

CPTS 223 Advanced Data Structure C/C++

Hashing

Overview

- Hash Table Data Structure : Purpose
 - To support insertion, deletion and search in average-case constant time
 - Assumption: Order of elements irrelevant
 - \rightarrow data structure **NOT** useful for if you want to maintain and retrieve an order of the elements

Average-case O(1)

- Hash function
 - Hash["string key"] \rightarrow integer value
- Hash table ADT
 - Implementations, analysis, applications
- Collision in hash tables

Hash table



Hash table

- Hash table is an array of fixed size TableSize
- Array elements indexed by a key, which is mapped to an array index (0...TableSize-1)
- Mapping (hash function) h from key to index
 - e.g., h("john") = 3-





Figure 5.1 An ideal hash table

Hash table operations

- hash function
 Insert
 T [h("john")] = <"john",25000>
- Delete
 - T [h(``john'')] = NULL
- Search
 - T [h("john")] returns the element hashed for "john"



0	
1	
2	
3	john 25000
4	phil 31250
5	
6	dave 27500
7	mary 28200
8	
9	

Figure 5.1 An ideal hash table

Factors of hash table design

- Hash function
- Table size
 - Usually fixed at the start
- Collision handling scheme

0	
1	
2	
3	john 25000
4	phil 31250
5	
6	dave 27500
7	mary 28200
8	
9	

Figure 5.1 An ideal hash table

Hash function

- A hash function is one which maps an element's key into a valid hash table index
 - $h(key) \rightarrow hash table index$
- Note that this is different from saying:
 - h(string) \rightarrow int
 - Because the key can be of any type
 - e.g., "h(int) => int" is also a hash function

Hash function properties

h(key) \rightarrow hash able index

- A hash function maps key to integer
 - Constraint: Integer should be between:
 [o, TableSize-1]



- A hash function can result in a many-to-one mapping (causing collision)
 - Collision occurs when hash function maps two or more keys to same array index
- Collisions cannot be avoided
- But its chances can be reduced using a "good" hash function

Hash function properties

- A "good" hash function should have the following properties:
 - Reduced chance of collision
 - Different keys should ideally map to different indices
 - Distribute keys uniformly over table
 - Should be fast to compute

Effective table size

- Simple hash function (assume integer keys)
 - h(Key) = Key mod TableSize
- For random keys, h() distributes keys evenly over table
 - What if TableSize = 100 and keys are ALL multiples of 10?
 - Better if TableSize is a prime number

Hash function for string keys

- Approach 1: A very simple function to map strings to integers
 - Add up character ASCII values (0-255) to produce integer keys
 - e.g., "abcd" = 97+98+99+100 = 394
 - → h("abcd") = 394 % TableSize
- Potential problems:
 - Anagrams will map to the same index
 - h("abcd") == h("dbac")
 - Small strings may not use all of table
 - Strlen(S) * 255 << TableSize
 - Time proportional to length of the string

Hash function for string keys

- Approach 2: treat first 3 characters of string as base-27 integer (26 letters plus space)
 - Key = S[0] + (27 * S[1]) + (272 * S[2])
- Potential problems:
 - Assumes first 3 characters randomly distributed?
 - Not true of English





Hash function for string keys

- Approach 3: use all N characters of string as an N-digit base-K number
- Choose K to be prime number larger than number of different digits (characters)
 - i.e., K = 29, 31, 37
- If L = length of string S, then
- $h(S) = \left[\sum_{i=0}^{L-1} S(L-1-i) * 37^{i}\right] \mod 1$ **TableSize**
- Use Horner's rule to compute h(S)
- Limit L for long strings

	1	/**					
	2	* A hash routine for string objects.					
	3	*/					
	4	unsigned int hash(const string & key, int tableSize)					
er	5	{					
	6	unsigned int hashVal = 0;					
	7						
	8	for(char ch : key)					
	9	hashVal = 37 * hashVal + ch;					
	10						
	11	return hashVal % tableSize;					
ho	12	}					
Ju	Figu	Jre 5.4 A good hash function					
Dala	_						
Polynomial function of 37							
Degree = L-1							
Prob	lems:	potential overflow larger runtime					

Techniques for collisions

With linked lists

Open addressing:

without linked lists

• Separate chaining

- Probing
- Double hashing

Techniques for collisions

- What happens when $h(k_1) = h(k_2)$?
 - \rightarrow collision
- Collision resolution strategies
 - Separate chaining
 - Store colliding keys in a linked list at the same hash table index
 - Open addressing
 - Store colliding keys elsewhere in the table
- Rehashing: when the hash table is too full
 - If too full
 - hash table operations take too long
 - insertions might fail for open addressing

Separate chaining

Insertion sequence: { 0, 1, 4, 9, 16, 25, 36, 49, 64, 81 }

- Strategy: maintains a linked list at every hash index for collided elements
- Hash table T is a vector of linked lists •
 - Insert element at the head (as shown here) or at the tail
- Key k is stored in list at T[h(k)]
- e.g., TableSize = 10
 - h(k) = k mod 10 •
 - Insert first 10 perfect squares •



Figure 5.5 A separate chaining hash table 16

Separate chaining: declaration



Figure 5.6 Type declaration for separate chaining hash table

Separate chaining: hash function



template <>
class hash<string>

};

```
public:
    size_t operator()( const string & key )
    {
```

```
size_t hashVal = 0;
```

```
for( char ch : key )
hashVal = 37 * hashVal + ch;
```

```
return hashVal;
```







Figure 5.10 insert routine for separate chaining hash table

Separate chaining



Separate chaining: remove



Separate chaining: analysis

- Load factor λ of a hash table T is defined as follows:
 - N = number of elements in T
 - Current size
 - M = size of T
 - TableSize
 - $\lambda = N/M$
 - load factor
 - i.e., λ is the average length of a chain
- Search time: $O(1 + \lambda)$
- Ideally, want $\lambda \leq 1$

Separate chaining: disadvantage

- Linked lists could get long
 - Especially when N approaches M
 - Longer linked lists could negatively impact performance
- More memory because of pointers
- Absolute worst-case (even if N << M):
 - All N elements in one linked list!
 - Typically the result of a bad hash function

Open addressing

- When a collision occurs, look elsewhere in the table for an empty slot
- Advantages over chaining
 - No need for list structures
 - No need to allocate/deallocate memory during insertion/deletion (slow)
- Disadvantages
 - Slower insertion May need several attempts to find an empty slot
 - Table needs to be bigger (than chaining-based table) to achieve average-case constant-time performance
 - Load factor $\lambda \approx 0.5$ _____ Load factor = N/M ____

N: # of elements in hash table M: size of T (TableSize)

Open addressing

- A "probe sequence" is a sequence of slots in hash table while searching for an element x
 - h_o(x), h_1(x), h_2(x), ...
 - Needs to visit each slot exactly once
 - Needs to be repeatable (so we can find/delete what we have inserted already)
- Hash function (f(i) can be linear, quadratic, etc.
 - $h_i(x) = (h(x) + f(i)) \mod TableSize$
 - $f(o) = o \rightarrow position for the o-th probe$
 - f(i) is "the distance to be traveled relative to the o-th probe position, during the i-th probe".

Open addressing: linear probing

- oth probe occupied 1st probe occupied 2nd probe occupied grd probe unoccupied insesrt x here
- f(i): a linear function of i
- e.g., f(i) = i
- h_i(x) = (h(x) + f(i)) mod TableSize
 = (h(x) + i) mod TableSize
- Probe sequence: +0, +1, +2, ...
- Continue until an empty slot is found
- #failed probes is a measure of performance (smaller \rightarrow better)

Open addressing: linear probing

- f(i): a linear function of i, e.g., f(i)=i
 - hi(x) := (h(x) + i) mod TableSize
 - Probe sequence: +0, +1, +2, +3, +4, ...
- Example: h(x) = x mod TableSize
 - ho(89) = (h(89)+f(0)) mod 10 = 9
 - ho(18) = (h(18)+f(o)) mod 10 = 8
 - ho(49) = (h(49)+f(0)) mod 10 = 9 (X)
 - h1(49) = (h(49)+f(1)) mod 10

= (h(49)+1) mod 10 = 0

Hash table at index 9 has been occupied by 89 (the first insert)

Open addressing: linear probing Insert sequence: 89, 18, 49, 58, 69 After 89 After 18 After 49 After 58 After 69 Empty Table Totally 7 #unsuccessful probes: Hash table with linear probing, after each insertion Figure 5.11

Linear probing: issues

- Probe sequences can get longer with time
- Primary clustering (blocks of occupied cells)
 - If probe occurs, keys tend to cluster in one part of table
 - Keys that hash into cluster will be added to the end of the cluster (making it even bigger)
 - Side effect: Other keys could also get affected if mapping to a crowded neighborhood

Linear probing: analysis

• Expected number of probes for insertion or unsuccessful search

 $\frac{1}{2}(1+1/(1-\lambda)^2)$

Target element is not present in hash table

 Expected number of probes for successful search _____

 $\frac{1}{2}(1+1/(1-\lambda))$

Target element is present in hash table

- Example ($\lambda = 0.5$)
 - Insert / unsuccessful search

2.5 probes

- Successful search
 1.5 probes
- Example (λ = 0.9)



Random probing

- Instead of using f(i)=i
- h_i(x) = (h(x) + probe_positions[i]) mod TableSize
 - probe_positions is an array that contains a list of random numbers
- The array should be generated beforehand and used in both search and insert
- Avoids primary clustering

Random probing

- Instead of using f(i)
- h_i(x) = (h(x) + probe_positions[i]) mod TableSize
 - probe_positions is an array that contains a list of random numbers
- The array should be generated beforehand and used in both search and insert
- Avoids primary clustering

Recall: probe sequence needs to be repeatable (so we can find/delete what we have inserted already)

Random probing

- Random probing does not suffer from clustering
- Expected number of probes for insertion or unsuccessful search:

$$\frac{1}{\lambda}\ln\frac{1}{1-\lambda}$$

- Example
 - $\lambda = 0.5 \rightarrow 1.4$ probes
 - $\lambda = 0.9 \rightarrow 2.6$ probes



Figure 5.12 Number of probes plotted against load factor for linear probing (dashed) and random strategy (*S* is successful search, *U* is unsuccessful search, and *I* is insertion)





- Avoids primary clustering
- f(i) is quadratic in i
 - e.g., f(i) = i^2
 - h_i(x) = ((h(x) + i^2) mode TableSize
- Probe sequence:
 - +0, +1, +4, +9, +16, ...
- Continue until an empty slot is found
- #failed probes is a measure of performance

Why?

Quadratic probing

- Avoids primary clustering
- f(i) is quadratic in i, e.g., f(i) = i
 - h_i(x) = (h(x) + i^2) mod TableSize
 - Probe sequence: +0, +1, +4, +9, +16, ...



37

Hashing

Quadratic probing

	Empty Table	After 89	After 18	After 49	After 58	After 69	
0				49	49	49	
1							
2					58	58	
3						69	
4							
5							
6							
7							
8			18	18	18		V F
9		89	89	89	89	89	
#unsuccessful p	robes:	0	0	1	2	2	
Figure 5.13 Hash table with quadratic probing, after each insertion							

Quadratic probing: analysis

- Theorem 5.1
 - New element can always be inserted into a table that is at least half empty and TableSize is prime
- Otherwise, may never find an empty slot, even is one exists
- Ensure table never gets half full
 - If close, then expand it
- May cause "secondary clustering":
 - elements that hash to the same position will probe the same alternative cells
 - Take long time to find empty slots

Quadratic probing: delete



- Deletion
 - Emptying slots can break probe sequence and could cause find stop prematurely
 - Lazy deletion
 - Differentiate between empty and deleted slot
 - When finding, skip and continue beyond deleted slots
 - If you hit a non-deleted empty slot, then stop find procedure returning "not found"
- When does it really (physically) delete the data?
 - Delete upon insert
 - Compaction at some time

Quadratic probing

- 1 template <typename HashedObj>
- 2 class HashTable
- 3 {
- 4 public:
- 5 explicit HashTable(int size = 101);
- б

7

8

- bool contains(const HashedObj & x) const;
- 9 void makeEmpty();
- 10 bool insert(const HashedObj & x);
- 11 bool insert(HashedObj && x);
- 12 bool remove(const HashedObj & x);
- 13

14

enum EntryType { ACTIVE, EMPTY, DELETED };

For lazy deletion

WSU

```
Hashing
          16
                 private:
         17
                   struct HashEntry
                                                                             Initialize by
                       HashedObj element;
          19
                                                                             EMPTY
          20
                       EntryType info;
          21
          22
                       HashEntry( const HashedObj & e = HashedObj{ }, EntryType i = EMPTY )
          23
                         : element{ e }, info{ i } { }
          24
                       HashEntry( HashedObj && e, EntryType i = EMPTY )
          25
                         : element{ std::move( e ) }, info{ i } { }
          26
                   };
          27
          28
                   vector<HashEntry> array;
          29
                   int currentSize;
          30
          31
                   bool isActive( int currentPos ) const;
          32
                   int findPos( const HashedObj & x ) const;
          33
                   void rehash( );
          34
                   size t myhash( const HashedObj & x ) const;
          35
              };
```

Figure 5.14 Class interface for hash tables using probing strategies, including the nested HashEntry class

Quadratic probing

```
explicit HashTable( int size = 101 ) : array( nextPrime( size ) )
1
2
          { makeEmpty( ); }
3
        void makeEmpty( )
4
5
                                       Clear hash table:
            currentSize = 0;
6
                                       set to EMPTY
            for( auto & entry : array
7
                entry.info = EMPTY;
8
9
```

Figure 5.15 Routines to initialize quadratic probing hash table

WSU



Figure 5.16 contains routine (and private helpers) for hashing with quadratic probing

WSU





Figure 5.17 Some insert and remove routines for hash tables with quadratic probing

Double hashing

- Use a second hash function
- Good choices for h_2(x)?
 - Should never evaluate to o
 - e.g., $h_2(x) = R (x \mod R)$
 - R is prime number less than TableSize
- Previous example with R=7
 - $h_2(x) = 7 (x \mod 7)$ $h_i(x) = (h(x) + f(i))$
 - f(i) = i * h_2(x)
 - $ho(49) = (h(49)+f(0)) \mod 10 = 9 (X)$
 - $h_1(49) = (h(49) + \frac{1}{7} 49 \mod 7) \mod 10 = 6$

WSU

Hashing

Double hashing: example



Double hashing: analysis

- TableSize must be prime
- Empirical tests show double hashing close to random probing
- Extra hash function takes extra time to compute



Probing techniques: review



Rehashing

- Increases the size of the hash table when load factor becomes "too high" (defined by a cutoff)
 - Anticipating that prob(collisions) would become higher
- Typically expand the table to twice its size (but still prime)
- Need to reinsert all existing elements into new hash table



Figure 5.21 Hash table after rehashing

5U

Rehashing: analysis

- Rehashing takes time to do N insertions
 - O(N)
- Therefore, should do it infrequently
- Specifically
 - Must have been N/2 insertions since last rehash
 - Amortizing the O(N) cost over the N/2 prior insertions yields only constant additional time per insertion

Rehashing: implementation

- When to rehash?
 - When load factor reaches some threshold (e.g., $\lambda \ge 0.5$), OR
 - When an insertion fails
- Applies across collision handling schemes

Rehashing for separate chaining



Figure 5.22 Rehashing for both separate chaining hash tables and pro

Rehashing for quadratic probing



Hash tables in C++ STL

- Some implementations of STL have hash tables (e.g., SGI STL)
 - hash_set
 - hash_map

STL unordered_map

- unordered_map: A better option for hash table
- Supported in C++ 11 and later
- Supports a standard interface for common hash table operations
- Operations on average are O(1)

Problems with large tables

- What if hash table is too large to store in main memory?
- Solution: Store hash table on disk
 - Minimize disk accesses
- But...
 - Collisions require disk accesses
 - Rehashing requires a lot of disk accesses



Hash table applications

- Symbol table in compilers
- Accessing tree or graph nodes by name
 - e.g., city names in Google maps
- Maintaining a transposition table in games
 - Remember previous game situations and the move taken (avoid recomputation)
- Dictionary lookups
 - Spelling checkers
 - Natural language understanding (word sense)
- Heavily used in text processing languages
 - e.g., Python dictionary

Hashing: summary

- Hash tables support fast insert and search
 - O(1) insert, search, delete avg performance
 - Deletion possible, but degrades performance
- Not suited if ordering of elements is important
- Many applications

Hashing: checklist to remember

- Table size prime
- Table size much larger than number of inputs (to maintain λ closer to o or < 0.5)
- Tradeoffs between chaining vs. probing
- Collision chances decrease in this order:

linear probing > quadratic probing > {random probing, double hashing}

- Rehashing required to resize hash table at a time when λ exceeds a threshold (usually 0.5)
- Good for searching. Not good if there is some order implied by data