

ADB Lecture 1

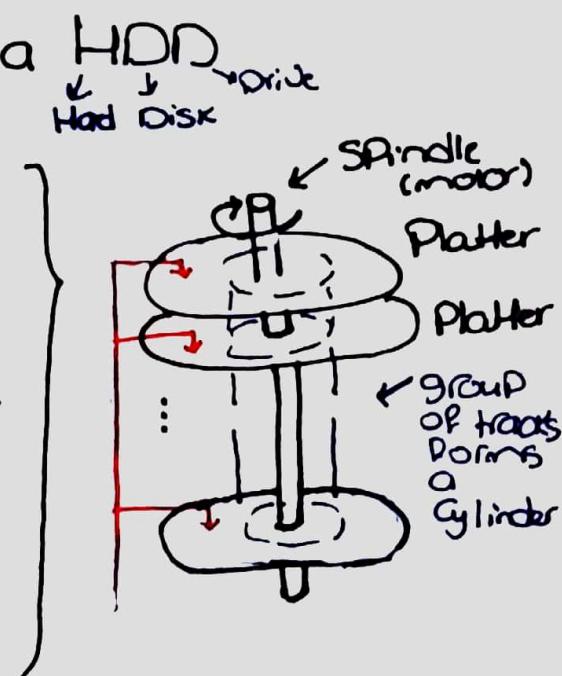
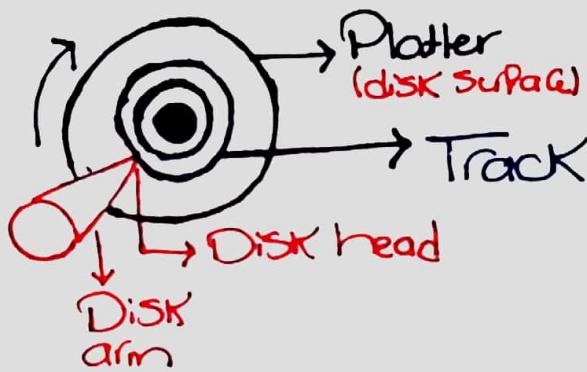
⇒ In this Lecture, we'll understand how the DBMS (database management system) stores records

→ Data corresponding to records is huge & must persist across program executions

Hence, the DBMS uses a HDD

→ Many layers of

- magnetic disk
- Stores OS and IS by using 2 diff. magnetic Polarity



2. Every Platter has a group of Concentric Circles (tracks)

→ its capacity is usually in the 4-50 KB range. Notice that outer tracks → more data
Zoning

1. The group of Platters (magnetic disk's) connected to the rotating Spindle is called a Disk Pack

3. Each track is further divided into Sectors.

1 Smallest addressable unit on HDD.

* The OS can choose to deal with the disk at a different basic unit (logical block)

- In this case, it must be a multiple of the Sector Size.
 - * Often Set as large as the Page Size (basic unit for main memory)

- Consequently, in a DBMS a Page
 - has the same size as a block
 - is the smallest unit for read/write from disk
 - corresponds to a set of consecutive records in the DB

- To retrieve a random Page

1. Disk Arm Positions the disk head on the track that has the page (via a seeking motion).
2. It then has to wait for the disk to rotate so that the head is on the right block
3. Then wait until reading (or writing) the block is done (by rotating over it)

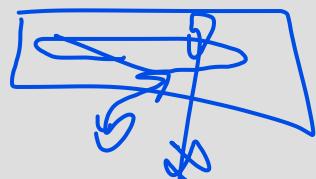
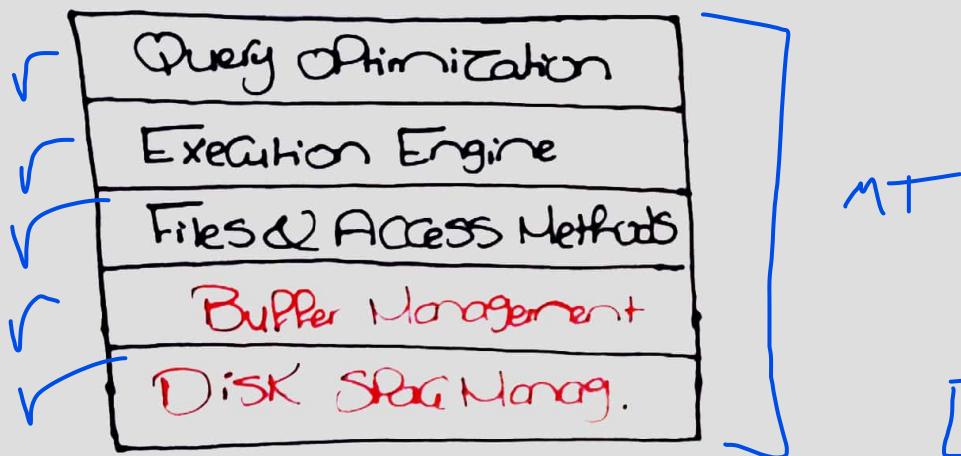
Hence, total access time is given by

$$T_{\text{access}} = T_{\text{seek}} + T_{\text{rot delay}} + T_{\text{transfer}}$$

• Slow due to moving parts

- We assume that a random Page can be retrieved at a fixed Cost (using average access time (avg. seek time and $T/2$ for rotation delay))
 - Still true that reading consecutive Pages is much cheaper than reading them in random
 - Seek and rotation time only once added.

- DBMS follows a layered architecture



1. Disk Space Management

- Lowest layer of DBMS Software
- Called by higher layers to read/write or allocate/deallocate a page
 - They only call the routines (no need to know how it's done or how free space is managed)

2. Buffer Management

- Brings pages to main memory (by calling lower layer) if they aren't already there.
- More importantly, decides which pages in the main memory should be replaced to make room for new pages
- Uses a replacement policy (such as FIFO)

- This is important because some queries (e.g., statistical) require going over many pages (or all). But at any time the main memory will never have enough room.

* Note that "file and access methods" (the layer above) can ask the buffer manager to access a page based on a record id (rid) it contains.

- Even higher* layers just ask for records and files of records
- A file is a logical collection of pages (each with its records)

The DBMS is expected to support

record id
Should be sufficient

→ insertion, deletion, modification of a record

→ reading a particular record

→ scanning all records (or a specific set)

- The way such files are organized on external storage is called File Organization

As we'll see next lecture, sometimes storing an extra index data structure based on some index search key can speed up specific queries

→ Index search key doesn't have to be a primary, foreign, candidate or superkey (can be any other column.)

more next lecture

• Many File Organization methods (alternatives) exist:

• Heap Files (random order)

• Sorted Files

• Hash or tree-based Indexes

1. Heap Files

(simplest)

→ Has nothing to do with Priority queues / heap tree (heap \leftrightarrow random)*

→ The records in the Pages of a heap file are not ordered in anyway.

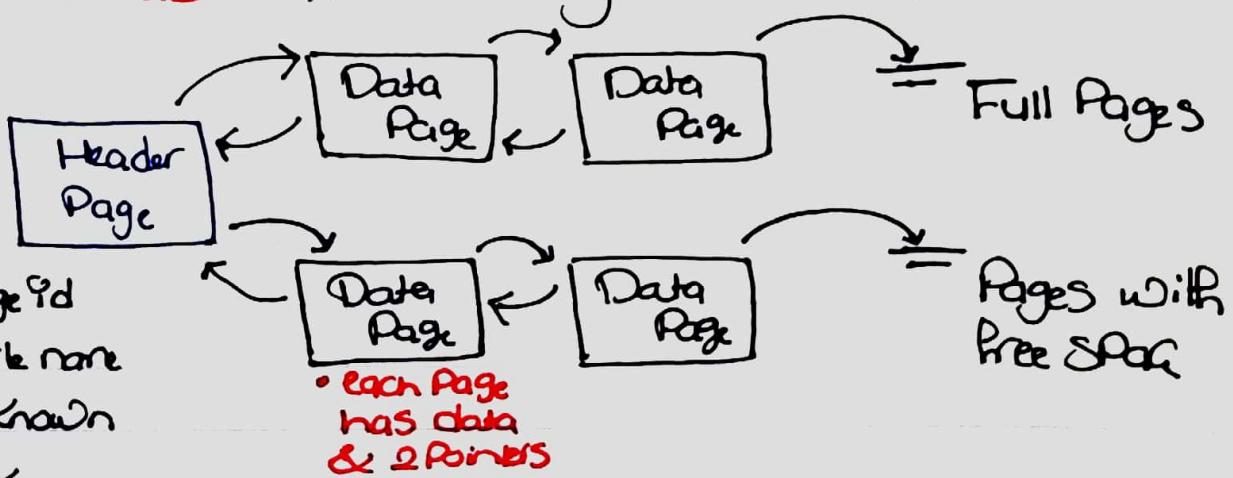
• To support record level operations, we need

→ To keep track of pages in a file

→ To keep track of records in a page

→ To keep track of free space on pages

• One way to achieve this is to implement it as two doubly-linked lists



• header Page Id
& heap file name
Stored at known
loc in disk

Notice that
Search time
Per Page is
Const.

* To Search

→ Just go over
all Pages

→ avg. time taken is $b/2$
($b = \# \text{Pages in File}$)

→ Full

→ Pages with
free space

* To Insert

- Start with the first free page
- Keep going until you find a page with enough free space to fit your record.

This can be a problem especially when records have varying fields and thus are of varying lengths

- may need to scan several pages in that case

→ Allocate a new one on the disk (and add its pointer)
if there's none with enough space (and deallocate if 100% free).

* Generally good when we tend to retrieve all the data (e.g. for statistics)

- Note: there are multiple methods to arrange the records themselves on a page

- Page contains a collection of slots (one for each record)
- Each record is hence identified by <Page-Id, Slot-num> = r^{id}

- Depend on whether data is made of fixed or varying length records
- In each case there are different options for

How Fields are organized in a record (Record Format)

- e.g., if they are varying we can decide to Put Columns next to each other & use **delimiters**
↳ (to Separate)



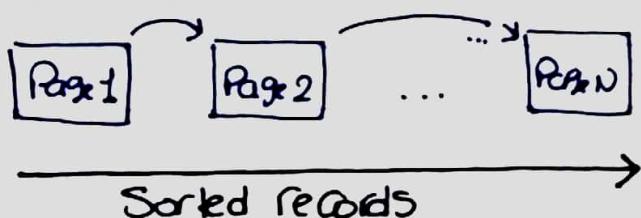
How to organize the records within a Page (Page Format)

- whether to **Span** or **unSpan**
 - + Don't allow records to SPAN more than one Page
- can lead to unused space



2. Ordered Files (Sequential Files)

- Defined by an ordering Field (records are always sorted w.r.t of it)
 - Does not involve data structure overhead
 - gives you **quick search** ($\log_2 b$ for b Pages in File) (binary S.)



- less data structure overhead

- Clearly, **insertion can be a headache** (Search for insertion position, Change Pointers after inserting Page)
- **Solutions**

- Maintain an overflow file and keep insertions there and merge Periodically (Searching done on both)

* Use ordered files when data is not expected to change frequently (class roster).

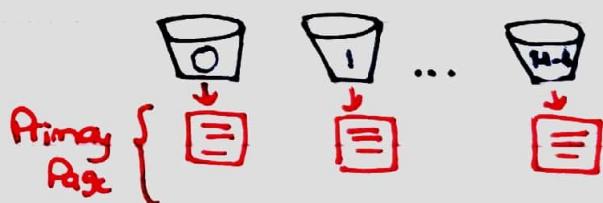
- Notice that sequential scan still takes $\frac{N}{M}$ even if file is ordered.

3. Hashed Files

(called external hashing)
• i.e., Per disk files

- Choose a field in the data to be a hash key

→ Divide file blocks into M ones of equal size



- For record in data:

$$h = K \bmod M \quad \# K \text{ is chosen field}$$

Insert_at_bucket(h) $\#$ goes to Primary Page

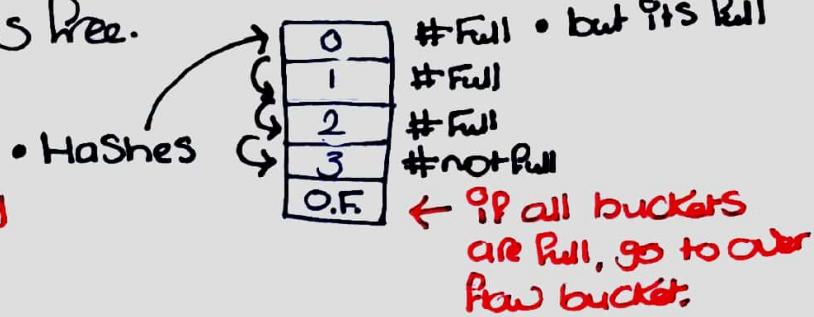
- $O(1)$
- They clearly speed up search (when it's in the form
.... WHERE Chosen.Field = '...')
 - Same applies to insertion (hash chosen field)

* Nothing special
for other fields

- When a Collision occurs (a record hashes to a bucket that's already full)

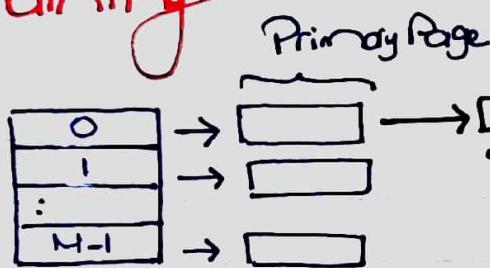
• Open Addressing

→ Keep Considering the next bucket until you find one that's free.



* In this case, it's called Linear Probing

• Chaining



• Overflow Page
(make it and Point to it on Primary
is full)

• More Complex
Versions
might exist

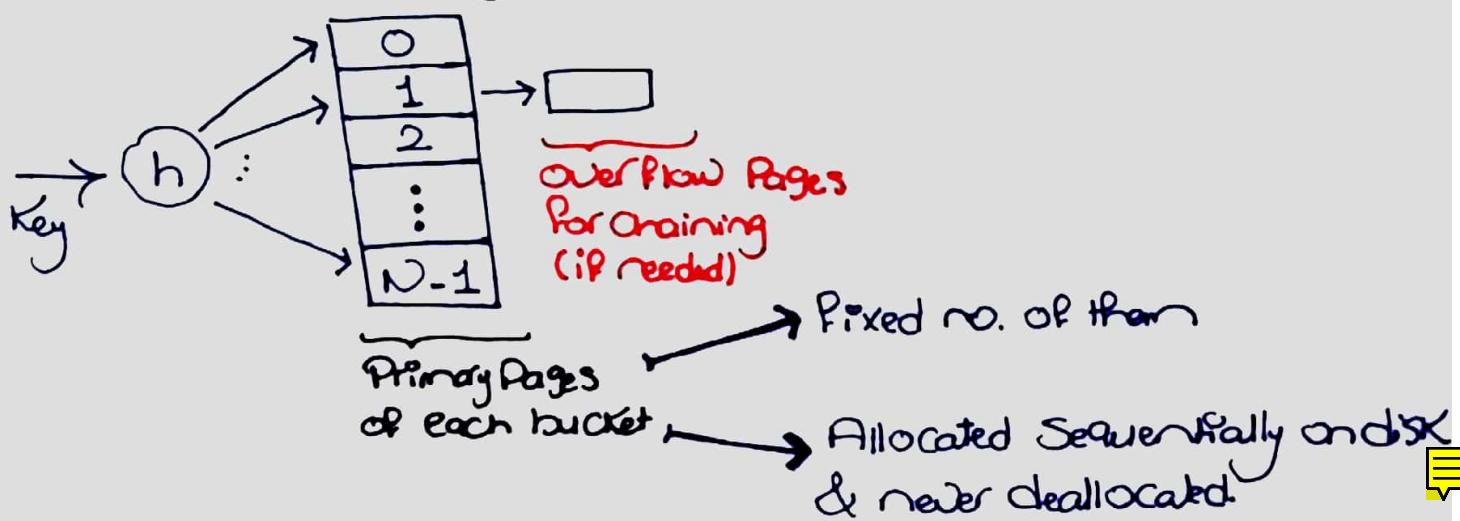
• Multiple Hashing

→ instead of choosing the bucket based on $h_1(K)$ everytime, if it results in a collision use another hash function $h_2(K)$ as an offset value



→ If there's a collision at the bucket ③ then we can use a 3rd hash function as offset or open use the 2nd again (Open addressing in general)

Static Hashing



* Need the Hash Function to distribute Keys over the Range $0, 1, \dots, N-1$

$$h(K) = (aK + b) \% N$$

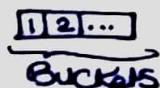
↓
Chosen Constants
(to tune h)

works good for that.

- Advantages

→ # Buckets are Fixed and Known → Stored Sequentially
Thus (assuming no overflows):

1. Search takes 1 I/O (read from disk)



2. Insertion takes 2 I/O

(read from disk → insert → write to disk)

①

②

!! Only Count disk IOs

(insertion is on main memory)

• Some applies for deletion

• Disadvantages

- If File Shrink's, can't make use of the allocated disk space that's free
- A lot of Collisions can lead to long overflow chains
 - Search / insertion no longer takes const. time

• Solutions

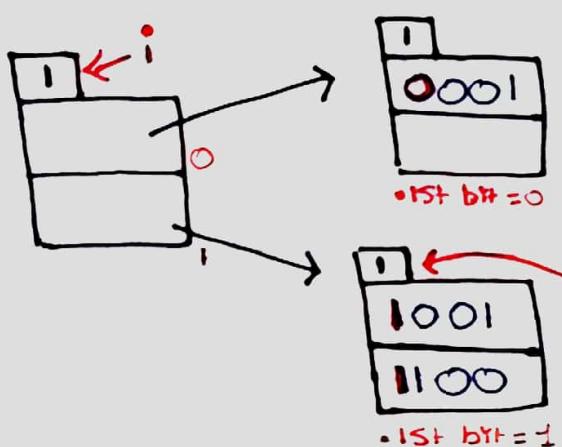
- Assuming file doesn't grow much, just try to assure that Pages are 80% full initially (Prevent overflow)
- When overflow pages build up, Periodically rehash (i.e., increase* N (no. of buckets) then reapply the hash function)
- Use dynamic hashing

Extendible Hashing

(type of dynamic hashing)

- Buckets (Primary Pages) are not fixed.
- One idea to implement is to reorganize the file by doubling the no. of buckets when an overflow will happen.
- Need to read previous configuration, double buckets and write again (rehash). . expensive operation

Rather, Use a directory (list) of pointers to buckets



* hash function h_k typically takes form:
 $h_i = h_k \% 2^i$

→ Thus, to decide bucket look at last(i) bits (global depth)

① unlike book, Slides takes First (i) bits

"local depth" = no. of bits to look at for this bucket.

- Here h_k is 4 bits (clearly)
- Each bucket holds 2 keys

Algorithm:

- If an overflow is about to happen upon insertion
 - 1. Make another bucket besides the one that was about to overflow and increase their local depth by 1.
 - Redistribute entries on both according to new local depth.
 - 2. If ① requires that the directory size is doubled then do that and fix pointers

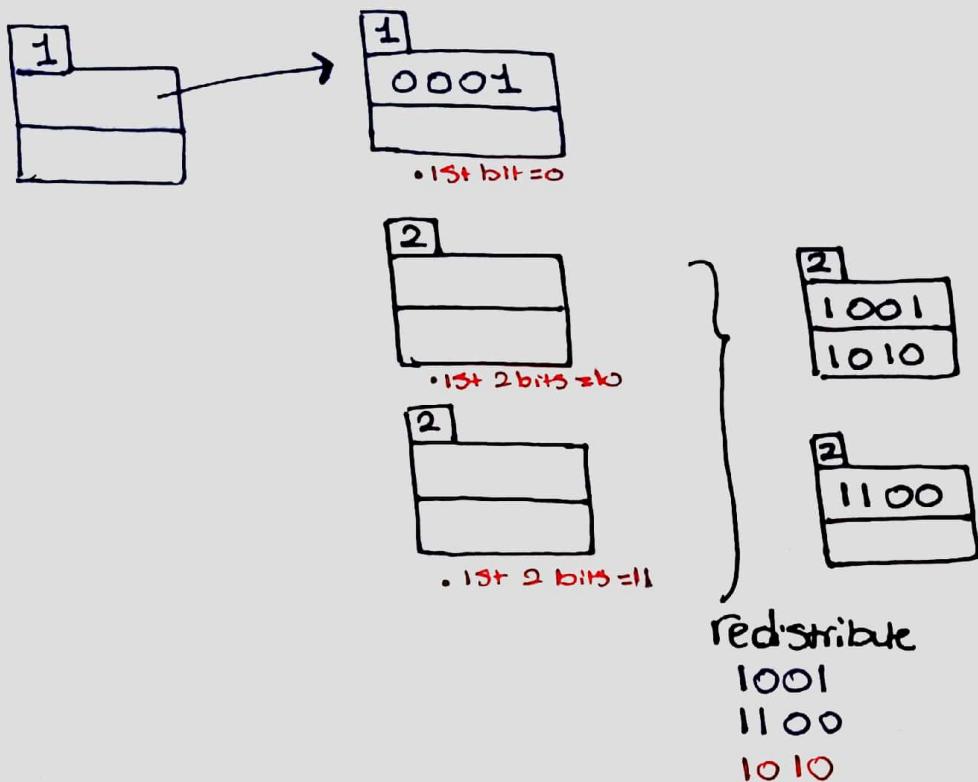
* Let's try to insert '1010' for the previous setup.

- 1st bit is 1 so it should go to 2nd bucket
- Will overflow!

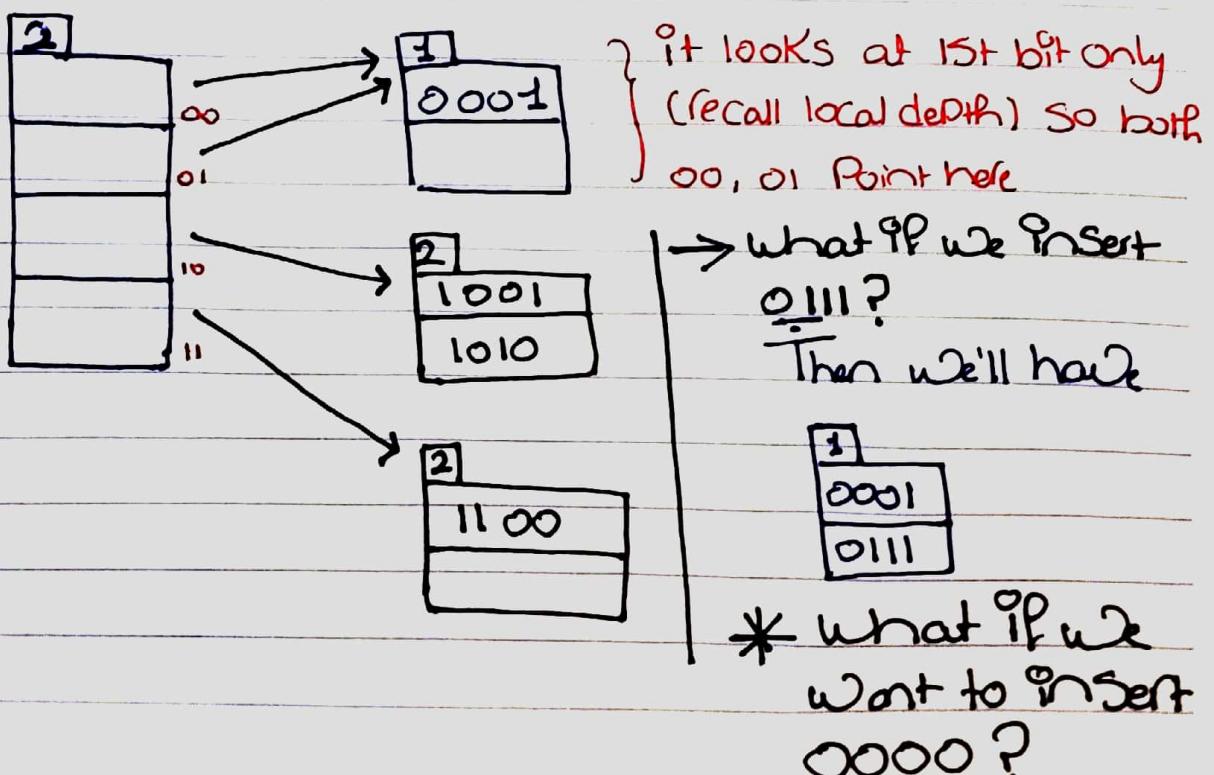
Thus,

Called 'Split Page'

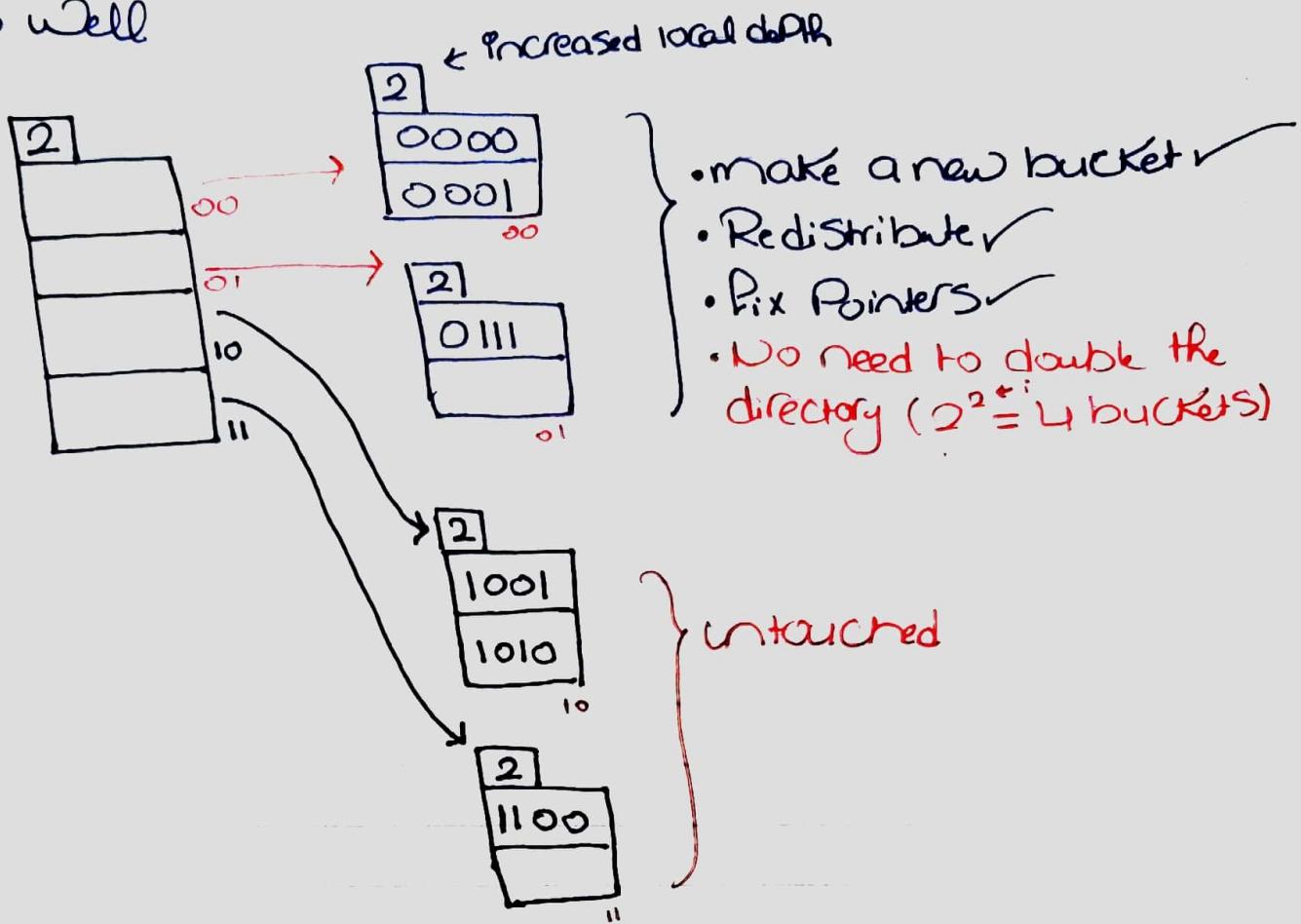
→ Let's make a new bucket and let both have increased local depth.



→ Now clearly, a directory of $i=1$ can't correspond to 3 buckets. So we need to double i .
• This will be true whenever local depth of some bucket(s) > global depth (i)



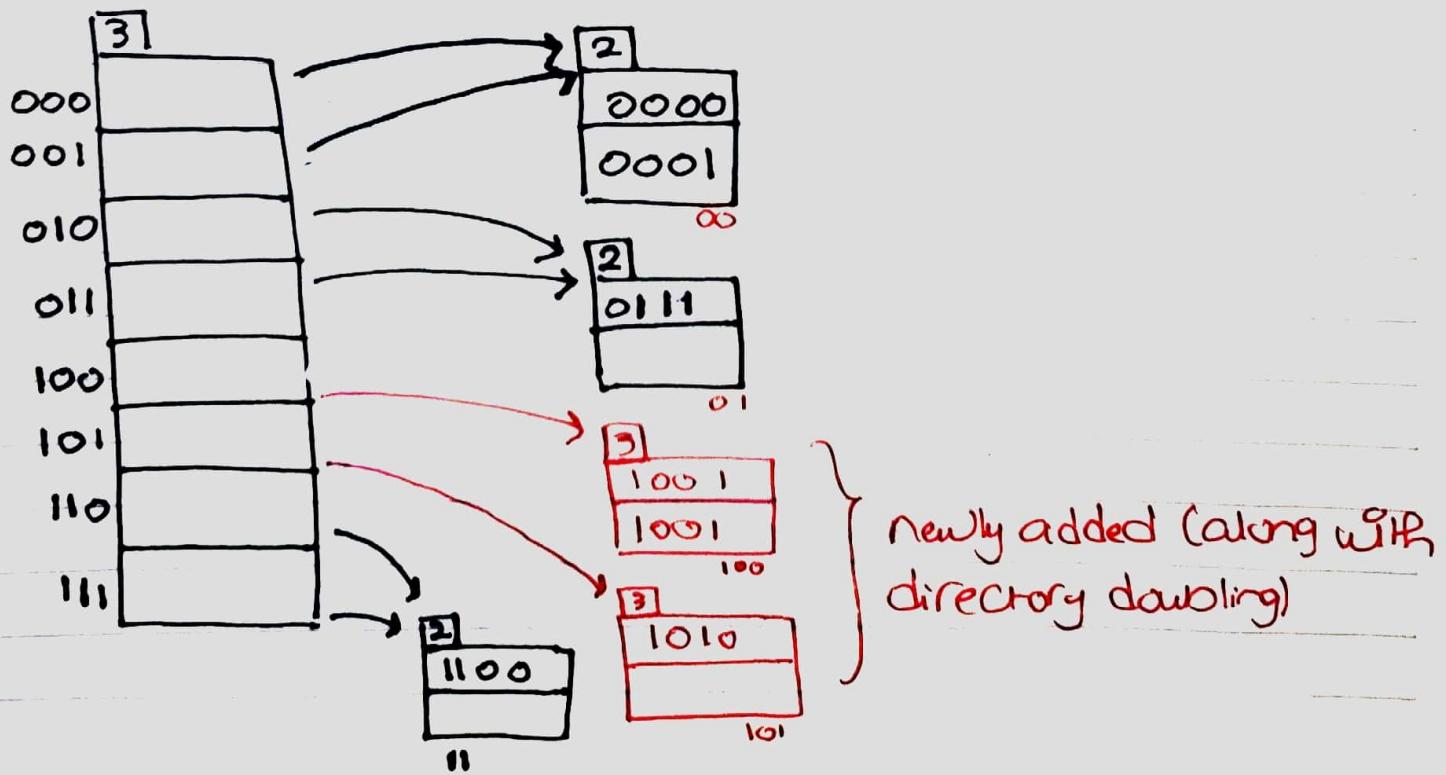
→ 1st bucket will overflow, let's make another as well



. Now suppose we want to insert 1001 (again)

→ Clearly need another bucket (3rd will overflow)

→ Clearly will need to double directory



Comments on extendible Hashing

- One disk access for sure when doing equality search (`WHERE = ' '`)
 - Because no overflows is guaranteed
 - assuming directory fits in ^{main} mem.
- Sudden / quick growth can occur due to many collisions (**Skewed hash values**)
 - entries with same hash aggravate the problem

* Note

- ✓ → Upon deletion, an empty bucket can be merged with its split image
 - The directory can be halved again if all directory elements point to both bucket and its split image (intuitive).

- ✓ → Dynamic Hashing also possible via a tree data structure (list of pointers at leaves and path to determine the binary corresponding to leaf element)

e.g.,

