



BITS Pilani
Hyderabad Campus

Database Management Systems

Dr.R.Gururaj
CS&IS Dept.

Lecture Session-12

Hashing Techniques



Content

1. Introduction to hashing
2. Internal hashing
3. Collision
4. External hashing
5. Static hashing
6. Dynamic hashing

Hashing



Hashing technique is an alternative to indexing, for fast retrieval of data records based on search key.

The search field is called as *hash field* of the file.

In most cases the hash field is also a key field of the file, in which case it is called as *hash key*.

The basic idea of hashing is that a hash function h , when supplied a hash field value K of a record produces the address B of the disk block that contains the record with specified key value.

$$\begin{array}{ccc} \underline{h} & : & \underline{K} \rightarrow \underline{B} \\ \uparrow & & \uparrow \\ \text{Hash} & & \text{Disk block} \\ \text{function} & & \text{address} \end{array}$$

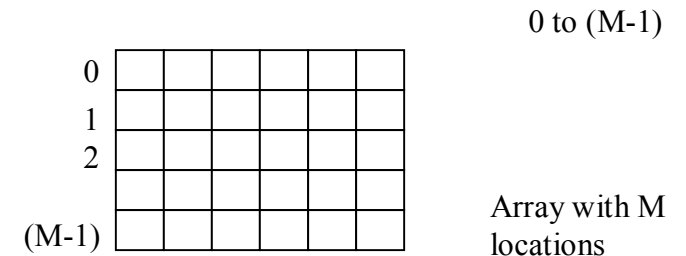
Once the disk block is known, the actual search for the record within the block is carried out in main memory buffer.

For most records we require only one block access.

Internal Hashing

Used for internal files.

A hash table is implemented through use of an array of records.



The most common hash function used is $h(k) = K \bmod M$

This gives the index of the location in the array.

For example- if $M = 10$ and the key value is 24

$$K \bmod M \Rightarrow 24 \bmod 10 = 4$$

Hence the record with key value 24 will be stored in 5th location of the array.

If two or more records are hashed to same location it is called as *collision*.

Then we need to find some other location for the new record. This process is known as *collision resolution*.



Methods for collision resolution

Open addressing: When collision occurs try with alternate cells until an empty cell is formed.

Chaining: for this various overflow locations are kept by extending the array by number of overflow positions. A pointer field is added to each record location. Collision is resolved by allocating an unused overflow position.

Multiple hashing: We apply a second hash function if the first hashing results in a collision.

The goal of a good hashing function is to distribute the records uniformly over the address space so as to minimize collisions while not leaving many unused locations.

External Hashing

Hashing used for disk files is called as *external hashing*. The disk block contains records. A single disk block or cluster of contiguous blocks is known as a *bucket*.

The hashing function maps a key value into a relative bucket number. A table maintained in the file header converts the bucket number into the corresponding disk block address, as shown in the figure below.

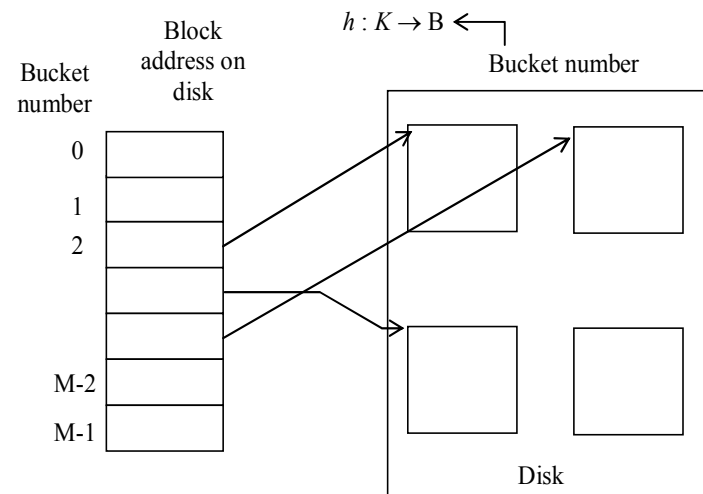
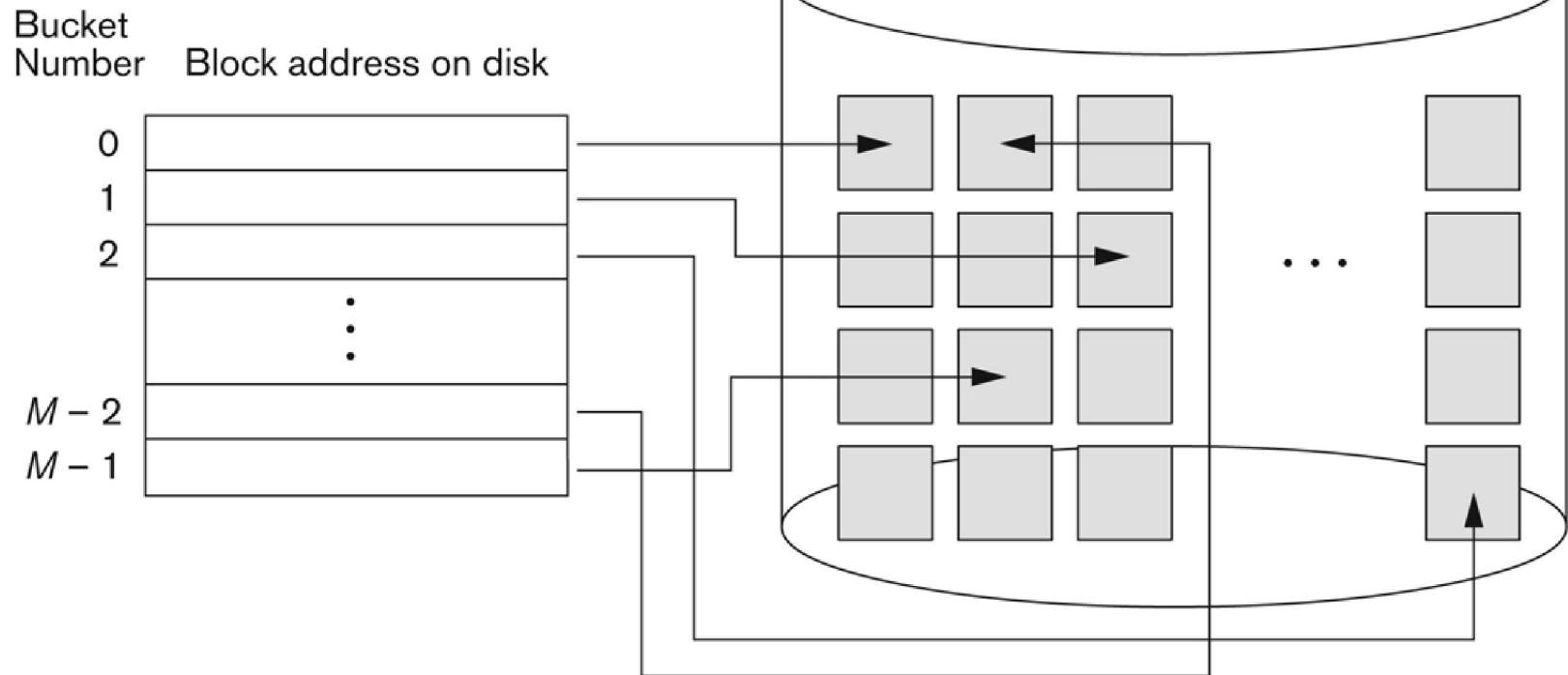


Figure 13.9

Matching bucket numbers to disk block addresses.





The above scheme is called as *static hashing* because the number of buckets allocated is fixed. This is a big constraint for files that are dynamic.

When a bucket is filled to capacity and if the new record is hashed on to the same bucket, then chaining is adopted, where a pointer is maintained in each bucket to a linked list of overflow records for the bucket.

The pointers are record pointers which include both block address and a relative record position within that block.

Handling overflows in Static External Hashing

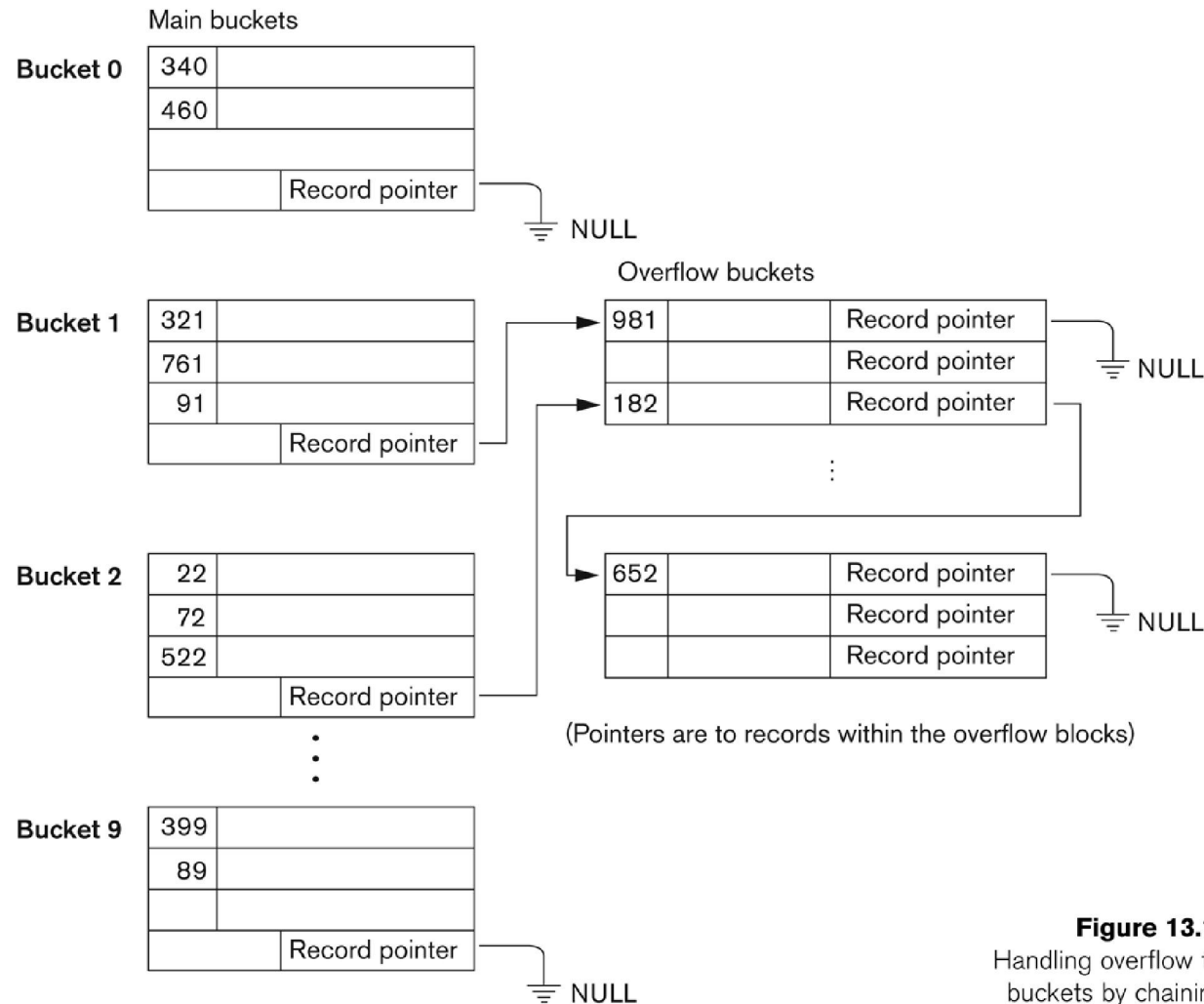


Figure 13.10
Handling overflow for
buckets by chaining.

Dynamic Hashing

This scheme allows us to expand or shrink the hash address space dynamically.

Each result of applying the hash function is a nonnegative integer and hence can be represented with a binary pattern. This we call it as *hash value* of the record.

Records are distributed among the buckets based on the values of the leading bits in their hash value.

Extendible Hashing

The first technique is called as *extendible hashing*. This scheme stores a directory structure in addition to the file. This access structure is based on the result of the hash function to the search field. The major advantage of extendible hashing is that performance does not degrade because of chaining, as the file grows as we have seen in static hashing. In extendible hashing no additional space is wasted towards the allocations for future growth, but additional buckets can be allocated dynamically as needed. The only overhead in this scheme is that a directory structure needs to be searched before the buckets are accessed.

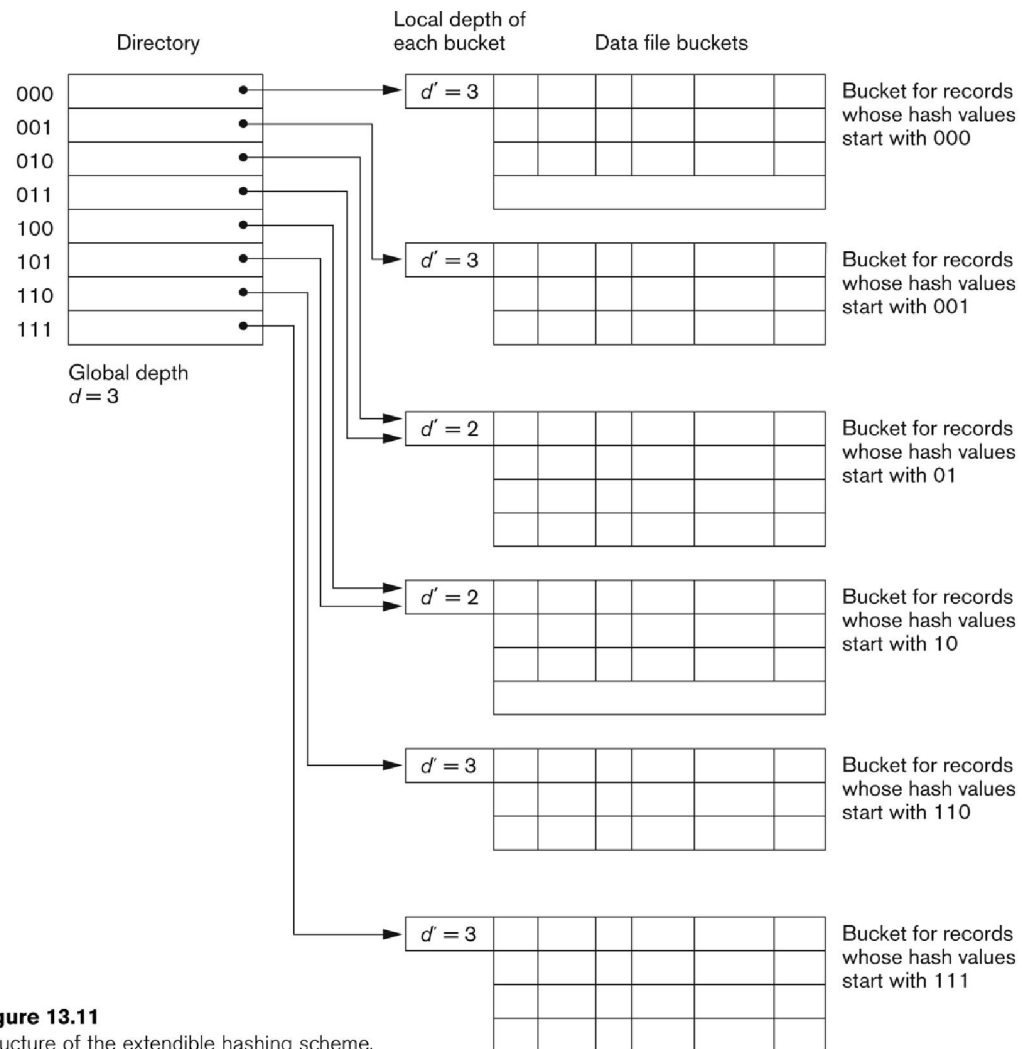


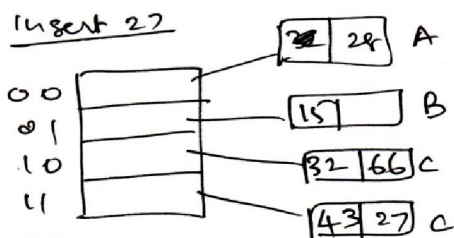
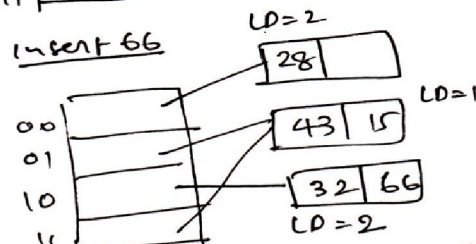
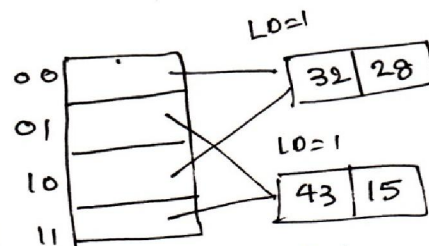
Figure 13.11
Structure of the extendible hashing scheme.

Q4: Key value

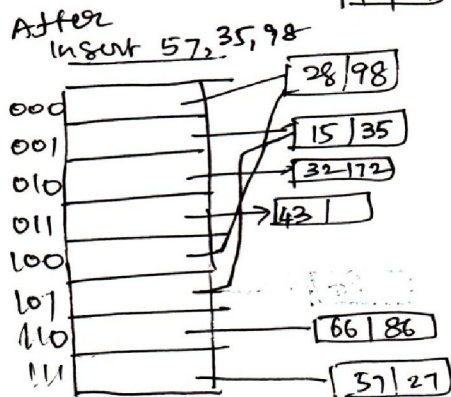
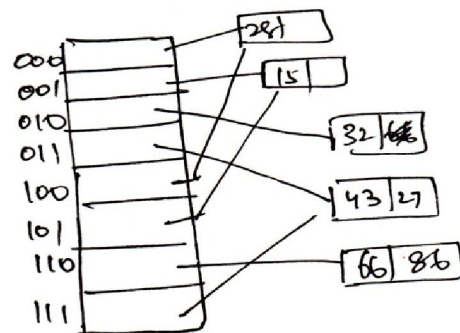
(a) Key values: 32, 28, 43, 15, 66, 27, 86, 57, 35, 98, 72.

| | Binary |
|------------------|--------|
| $32 \div 10 = 2$ | 0010 |
| $28 \div 10 = 8$ | 1000 |
| $43 \div 10 = 3$ | 0011 |
| $15 \div 10 = 5$ | 0101 |
| $66 \div 10 = 6$ | 0110 |
| $27 \div 10 = 7$ | 0111 |
| $86 \div 10 = 6$ | 0110 |
| $57 \div 10 = 7$ | 0111 |
| $35 \div 10 = 5$ | 0101 |
| $98 \div 10 = 8$ | 1000 |
| $72 \div 10 = 2$ | 0010 |

Global depth = 2



When insert 86 increase Global depth



Linear Hashing



In the second scheme called *linear hashing*, no directory structure is used. Here instead of one hash function, multiple hash functions are used. When collision occurs with one hash function, the bucket that overflows is split in to two and the records in the original bucket are distributed among two buckets using the next hash function $h_{(i+1)}(k)$. Hence we have multiple hash functions.

Q5 @ Linear Hashing.

$$h_0 = K \text{ Mod } 2$$

$$h_1 = K \text{ Mod } 4$$

$$h_2 = K \text{ Mod } 8$$

| keys | 14 | 11 | 28 | 15 | 33 | 46 | 8 | 21 | 4 |
|--------------------------|----|----|----|----|----|----|---|----|---|
| $h_0 (K \text{ Mod } 2)$ | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 |
| $h_1 (K \text{ Mod } 4)$ | 2 | 3 | 0 | 3 | 1 | 2 | 0 | 1 | 0 |

$n=0$
Initially we have only 2 buckets

Insert 14 0 [14 | 28]

Insert 11

Insert 28 1 [11 | 15]

Insert 15

Insert 33 - overflow
Split bucket 0

$n=1$
Insert 46 (h_1) 0 [28 | 8]

Insert 8 (h_1) 1 [11 | 15] → [33 | 21]

2 [14 | 46]

Insert 21

Insert 4

overflow
split '1'

0 [28 | 8] → [4]

1 [33 | 21]

2 [14 | 46]

3 [11 | 15]

Now all buckets are
split in 1st round

$n=0$

we use h_1 for
next insertions.



Summary



- What is hashing
- Internal hashing
- External hashing
- What is static external hashing
- What is dynamic hashing
- How Extendible and Linear hashing techniques work