## Hash Tables

Hash functions and hash tables Idea and Examples Hash function details Address Generation (Hash code + Compression code)

**Collision Resolution** 

Linear probing Quadratic probing Double hashing

## Idea

Hash tables are an example of a dictionary.

Data is stored and retrieved by use of a function of the key.

It is stored, but not sorted!

## Example

Student records are stored in an array using a 7 digit student i.d. the index.

If the i.d. were used unmodified, the array would have to have enough room for 10,000,000 student records.

Instead, student i.d.'s are hashed to produce an integer between, say 1 and 100,000 which indexes into an array.

## Problem

Since a possible 10,000,000 numbers are being compressed into just 100,000, how can we guarantee that no 2 i.d.'s end up stored in the same place?

















#### Hash Code Maps

Hash codes reinterpret the key as an integer. They should: 1. Provide good "spread"

2. Give the same result for the same key

#### Examples:

- Memory address:
- We reinterpret the memory address of the key object as an integer (default hash code of all Java objects)
- Good in general, except for numeric and string keys
- Integer cast:
  - We reinterpret the bits of the key as an integer
  - Suitable for keys of length less than or equal to the number of bits of the integer type (e.g., byte, short, int and float in Java)

#### Hash Code Maps

- Component sum:
  - We partition the bits of the key into components of fixed length (e.g., 16 or 32 bits) and we sum the components (ignoring overflows)
  - Suitable for numeric keys of fixed length greater than or equal to the number of bits of the integer type (e.g., long and double in Java)

#### Hash Code Maps (cont.)

- Polynomial accumulation:
  - We partition the bits of the key into a sequence of components of fixed length (e.g., 8, 16 or 32 bits)
    a<sub>0</sub> a<sub>1</sub> ... a<sub>n-1</sub>
  - We evaluate the polynomial  $p(z) = a_0 + a_1 z + a_2 z^2 + ... + a_{n-1} z^{n-1}$ at a fixed value z, ignoring overflows
  - Especially suitable for strings (e.g., the choice z = 33 gives at most 6 collisions on a set of 50,000 English words)

# Hash Code Maps Cont.

- Polynomial p(z) can be evaluated in O(n) time using Horner's rule:
  - The following polynomials are successively computed, each from the previous one in O(1) time  $p_0(z) = a_{r,1}$

$$p_i(z) = a_{n-i-1} + zp_{i-1}(z)$$
  
(i = 1, 2, ..., n -1)

• We have  $p(z) = p_{n-1}(z)$ 

#### **Compression Maps**

Compression Maps:

• Take the output of the hash code and compress it into the desired range.

• If the result of the hash code was the same, the result of the compression map should be the same.

•Compression maps should maximize "spread" so as to minimize collisions.

#### **Compression Maps Examples**

- Division:
  - $h_2(y) = y \mod N$
  - The size  $\boldsymbol{N}$  of the hash table is usually chosen to be a prime
  - The reason has to do with number theory and is beyond the scope of this course

#### **Compression Maps Examples**

- Multiply, Add and Divide (MAD):
  - $h_2(y) = (ay + b) \mod N$
  - *a* and *b* are nonnegative integers such that  $a \mod N \neq 0$
  - Otherwise, every integer would map to the same value  ${\pmb b}$

**Address Generation** 

Some examples ...





### Address Generation (b)

 $h_1(x)$ : integer cast

**b)**  $h_2(h_1(x))$ : sum of subset of bits of  $h_1(x)$ 

 $\rightarrow$  Simple to calculate

→ More random than a)



#### Address Generation (c)

**h**<sub>1</sub>(**x**): integer cast

- c)  $h_2(h_1(x))$ : subset (of r bits) of  $h_1(x)^2$
- $\rightarrow$  Multiplication is involved
- $\rightarrow$  More random than **a**) and **b**)

#### Address Generation (d)

**d)**  $h_2(h_1(x))$ : =  $h_1(x)$  MOD N

- $\rightarrow$  Division is involved!
- $\rightarrow$  Very random (if N is odd)







Collision Resolution (1) Linear Probing		
$ \underbrace{h(K_{i}), h(K_{i}) + 1, h(K_{i}) + 2, h(K_{i}) + 3}_{h_{0}(K_{i}) h_{1}(K_{i}) h_{2}(K_{i}) h_{3}(K_{i}) $		
Let $h_0(K_i) = h(K_i)$		
h <sub>j</sub> ( K <sub>i</sub> ) = [h ( K <sub>i</sub> ) + j ] mod N		





















#### Performances of Double Hashing

Experimental results for a hash table with load factor  $\ \alpha$ 

Search

α = n/N	<b>C(α )</b>
0.1 (10%)	1.05
0.5 (50%)	1.38
0.75 (75%)	1.83
0.9 (90%)	2.55

#### Performance of Hashing: Summary

- In the worst case, searches, insertions and removals on a hash table take O(n) time
   The worst case occurs when a
- The worst case occurs when all the keys inserted into the dictionary collide
   The load factor a = n/N
- The load factor a = n/N affects the performance of a hash table
- hash table • Assuming that the hash values are like random numbers, it can be shown that the expected number of probes for an insertion with open addressing is approximately  $1 / (1 - \alpha)$
- The expected running time of all the dictionary ADT operations in a hash table is  $\mathcal{O}(1)$
- In practice, hashing is very fast provided the load factor is not close to 100%
- Applications of hash tables: - small databases
  - compilers
  - browser caches