

# Lecture 28: HashMap & Collections

CS 62  
Spring 2013  
Kim Bruce & Kevin Coogan

## Map<K,V>

- Collection of associations between a key and associated value, e.g. name & phone number
  - Though doesn't use Bailey's Association class
- As usual lots of implementations
- Also called dictionaries after example
  - Look up table!

## Hash Functions

- Want  $H: \text{ElType} \rightarrow \text{Subscripts}$ , where
  - $H(\text{elt})$  can be computed quickly
  - if  $e_1 \neq e_2$  then  $H(e_1) \neq H(e_2)$ 
    - $H$  is one-to-one
  - Usually difficult to achieve
  - Looked at examples Wednesday
- if redefine equals then must redefine hashCode  
so  $x.\text{equals}(y) \Rightarrow x.\text{hashCode}() == y.\text{hashCode}()$

## What if get Hash Clashes?

- Home address of key  $K$  is  $H(K)$ .
- Suppose have two keys  $K_1 \neq K_2$ ,
  - but  $H(K_1) = H(K_2)$ , i.e., have same home address
- What happens when insert both into hash table?
  - Note original key and value must both be stored!!
- Two ways out:
  1. Rehash as needed to find an empty slot (open addressing)
  2. External chaining

## Quadratic Probing

- Use  $(\text{home} + j^2) \% \text{TableSize}$  on  $j$ th rehash
  - Helps with secondary clustering, but not primary
- Can result in case where don't try all slots
  - E.g.,  $\text{TableSize} = 5$ , and start with  $H = 1$ . Rehashings give 2, 0, 0, 2, 1, 2, 0, 0, ...
  - The slots 3 and 4 will never be examined to see if they have room.

## Double Hashing

- Use second hash function on key to determine delta for next try.
  - E.g.,  $\text{delta}(\text{Key}) = (\text{Key} \% (\text{TableSize} - 2)) + 1$
  - Should help with primary and secondary clustering.
  - Ex: Suppose  $H(n) = n \% 5$ . Then  $H(1) = H(6) = H(11)$ .
    - However,  $\text{delta}(1) = 2$ ,  $\text{delta}(6) = 1$ , and  $\text{delta}(11) = 3$ .

## External Chaining

- Each slot in table holds unlimited # elts
  - Each slot is list -- implemented as desired
  - For good performance, list should be short
    - so no need for balanced binary search tree -- waste of time
- Advantages
  - Deleting simple
  - # elts in table can be > # slots
  - Avoids problems of secondary clustering
  - Primary clustering?

## Analysis

- Behavior of the hash clash strategies depends on the *load factor* of the table.
- Load factor  $\alpha = \# \text{ elts in table} / \text{size of table}$ 
  - ranges between 0 and 1 with open addressing
  - can be > 1 with external chaining.
- Higher the load factor, the more likely you are to have clashes.

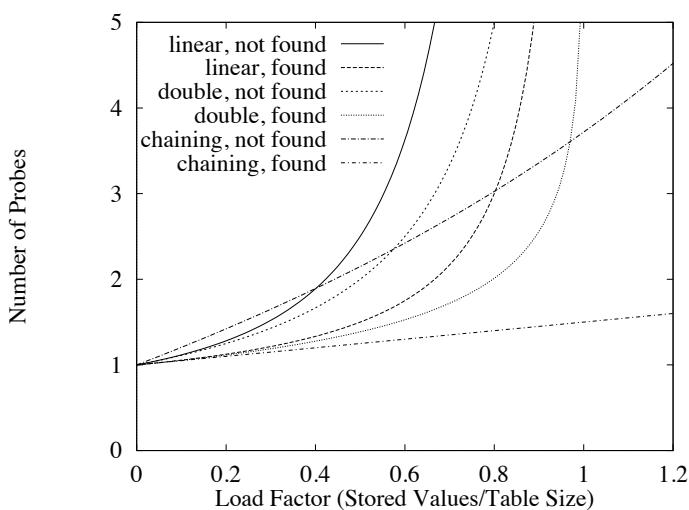
## Performance

Strategy	Unsuccessful	Successful
Linear rehashing	$1/2 (1 + 1/(1-\alpha)^2)$	$1/2 (1 + 1/(1-\alpha))$
Double hashing	$1/(1-\alpha)$	$- (1/\alpha) \log(1-\alpha)$
External hashing	$\alpha + e^{-\alpha}$	$1 + 1/2 \alpha$

*Entries represent number of compares needed to find elt or demonstrate not there.*

## Performance for $\alpha = .9$

Strategy	Unsuccessful	Successful
Linear rehashing	55	5.5
Double hashing	10	~4
External hashing	3	1.45



## Space requirements

- Open addressing:  $\text{TableSize} + n * \text{objectsize}$
- External chaining:  $\text{TableSize} + n * (\text{objectsize} + 1)$
- Rule of thumb:
  - Small elts, small load factor -- use open addressing
  - Large elts, large load factor -- use external chaining

## Using Hashcodes in Java

- HashMap and Hashtable both implement Map
  - Hashtable has all ops synchronized!
  - HashMap allows null keys and values - HT doesn't
  - HashSet is hashtable based implementation of sets.

## HashMap<K,V>

- HashMap constructor
  - Hashtable(int initialCapacity, float loadFactor)
  - Default load factor is .75 if not specified, default capacity 11.
  - If loadFactor exceeded then create larger table and rehash all old values – expensive!
- Implementation seems to use external chaining

## Capacity

- Don't want to set capacity too high as wastes space, though resizing expensive.
- Iterators through table require space proportional to capacity and current size.

## Collections Framework

- Java library implementations of most useful general data structures.
- Description at <http://docs.oracle.com/javase/6/docs/technotes/guides/collections/reference.html>
- Includes concurrent implementations of data structures

## C++

## C/C++

- Designed in 1972 as a “high-level assembly language” for writing UNIX.
  - Close to the machine, but portable
  - Designed for speed
- C++ mid-70's by Bjarne Stroustrup
  - O-O extensions based on Simulat 67
  - Backwards compatible with C
  - Speed still a goal.
    - Only pay for features you use
  - Java also uses C/C++ syntax

## Why C++

- C/C++ still popular out in world, especially if need to program close to the machine
- Fast
  - Though Java now nearly as fast
- Support for interacting directly w/hardware
- More control
- You need it for CS 105

## First Program

- See first.cpp
- Similarities w/Java:
  - Main method, like Java,
  - primitive types. very similar.
    - Java has 8 – short, int, long, double, float, boolean, char, and byte.
    - C++ has many of the same primitives – short, int, char.
    - In C++ they can be modified as “signed” or “unsigned.”
  - Library types. C++ has extensive standard library with many of the same types as Java.
  - Similar syntax. Curly braces. method construction

## First Program

- Similarities w/Java:
  - Similar syntax. Curly braces. method construction (parameters, return types, name).
  - Control constructs are basically the same – for loops, while loops, if..if-else, etc.
  - Use of generics. C++ has had generics for a long time, Java only got them recently with Java 5.
    - But C++ generics very different from Java
  - Calling methods of the object with the ‘.’ operator.

## Differences from Java

- C++ doesn’t require classes
- main method – no arguments, return type int.
  - See other versions of the main function in C++ later.
- #include statements.
  - Look similar to Java’s import statement, but not.
  - Equivalent of cutting and pasting file at thst location.
  - Say #include <file> to using a built in system file.
  - Say #include “file” we are referring to a user class file, and we give an appropriate path name.
  - # indicates preprocessor directive – before compilation

## first.cpp

- using std namespace
  - Equivalent of Java’s import statement.
  - Always have to specify.
  - In C++, the vector type is in the “std” namespace.
    - Get access by saying “use std namespace”
    - Now can write: vector<int> nums;
- Without “using” command, refer to std::vector or std::cout  
:: operator is the “scope resolution operator.”  
The sample file first\_no\_using.cpp gives a C++

## first.cpp

- The sample file first\_no\_using.cpp gives C++ program equivalent to first.cpp, but uses only the scope resolution operator.
- Function prototypes
  - C++ requires that declare function or vble before use it.
  - C++ allows function prototypes: Allow us to declare the name before we implement the code, so that can refer to the function before we implement it.

# first.cpp

- Operator overloading.
  - Notice the “nums[i]” notation. nums is of type vector, like the Java Vector or ArrayList.
  - The C++ documentation for vector shows the “operator[]” entry in the member function list.
  - overloaded the [] operator, so that it acts like “at” member function.
    - The “at” function returns value at a particular index in array.