

## **Hashing**

---

---

Must devise an appropriate hash function.  
Because the hash function maps a large space  
to a small space, we will have “collisions” .

---

We can make each location a “bucket” that  
can store lots of records.

- But buckets must have fixed size,  
thus they can still overflow.

We will need a scheme to handle this.

Must decide on the # and size of buckets.

When file gets very full, collisions can be too  
numerous. May be worthwhile re-organizing  
the file layout to have more buckets

## **Reading:**

- Chapter 9 on searching, including hashing

## Hashing Performance

If everything is well designed, retrieval can be very fast — just a few file accesses.

One operation is really slow:

A hash function is a mathematical function that maps from keys  $\Rightarrow$  locations.

There are some standard types of hash function, including

- **mid-square:** square the number and then take some digits from the middle.
- **folding:** Divide the number in half and combine the two halves, *e.g.*, add them together.
- **modular division:** Mod by some number, preferably a prime.

See the text for more about hash functions. Note that there is a lot of interesting theory about hash functions and their properties. (csc 378 covers this.)

## Examples of hash functions

Say our key is a string. Before we hash it, we need to turn it into an integer.

One solution: Concatenate together the alphabetic position of the 1st and the second character. E.g. “Toyota”  $\Rightarrow$  2015.

Now we need a hash function to hash up the integer. Examples of the three general types:

- Mid-square, taking the middle 2 digits. E.g.  $2015^2 = 406\boxed{02}25$ .

- Folding: adding the two halves.  
E.g.  $2015 \Rightarrow 20 + 15 = 35$ .

- Mod by 97.  
E.g.  $2015 \bmod 97 = 75$ .

Two kinds of collisions:

- collisions that occur because we begin with the same 2 integers. Are collisions for every hash function we might choose.
- collisions that occur even though we begin with 2 different integers. Not necessarily collisions with a different hash function.

Couldn't we use the 2015 as the hashed value?

This would be a bad idea. With 4 digits, there are 10,000 possible values (0..9999), yet only a few will be used.

- Some are unlikely to crop up  
**E.g.** "aa"  $\Rightarrow$  0101.
- Some cannot crop up  
**E.g.** ??  $\Rightarrow$  2701; 27 is out of range.

Yet we need our hash table to be continuous, and therefore to have all 10,000 slots. So our hash table will be largely empty.

With each of the 3 hash functions we looked at, the range of  $h(key)$  is 0..99 (or less with mod 97). So our hash table only needs 100 slots.

Of course we might overflow it, but we have to deal with this anyway.

## Avoiding collisions?

Upon doing a new insertion, how likely is a collision?

It's certainly more likely when many items have already been inserted.

**Loading factor:** =  $\frac{\#records\ currently\ in\ the\ file}{\#records\ that\ the\ file\ can\ hold}$

In our cars example, the loading factor is only  $\frac{13}{100} = 0.13$ , yet we already have collisions!

We could reduce collisions by making the capacity of the file bigger (and hence the loading factor smaller). But ...

## Exactly how likely are collisions?

For a given file capacity, how likely are collisions as file gets more loaded?

Example: A file of 365 buckets. Let  $Q(n)$  be the probability that NO collisions occur during  $n$  insertions.

$$Q(1) =$$

$$Q(2) =$$

$$Q(3) =$$

In general,  
 $Q(n) =$

$$Q(1) =$$

Solution to the recurrence relation:

$$Q(n) = \frac{365!}{365^n(365-n)!}$$

The probability that collisions DO occur is  
 $\frac{1 - Q(n)}{1 - Q(n)}$ .

$n$	$\frac{1 - Q(n)}{1 - Q(n)}$
10	0.1169
20	0.4114
23	0.5073
30	0.7063
40	0.8912
50	0.9704
60	0.9941

If the loading factor is only  $\frac{23}{365}$ , 50% chance.

If the loading factor is only  $\frac{47}{365}$ , > 95% chance!

So yes, collisions are a problem!

## Buckets

The hash function  $h(k)$  tells us where to store (or retrieve) a record with key  $k$ .

This could be a slot in an array in memory, or a slot in a file. (For this course, a file.) Either way, we often use the term “hash table”.

The slots are called “buckets”, because they have capacity for  $> 1$  record. We choose the capacity based on the number of records that can be read or written in one file access.

Analogy: your address book.  
**key:**

**hash function:**

**bucket:**

148

What do we do when a collision occurs?

Easy case: the bucket has room for the record.

Hard case: the bucket doesn’t have room for the record. We call this “overflow”.  
We need to figure out two things:

- A place to put the record that won’t fit it its home bucket.
- A way to find that record later!

How do you handle overflow in your address book?

149

## Handling overflow

Two kinds of approach: either compute or store where to try next. (Gee, where have we heard that before?)

### Open addressing

Compute another bucket to try, based on some rule.

Compute where to look, based on the key.  
General method for insertion:  
(search is analogous; why?)

### Open Addressing

Let  $A_i$  be where to look on the  $i$ th try.

- Use the hash function ( $h$ ) to find the first bucket where the record might go,  $A_0$ .

$$A_0 = h(key)$$

### Closed addressing (or chaining)

Store the location of another bucket to try, using some sort of pointer.

- If that bucket is full, use a new function ( $f$ ) to find the next bucket to try. Repeat as necessary.
- Stop when we hit a bucket with room, or the sequence of  $A_i$ 's starts to repeat (*i.e.*, there is no room anywhere).

Many ways to design the function  $f$  ...

## I. Linear Probing

## Linear probing example

$n$  (# buckets) = 13  
bucket size = 1  
step size = 2

$$h(\text{key}) = (\text{sum 1st } 3) \bmod 13$$

$$\text{So } A_i = (A_{i-1} + 2) \bmod 13$$

Step through a sequence of buckets always using the same step size.

Use mod to wrap around when we hit the end of the hash table.

$$A_i = (A_{i-1} + stepSize) \bmod n$$

(We can also express  $A_i$  in terms of  $A_0$ .)

### Example:

$$\begin{aligned} A_i &= (A_{i-1} + 2) \bmod n \\ &= (A_0 + 2^i) \bmod n \end{aligned}$$

The sequence of buckets considered is called the “probe sequence”.

key	$h(\text{key})$	probe sequence
Chevrolet	$16 \bmod 13 = 3$	
Chrysler	$29 \bmod 13 = 3$	
Jaguar	$18 \bmod 13 = 5$	
Nissan	$42 \bmod 13 = 3$	
Karmann Ghia	$30 \bmod 13 = 4$	

Bucket #	Bucket Contents
0	
1	
2	
3	
4	
5	
6	
7	
8	
9	
10	
11	
12	

## II. Non-linear Probing

Problem: **Primary clustering**.

If several records hash to the same spot, or even *any* spot along the probe sequence, they will all follow that same probe sequence.

Example:  $n = 100$ ; step size = 2.  
hash to: 5    probe seq:  
              9    probe seq:

Solution: Make the step size depend on the step number,  $i$ .  
This is called **non-linear probing**.

E.g.,  $A_i = (A_{i-1} + 2i^2) \bmod n$

Example:  $n = 100$ ; step size =  $2i^2$ .  
hash to: 5    probe seq:  
              9    probe seq:

154

## III. Double Hashing

Problem: **Secondary clustering**.

All records that hash to the same spot still have the same probe sequence.

Example:  $n = 100$ ; step size =  $2i^2$ .  
key 1; hash to: 5    probe seq:  
key 2; hash to: 5    probe seq:

Solution: Make step size depend on the *key* (but differently than in original hash function).  
This is called **double hashing**.

E.g.,  $A_i = (A_{i-1} + h_2(key)) \bmod n$

155

## Closed Addressing

Instead of computing where to look next, store it, using a “pointer”.

Each full bucket has a pointer to an overflow area,

- in the same file (for example at the end)
- Or in another file

There are many ways to organize this, since pointers are so flexible.

Cost vs open addressing:

Savings:

Do deletions introduce problems?