

A **database** stores **records** with various **attributes**.

The database can be represented as a **table**, where each **row** is a record, and each **column** is an attribute.

Number	Name	Dept	Alias
20090612	오재훈	산디과	alpha0401
20100202	강상익	무학	scala
20100311	손호진	무학	python_is_great

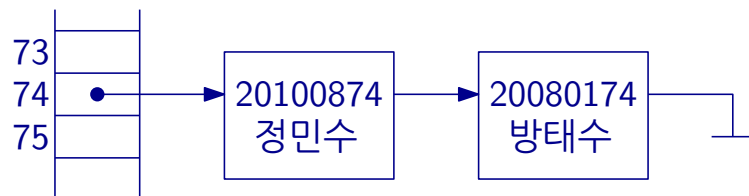
row

column

Databases often designate one attribute as the **key**. The key has to be unique—every key appears on only one row. A table with keys is a **keyed table**.

We want to find records (rows) by key, so the keyed table is a map: key  $\rightarrow$  record.

**Chaining:** Each slot is actually a linked list of (key, value) pairs stored in this slot. (We need the key!)



To search for a key 20080174, we access the table at index 74, and then search through the linked list.

Let's make a keyed table of all the students in the class, with the student number as the key.

```
case class Student(name: String, id: Int,
                  dept: String, alias: String)
```

Using an array with 100 slots, we can use the last two digits of the student number as the index.

But the last two digits are not unique — we have **collisions**:

Number	Name	Dept	Alias
20100874	정민수	무학	ubuntu
20080174	방태수	산디과	apple

We assume the hash function is good: It should distribute the items on the slots **uniformly**.

Analysis of hash tables assumes that the hash function is **random**: Each slot is equally likely to be chosen. The choices for two different items are **independent**.

Consider insertion/deletion/searching an item  $x$ . The running time is proportional to the length of the chain for  $x$ .

This is equal to the number of items  $y$  for which  $h(y) = h(x)$ . For given  $y$ , this happens with probability  $1/N$ . The expected value for all  $y$  is  $n/N$ .

Here  $n$  is the number of items, and  $N$  is the table size.

**Load factor:** The load factor  $\lambda$  of a hash table is  $n/N$ . Running time is  $O(\lambda)$ .

We could make the data structure much more compact if we could avoid the linked lists and store all data in the table.

**Open addressing:** allow to store items at a slot different from its hash code.

**Closed addressing:** items must be stored at the slot given by its hash code: chaining.

Easiest form of open addressing: **Linear probing.**

Start at the slot given by the hash code.

If it is already in use, try the next, and continue until a free slot is found.

0	
1	
2	
3	
4	
5	
6	
7	
8	18
9	89

insert: 89 18

5

0	49
1	
2	
3	
4	
5	
6	
7	
8	18
9	89

insert: 89 18 49

6-3

0	49
1	58
2	
3	
4	
5	
6	
7	
8	18
9	89

insert: 89 18 49 58



Assuming that the hash function behaves randomly, the expected number of probes for an insertion (or unsuccessful search) is (for  $N \rightarrow \infty$ ):

$$\frac{1}{2} \left( 1 + \frac{1}{(1 - \lambda)^2} \right)$$

Linear probing works very well when the hash function is good and the load factor  $\lambda$  is small, say  $\lambda \leq 0.5$ .

Linear probing is more sensitive to bad hash functions than chaining.

Load factor includes items that have been deleted! When there are too many deleted items, we need to **rehash** the table.

9

Hash codes and compression functions are a bit of a black art. It is easy to mess up.

An obvious compression function is  $h_2(x) = x \bmod N$ .

It only works well if  $N$  is a prime number.

A better compression function is

$$h(x) = ((ax + b) \bmod p) \bmod N,$$

where  $a$ ,  $b$ , and  $p$  are positive integers,  $p$  is a large prime, and  $p \gg N$ .  $N$  does not need to be prime.

11

We typically use two functions:

### Hash code

$h_1 : \text{keys} \rightarrow \text{integers}$

### Compression function

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

Index in hash table is computed as  $h_2(h_1(\text{key}))$ .

Ideally, the hash function should map keys uniformly at random to an index into the hash table.

**Resizing hash tables:** We change the compression function only, and then need to **rehash** all elements.

10

A good **hash code** for strings:

```
def hashCode(key: String): Int = {
  var hash = 0
  for (ch <- key)
    hash = (127 * hash + ch) % 16908799
  hash
}
```

Mix up the bits

Each character has different effect.

Bad hash codes:

- Sum up the codes of the letters (too small, and anagrams collide).
- Take the first three letters (“pre” is common, “xzq” never occurs).

Why is the function above good? Because it works in practice...

12

Scala `HashSet` and `HashMap` compute a hash code by calling the element's `##` method.

Every Scala object implements `##`.

`HashSet` and `HashMap` only work correctly if the following “contract” is observed:

If `obj1 == obj2` then `obj1.## == obj2.##`.

This is true for all standard types, but needs to be done by the programmer for new types!

The default implementation of `##` simply returns the memory address of the object on the heap.

**Mutable** keys are dangerous! If you change a key in the hash table, you cannot find it anymore.

13

Hashing with guaranteed **constant** search time!

We need **two** independent hash functions  $h(x)$  and  $g(x)$ .

To insert item  $x$ , check if slot  $h(x)$  or  $g(x)$  is empty. If so, insert the item there.

Otherwise, let  $y$  be the item at position  $g(x)$ . Insert  $x$  at slot  $g(x)$ , and move  $y$  to its other possible slot.

To find an item, we only need to check two slots!

If the load factor is small enough ( $\lambda \approx 1/3$ ), then insertions take expected constant time.

15

**Quadratic probing:** Try slots  $i + j^2$ , for  $j = 0, 1, 2, \dots$

**Secondary hash function:** Try slots  $i + jd(k)$ , where  $d(k)$  is a secondary hash function.

The details are tricky, because we need to make sure that the probing will find an empty slot.

Before you implement this, read a good book!

14

Hash tables do not support order on the items.

Hashing is fast if the hash function can be computed quickly.

Typical applications of hashing:

- symbol tables (in a compiler etc.),
- small databases,
- remembering positions (in a game tree),
- caching data (in a browser etc.),
- dictionaries.

16