

## 7- Hashing

Bruno MARTIN,  
University of Nice - Sophia Antipolis  
<mailto: Bruno.Martin@unice.fr>  
<http://www.i3s.unice.fr/~bmartin/mathmods.html>

## Hashing

The steps in hashing:

- 1 compute a **hash function** which maps keys in table addresses

Since there are more records ( $N$ ) than indexes ( $M$ ) in the table, two or more keys may hash to the same table address : it's the **collision** problem

- 2 the **collision resolution** process

Good hash functions should uniformly distribute entries in the table

Since, if the function uniformly distributes the keys, the complexity of searching is approx. divided by the table's size

## Hashing

Hashing is a **completely different** method of searching

**The idea is to access directly the record in a table using its key - the same way an index accesses an entry in an array -**

We use a hash function that computes a table index from the key

**Basic operations:** insert, remove, search

## Transform Keys into Integers in $[0, M - 1]$

If the key is already a large integer

- choose  $M$  to be a prime and compute  $key \bmod M$

If the key is an uppercase character string

- encode each char in a 5-bit code (5 bits ( $2^5$ ) are required to encode 26 items): each letter is encoded by the binary value of its rank in the alphabet
- compute the modulo of the corresponding decimal value

### Example

$ABC \Rightarrow 00001\ 00010\ 00011 \Rightarrow$   
 $1 * (2^5)^2 + 2 * (2^5)^1 + 3 * (2^5)^0 = 1091 \Rightarrow 1091 \bmod M \Rightarrow$   
*index table*

## Why does $M$ have to be prime ?

An example of hash function is

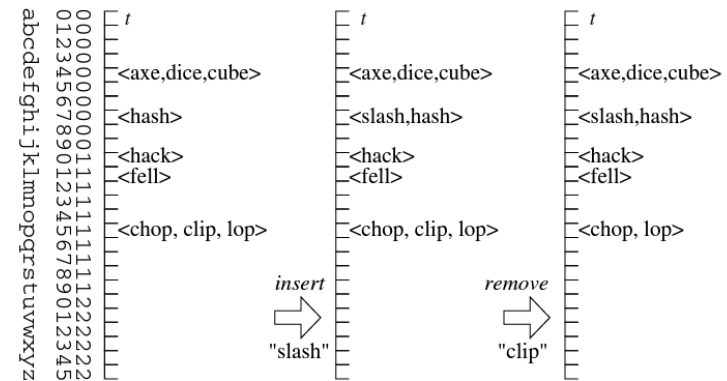
$$\text{hash}(key) = (key[0] \times (2^k)^0 + key[1] \times (2^k)^1 + \dots + key[n] \times (2^k)^n) \bmod M$$

Suppose you choose  $M = 2^k$  then

- $XXX \bmod M$  is unaffected by adding to  $XXX$  multiples of  $2^k$
- $\text{hash}(key) = key[0]$  :  $\text{hash}$  only depends on the 1<sup>st</sup> char of  $key$

The simplest way to ensure that the  $\text{hash}$  function takes all the characters of a key into account is to take  $M$  **prime**

## Example



**Fig. 4.1.** Hashing with chaining. We have a table  $t$  of sequences. The figure shows an example where a set of words (short synonyms of “hash”) is stored using a hash function that maps the last character to the integers 0..25. We see that this hash function is not very good

## How to Handle the Collision Process

We have an array of size  $M$  - called the hash table - and a  $\text{hash}$  function which gives for any key a possible entry in this array

**Problem:** decide what to do when 2 keys hash to the same address

A first simple method is to build for each table entry a **linked list** of records whose keys hash to the same entry

Colliding records are chained together we call it **separate chaining**

At the initialization, the hash table will be an array of  $M$  pointers to empty linked lists

## Searching a record in a Hash Table with linked lists

Main operation on a *HashTable*: **search** a record with its *key*:

- compute the  $\text{hash}$  value of the  $key$  :  $\text{hash}(key) = i$
- access to the linked list at position  $i$  :  $\text{HashTable}[i]$
- if there's more than your record in the list you have collisions
- searching becomes a search in a list: iterate on each record comparing the keys
- unsuccessful search: you iterate down the list without finding your record
- Operations of insertion and removal of records in a Hash Table become linked list operations

## Searching Performances

Good hash functions uniformly distribute  $N$  entries over the  $M$  positions of the table

Searching expected values in  $O(\alpha)$  ( $\alpha = \frac{N}{M}$  table's filling rate):

- Unsuccessful:  $\frac{1}{M} \sum_M (1 + \#L_i)$  since the element  $\notin L_i$

$$\overline{Q}^-(M, N) = \alpha + 1 \quad \text{since } \sum \#L_i = N$$

- Successful: searching for an element in the table equals the cost of inserting it when only the inserted elements before it were already in the table:

$$\overline{Q}^+(M, N) = \frac{1}{N} \sum_{i=0}^{N-1} \overline{Q}^-(M, i) = \frac{1}{N} \sum_{i=0}^{N-1} \left(1 + \frac{i}{M}\right) = 1 + \frac{\alpha}{2} = \frac{1}{2M}$$

The interest of hashing is that it is efficient and easy to program

## Expected cost – interpretation

- if  $N = O(M)$ , then  $\alpha = N/M = O(M)/M = O(1)$
- searching takes constant time on the average
- insertion is  $O(1)$  in the worst case
- deletion takes  $O(1)$  worst-case time for doubly linked lists
- hence, all dictionary operations take  $O(1)$  time on average with hash tables with chaining

## Alternative proof for successful search

- $x_i$  is the  $i^{\text{th}}$  element inserted into the table and  $k_i = \text{key}[x_i]$
- $X_{ij} = \mathbf{1}\{h(k_i) = h(k_j)\}$  for all  $i, j$  (indicator R.V.)
- simple uniform hashing:  $\Pr\{h(k_i) = h(k_j)\} = 1/M \Rightarrow E[X_{ij}] = 1/M$
- expected number of elements examined in a successful search:

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

$\sum_{j=i+1}^N X_{ij} = \#$  of elements inserted after  $x_i$  into the same slot as  $x_i$ .

$$\begin{aligned} (1) &= \frac{1}{N} \sum_{i=1}^N \left( 1 + \sum_{j=i+1}^N E[X_{ij}] \right) = \frac{1}{N} \sum_{i=1}^N \left( 1 + \sum_{j=i+1}^N \frac{1}{M} \right) = \\ &= 1 + \frac{1}{NM} \sum_{i=1}^N (N - i) = 1 + \frac{1}{NM} \left( \sum_{i=1}^N N - \sum_{i=1}^N i \right) = \\ &= 1 + \frac{1}{NM} \left( N^2 - \frac{N(N+1)}{2} \right) = 1 + \frac{N-1}{2M} \end{aligned}$$

## Another structure for Hash Table: Linear Probing

When the **number** of elements  $N$  can be **estimated** in advance  
You can **avoid** using any **linked list**  
You **store**  $N$  records in a table of size  $M > N$   
**Empty places** in the table **help you** for collision resolution  
It is called the **linear probing**



## Eliminating the Clustering Problem

Instead of examining each successive entry, we use a **second hash function** to compute a fixed increment to use for the sequence (instead of using 1 in linear probing)

Depending on the choice of the second hash function, the program may not work : obviously 0 leads to an infinite loop

## Hashing in Ruby

```
zip=Hash.new
zip={"06000" => "Nice", "06100" => "Nice", "06110" => "Le Cagnet",
"06130" => "Grasse", "06140" => "Coursegoules", "06140" => "Tourrettes
sur Loup", "06140" => "Vence", "06190" => "Rocquebrune Cap Martin",
"06200" => "Nice", "06230" => "Saint Jean Cap Ferrat", "06230" =>
"Villefranche sur Mer"}

zip["06300"]="Nice" # adds a new entry
zip.keys=>["06140", "06130", "06230", "06110", "06000", "06100",
"06200", "06300", "06190"]

zip.values=>["Vence", "Grasse", "Villefranche sur Mer", "Le Cagnet",
"Nice", "Nice", "Nice", "Nice", "Rocquebrune Cap Martin"]

zip.select { |key,val| val="Nice"}=>[["06000", "Nice"], ["06100",
"Nice"], ["06200", "Nice"], ["06300", "Nice"]]

zip.index "Nice" => "06000"

zip.each {|k,v| puts "#{k}/#{v}"=>
06140/Vence
06130/Grasse
06230/Villefranche sur Mer
06110/Le Cagnet
06000/Nice
06100/Nice
06200/Nice
06300/Nice
06190/Rocquebrune Cap Martin
```

## Conclusion on Hashing

Hashing is a classical problem in CS: various algorithms have been studied and are widely used

There are many empirical and analytic results that make utility of Hashing evident for a broad variety of applications

**Hashing is preferred to binary tree searches** for many applications because it is **simple to implement** and can provide **very fast constant searching times** when space is available for a large enough table