keys:

(15, 23, 7, 9)

Intuition:

- number of elements we have not yet seen: $n - (j-1)$
- number of slots we have not yet seen: $m - (j-1)$

Math: for $1 \leq j \leq m$:

$$P = \frac{4}{6} \cdot \frac{3}{5} \cdot \frac{2}{4} \cdot \frac{1}{3}$$

$$P(A_j | A_1 \cap \dots \cap A_{j-1})$$

$$= \frac{\# \text{elements not yet seen}}{\# \text{slots not yet seen}} = \frac{n-j+1}{m-j+1}$$

open addressing: unsuccessful search

Now we can calculate the expected value of T , or the average-case complexity of unsuccessful search(k).

$$E(T) =$$

$$\sum_{i=0}^{\infty} P(T \geq i)$$

$$= \sum_{i=0}^m P(T \geq i) + \sum_{i=m+1}^{\infty} P(T \geq i) \rightarrow 0$$

$$\leq \sum_{i=1}^m d^{i-1} \leq \sum_{j=0}^{\infty} d^j = \frac{1}{1-d}$$

open addressing: insert

To insert a new element:

- perform an unsuccessful search (for an available location)
- insert:

Thus, $\text{insert}(k, v)$ requires at most $\frac{1}{1-\alpha}$ probes on average.

open addressing: successful search

Let T be the number of probes performed in a **successful search**.

Idea: successful search(k) reproduces the same probing sequence as insert(k, v).

i unsuccessful search, $(i+1)^{st}$ successful

If k was the $(i + 1)^{st}$ key inserted into the table, then the expected number of probes made is at most

$$\frac{1}{T-i} = \frac{1}{T-\frac{i}{m}} \quad (\text{from previous result})$$

Then, averaging over all n keys in the table:

$$E(T) = \frac{1}{n} \cdot \sum_{i=0}^{n-1} \frac{1}{T-\frac{i}{m}}$$

open addressing: successful search

$$E(T) =$$

open addressing: successful search

$$E(T) \leq \frac{1}{\alpha} \ln \frac{1}{1 - \alpha}$$

This is pretty good!

- if the table is half full, the expected number of probes is < 1.387
- if the table is 90% full, this number is < 2.559

open addressing: delete

What about delete?

- with closed addressing: easy
 - first do search then
 - $\mathcal{O}(1)$ un-link
- with open addressing: two approaches
 - find an existing key to fill the hole
 -
 - over time slows down all operations
 - delete is problematic under open addressing