CS171 Introduction to Computer Science II

Hash Tables

Review

- Sequential search using linked list
- Binary search using ordered arrays
- Binary search tree (BST)
- 2-3 tree (balanced search tree, implemented as red black tree)

ST implementations: summary

implementation	guarantee			average case			ordered	operations
	search	insert	delete	search hit	insert	delete	iteration?	on keys
sequential search (linked list)	Ν	N	N	N/2	Ν	N/2	no	equals()
binary search (ordered array)	lg N	N	N	lg N	N/2	N/2	yes	compareTo()
BST	Ν	Ν	Ν	1.38 lg N	1.38 lg N	?	yes	compareTo()
red-black tree	2 lg N	2 lg N	2 lg N	1.00 lg N	1.00 lg N	1.00 lg N	yes	compareTo()

Hashing: basic plan

Save items in a key-indexed table (index is a function of the key).



• Equality test: Method for checking whether two keys are equal.

Hashing: basic plan

Save items in a key-indexed table (index is a function of the key).

Hash function. Method for computing array index from key.

Issues.

- Computing the hash function.
- Equality test: Method for checking whether two keys are equal.
- Collision resolution: Algorithm and data structure to handle two keys that hash to the same array index.

Classic space-time tradeoff.

- No space limitation: trivial hash function with key as index.
- No time limitation: trivial collision resolution with sequential search.
- Space and time limitations: hashing (the real world).



Hash Tables



- A hash table for a given key type consists of – Hash function h
 - Array (called table) of size N

Computing the hash function

Idealistic goal. Scramble the keys uniformly to produce a table index.

- Efficiently computable.
- Each table index equally likely for each key.

thoroughly researched problem, still problematic in practical applications

Ex 1. Phone numbers.

- Bad: first three digits.
- Better: last three digits.

Ex 2. Social Security numbers. +

- Bad: first three digits.
- Better: last three digits.

Practical challenge. Need different approach for each key type.

573 = California, 574 = Alaska (assigned in chronological order within geographic region)



Example

- We design a hash table for storing entries as (SSN, Name), where SSN (social security number) is a nine-digit positive integer
- Our hash table uses an array of size N = 10,000 and the hash function
 h(x) = last four digits of x



Hash Functions

 A hash function is usually specified as the composition of two functions:

Hash code:

 h_1 : keys \rightarrow integers

Compression function:

 h_2 : integers $\rightarrow [0, N-1]$

- The hash code is applied first, and the compression function is applied next on the result, i.e., h(x) = h₂(h₁(x))

Java's hash code conventions

All Java classes inherit a method hashcode (), which returns a 32-bit int.

Requirement. If x.equals(y), then (x.hashCode() == y.hashCode()).

Highly desirable. If !x.equals(y), then (x.hashCode() != y.hashCode()).



Default implementation. Memory address of x. Trivial (but poor) implementation. Always return 17. Customized implementations. Integer, Double, String, File, URL, Date, ... User-defined types. Users are on their own.

Implementing hash code: strings



char	Unicode		
'a'	97		
'b'	98		
'c'	99		

- Horner's method to hash string of length L: L multiplies/adds.
- Equivalent to $h = 31^{L-1} \cdot s^0 + \ldots + 31^2 \cdot s^{L-3} + 31^1 \cdot s^{L-2} + 31^0 \cdot s^{L-1}$.



Hash code design

"Standard" recipe for user-defined types.

- Combine each significant field using the 31x + y rule.
- If field is a primitive type, use wrapper type hashcode().
- If field is an array, apply to each element. <---- or use Arrays.deepHashCode()
- If field is a reference type, use hashCode(). < ____ applies rule recursively

In practice. Recipe works reasonably well; used in Java libraries. In theory. Need a theorem for each type to ensure reliability.

Basic rule. Need to use the whole key to compute hash code; consult an expert for state-of-the-art hash codes.

Modular hashing

Hash code. An int between -2^{31} and $2^{31}-1$. Hash function. An int between 0 and M-1 (for use as array index). typically a prime or power of 2

```
private int hash(Key key)
{ return key.hashCode() % M; }
```

bug



Uniform hashing assumption

Uniform hashing assumption. Each key is equally likely to hash to an integer between 0 and M - 1.

Bins and balls. Throw balls uniformly at random into M bins.





Java's string data uniformly distribute the keys of Tale of Two Cities

Collision Handling

- Separate Chaining
- Linear Probing



Separate chaining ST

Use an array of *M* < *N* linked lists. [H. P. Luhn, IBM 1953]

- Hash: map key to integer i between 0 and M 1.
- Insert: put at front of *i*th chain (if not already there).
- Search: only need to search ith chain.



ST implementations: summary

implementation	guarantee				average case	ordered	operations	
	search	insert	delete	search hit	insert	delete	iteration?	on keys
sequential search (linked list)	Ν	N	N	N/2	N	N/2	no	equals()
binary search (ordered array)	lg N	N	Ν	lg N	N/2	N/2	yes	compareTo()
BST	Ν	N	N	1.38 lg N	1.38 lg N	?	yes	compareTo()
red-black tree	2 lg N	2 lg N	2 lg N	1.00 lg N	1.00 lg N	1.00 lg N	yes	compareTo()
separate chaining	lg N *	lg N *	lg N *	3-5 *	3-5 *	3-5 *	no	equals()

* under uniform hashing assumption

Collision resolution: open addressing

Open addressing. [Amdahl-Boehme-Rocherster-Samuel, IBM 1953] When a new key collides, find next empty slot, and put it there.



linear probing (M = 30001, N = 15000)

Linear probing

Use an array of size M > N.

- Hash: map key to integer i between 0 and M 1.
- Insert: put at table index i if free; if not try i + 1, i + 2, etc.
- Search: search table index *i*; if occupied but no match, try i + 1, i + 2, etc.



Linear probing: trace of standard indexing client



Linear probing ST implementation

```
public class LinearProbingHashST<Key, Value>
   private int M = 30001;
                                                                       array doubling
   private Value[] vals = (Value[]) new Object[M];
                                                                        and halving
   private Key[] keys = (Key[]) new Object[M];
                                                                        code omitted
   private int hash (Key key) { /* as before */ }
   public void put(Key key, Value val)
   {
      int i;
      for (i = hash(key); keys[i] != null; i = (i+1) % M)
         if (keys[i].equals(key))
             break;
      keys[i] = key;
      vals[i] = val;
   ł
   public Value get(Key key)
   {
      for (int i = hash(key); keys[i] != null; i = (i+1) % M)
         if (key.equals(keys[i]))
             return vals[i];
      return null;
   }
```

Separate chaining vs. linear probing

Separate chaining.

- Easier to implement delete.
- Performance degrades gracefully.
- Clustering less sensitive to poorly-designed hash function.

Linear probing.

- Less wasted space.
- Better cache performance.

ST implementations: summary

implementation	guarantee			average case			ordered	operations
	search	insert	delete	search hit	insert	delete	iteration?	on keys
sequential search (linked list)	N	N	N	N/2	N	N/2	no	equals()
binary search (ordered array)	lg N	N	N	lg N	N/2	N/2	yes	compareTo()
BST	Ν	N	N	1.38 lg N	1.38 lg N	?	yes	compareTo()
red-black tree	2 lg N	2 lg N	2 lg N	1.00 lg N	1.00 lg N	1.00 lg N	yes	compareTo()
separate chaining	lg N *	lg N *	<mark>l</mark> g N *	3-5 *	3-5 *	3-5 *	no	equals()
linear probing	lg N *	lg N *	lg N *	3-5 *	3-5 *	3-5 *	no	equals()

* under uniform hashing assumption

Hashing vs. balanced search trees

Hashing.

- Simpler to code.
- No effective alternative for unordered keys.
- Faster for simple keys (a few arithmetic ops versus log N compares).
- Better system support in Java for strings (e.g., cached hash code).

Balanced search trees.

- Stronger performance guarantee.
- Support for ordered ST operations.
- Easier to implement compareTo() correctly than equals() and hashcode().

Java system includes both.

- Red-black trees: java.util.TreeMap, java.util.TreeSet.
- Hashing: java.util.HashMap, java.util.IdentityHashMap.