

Unit 3: Storage Strategies

Why indexing is required

There may be two types of data (Ordered, Unordered) (Sorted, Unsorted).

To reduce the I/O cost

Indexing is required

Consider a HDD in which

block size = 1000 bytes

Each record size = 250 bytes

If total number of records are 10,000 and the data entered in HDD without any order. What is average time complexity to search a record from hard disk.

$$\text{no of records in a block} = \frac{1000}{250} = 4$$

$$\text{no of blocks required} = \frac{10,000}{4} = 2500$$

$$\begin{aligned} \text{No of records in each block} \\ = \text{One block size / Each record size} \\ = 4 \end{aligned}$$

$$\text{Total no of blocks required} =$$

$$\text{Total no of records / records in one block} = 2500$$

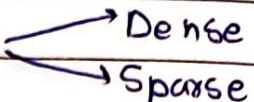
$$\text{Best Case} = 1$$

$$\text{Avg} = \frac{1 + 2500}{2} = 1250$$

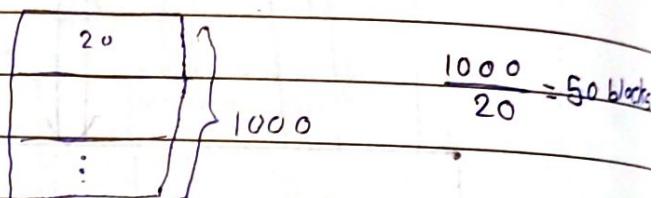
$$\text{Worst} = 2500$$

For Sorted

What is avg time complexity to search a record from index table. Index table entry = ^{Done} 20 bytes (10 for key + 10 for pointer).

Index
 

Dense is used when data is unordered it consist records for all records in HDD.



Size of index is same as a block
Sparse only stores initial value of each block. Thus reducing size.

Each block in index table contains total = total size of index / index entries. = 50 records

In Sparse data should be ordered

$$5 \cdot 2500 = 50$$

$$\log_2 50 = 6$$

No of Searches in index = 6
No of Searches to go to record = $6+1=7$

For dense storage

$$S = 10.000 = 200$$

50

~~$\log_2 200 = 8$~~

$$\text{Total Searches} = 8 + 1 = 9.$$

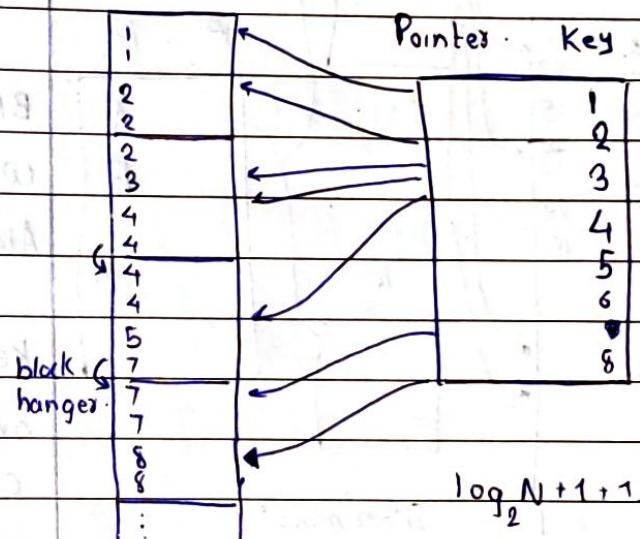
We use indexing to reduce input output cost.

* Types of Indexes :

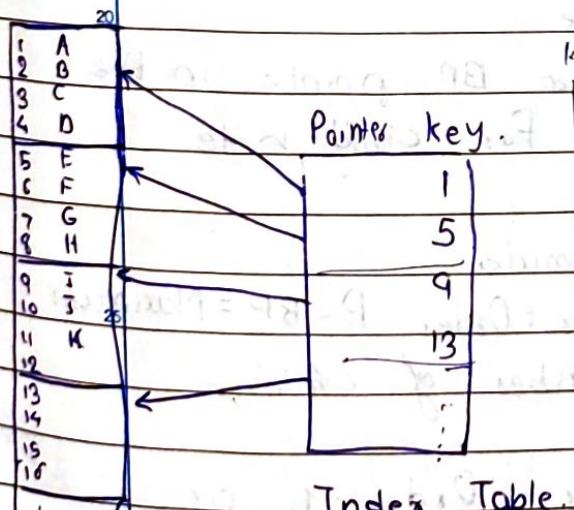
Sorted	Primary	Clustered
Unsorted	Secondary	Secondary
	Key	Non key

- Primary index has Sparse key values.
- No of entries in index table is equals to number of blocks in hard disk.
- Time complexity is $\log_2 N + 1$, where $N = \text{No of blocks in index tables}$.

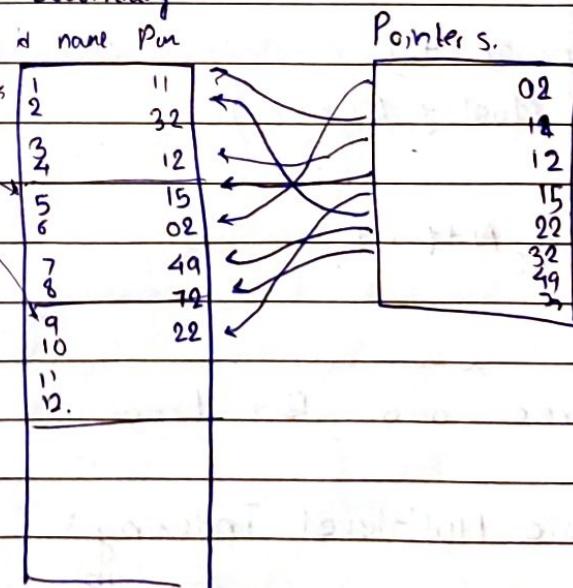
2] Clustered :



1] Primary :



3] Secondary :



Sparse

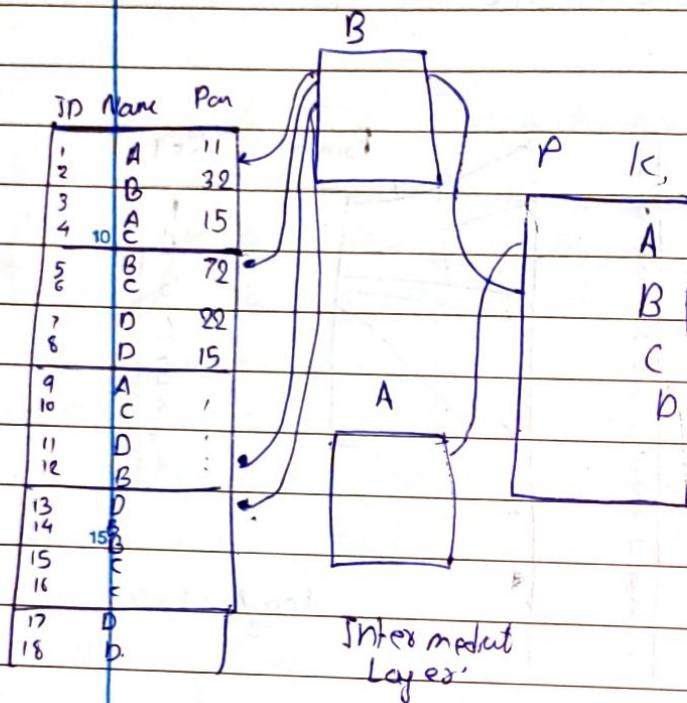
$$\log_2 N + 1$$

Index Data
Search Search

No of records in index
is equals to no of records
in hard disk.

$$\log_2 N + 1$$

N = Total No of blocks in
index table.



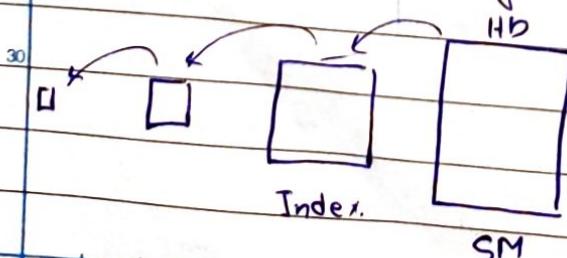
Sparse + Dense.

(Mostly Dense)

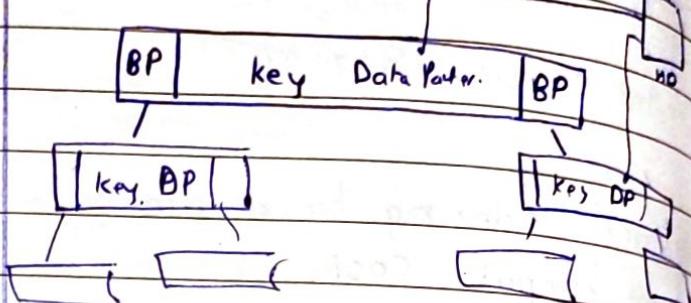
$$\log_2 N + 1 + 1$$

25
X X
B tree and B+ trees

Dynamic Multilevel Indexing:



B trees are used for dynamic multilevel indexing
Key points to remember:
TP/BP \rightarrow Block Pointer.



BP Called as block pointer
which is number of children.
Also called as Tree Pointer

Keys: key denotes the node value
and DP denotes data pointer also
called as record pointer.

Record Pointer points to Secondary
memory where the data is
stored.

Whereas, BP points in the
tree for child node

* Formula

• Order : Order $P = BP = \text{Maximum}$
number of children.

• Key = Order - 1, or
 $= BP - 1$

Create B tree order = 4.

Find maximum number of children, keys, minimum key.

1 minimum number of key

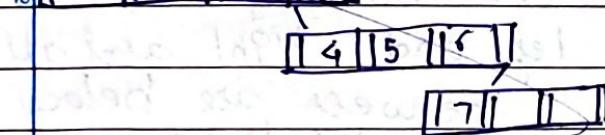
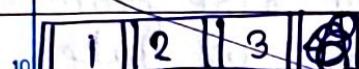
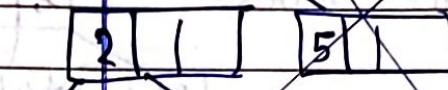
$$\min \text{ no of children} = \lceil \frac{P}{2} \rceil = 2$$

$$\min \text{ no of key} = \lceil \frac{P}{2} \rceil - 1 = 1$$

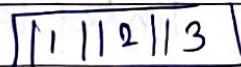
Keys are 1, 2, 3, 4, 5, 6, 7, 8, 9, 10

4/2s

Page :
Date :

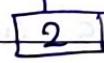
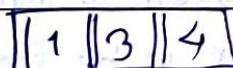


1 3 2



15

4 3 2



20

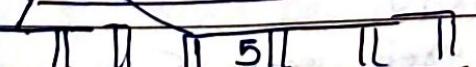
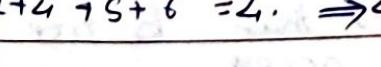


$$1+2+3+4/2 = 2.5 \rightarrow ?$$

25



30



* Consider a B-tree with key size = 10 bytes block size = 512 bytes
Data pointer 8 bytes
Block pointer 5 bytes

Find the order of b-tree

$$n \times BP + (n-1)k + (n-1)BP \leq 512$$

$$n \times 5 + (n-1)10 + (n-1)8 \leq 512$$

$$5n + 10n - 10 + 8n - 8 \leq 512$$

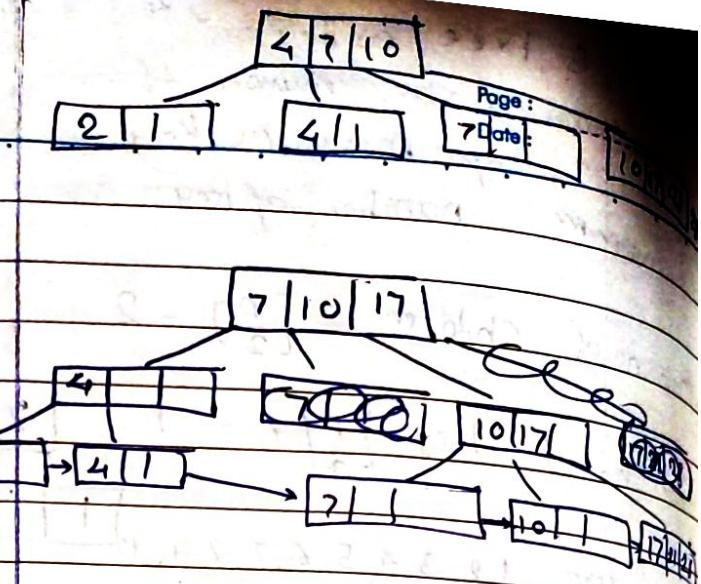
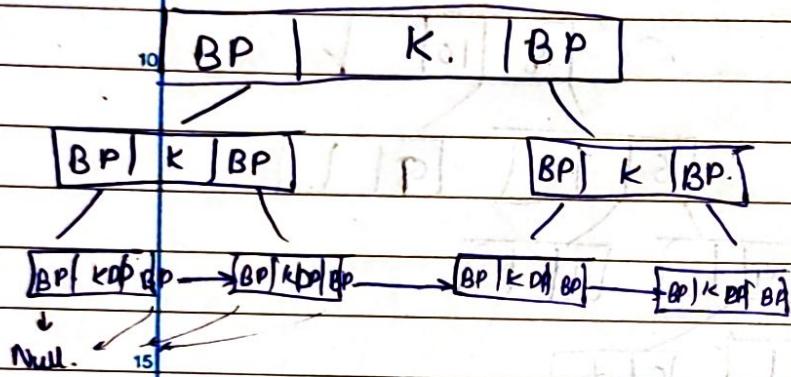
$$23n \leq 530 \Rightarrow n \leq 23$$

$$\frac{1}{9} \leq \frac{b-a}{a}$$

* B+ trees

Data Pointer DP is only present in leaf nodes.

- B trees grow height wise
B+ trees grows widthwise.



When you rise from root node just take left and right and all between are below that is the tree.
(Don't forget linked list arrows)

- LHS Child is lesser
RHS is greater or equals to.
- All keys and Data pointers are present in leaf nodes
- Each nodes in leaf node level are connected by linked list.

Q keys 2, 4, 7, 10, 17, 21, 28
order = 4.

B+ tree is used to implement database index in B+ trees
leaf nodes denotes actual data pointers.

All leaf nodes remains at same height or level
The leaf nodes are linked using linked list.

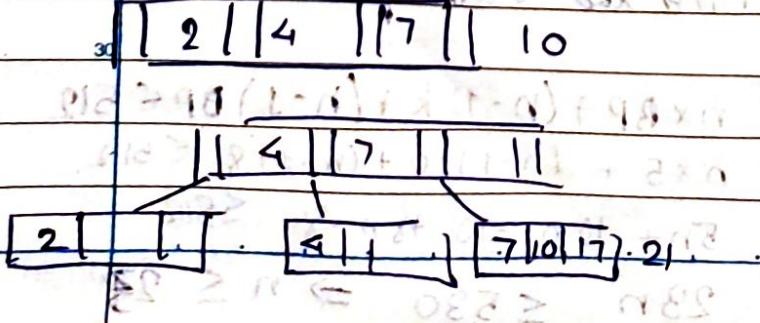
It supports random and sequential access
Internal nodes store key and block pointers

* Advantages

- faster than B tree
- easy access support random & sequential access

* Disadvantages

- req. more space



Consider a B+ tree with
 key size = 10 bytes
 block size = 512 bytes
 $DP = 8$ bytes & $BP = 5$ bytes

Find the order of leaf
 and non leaf nodes.

For leaf order is maximum
 number of key + 1
 data pointer as pair

* Difference between B trees &
 B+ trees.



$$\downarrow \\ S$$

512.

Non leaf (Internal)

$$5n + (n-1)10 \leq 512$$

$$5n + 10n - 10 \leq 512$$

$$15n \leq 522$$

$$n \leq 34.8$$

$$\underline{n \leq 34}$$

$$\text{Keys} = 33$$

leaf node



$$X(K+DP) + BP \leq 512$$

$$X(10 + 8) + 5 \leq 512$$

$$(10X + 8X) \leq 507$$

$$X \leq 28.16$$

$$\underline{X \leq 28}$$

For Non leaf order is maximum
 number of block pointer which
 is maximum number of child

HW

10 20 30 50

20 50 60 70

10 1 30 1

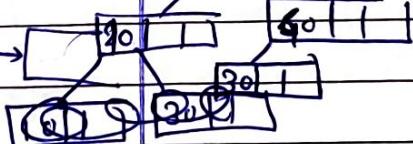
50 70 80 80

20 1 60 1

10 1 30 1

50 70 80

20 1 60 1



$$X(K+DP) + BP \leq 512$$

$$X(10 + 8) + 5 \leq 512$$

$$(10X + 8X) \leq 507$$

$$X \leq 28.16$$

$$\underline{X \leq 28}$$

10 20 40 50

20 1 1 string stop

10 1

40 50 60 70

20 50

10 1

40 1

60 70 80

20 50 1

10 1

40 1

60 70 80

30 1

20 50 1

10 1

40 1

60 70 80

51 1
30 1 25 1

20

25

30

20 50 70

Page :

Date :

51 101

30 50 70

60 70

80 90

Cascading VS Cascadeless Schedule.

In DBMS, schedules define the order of executions, ensuring data consistency.

1. Cascading Schedule :

- In a Cascading Schedule, the rollback of one transaction can lead to the rollback of other dependent transaction.
- If a transaction T_1 gets failed all other transactions dependent on T_1 (working on same data item as T_1) are aborted.

	T_1	T_2	T_3
	$R(A)$		
	$W(A)$		
		$R(A)$	
			$R(A)$
rollback	$R(B)$		
	$W(B)$		
	$Fail.$		
			$\rightarrow abort$

2. Cascadeless Schedule :

- In a cascadeless Schedule, transactions are not allowed to read uncommitted data from other transaction. This eliminates the possibility of cascading rollback.

ex.

<u>T₁</u>	<u>T₂</u>
R(A)	
W(A)	
Commit	R(A) X Not allowed till Commit
	R(A) ✓ Allowed

10 Cascade

- If there is a read after write operation on same data item but in different transaction then after failure you have to abort all the transactions which comes after write.

- CPU utilization is not proper as all transactions are aborted and rollbacked.

- Performance is low

Cascadeless

- You are not allowed to read a data item after write on same account (data item) in different transaction until it commits, fails or aborts.

- CPU utilization is relatively better

- Performance is better

* Strict Schedule :

- It states that after write operation the transaction must commit or fail.

- It guarantees recoverability and consistency by imposing strict rules on timing of read and write operation.

- A transaction can read or write a data item only after the transaction that last wrote to the data item has committed or aborted.
- No transaction can read or write from another transaction that is still in progress.

	T ₁	T ₂
10	R(A)	
		R(A)
commit of T ₁	W(A)	

Commit

* Serializability:

- It is checking of ability of a schedule to become serializable.
- A non-serial schedule is said to be serializable if it can be transformed to its corresponding serial schedule.

	T ₁	T ₂	T ₁	T ₂
25	R(A)			R(A)
	W(A)			W(A)
		R(A)	R(A)	
		W(A)	W(A)	

Serial.

Serializability

Conflict

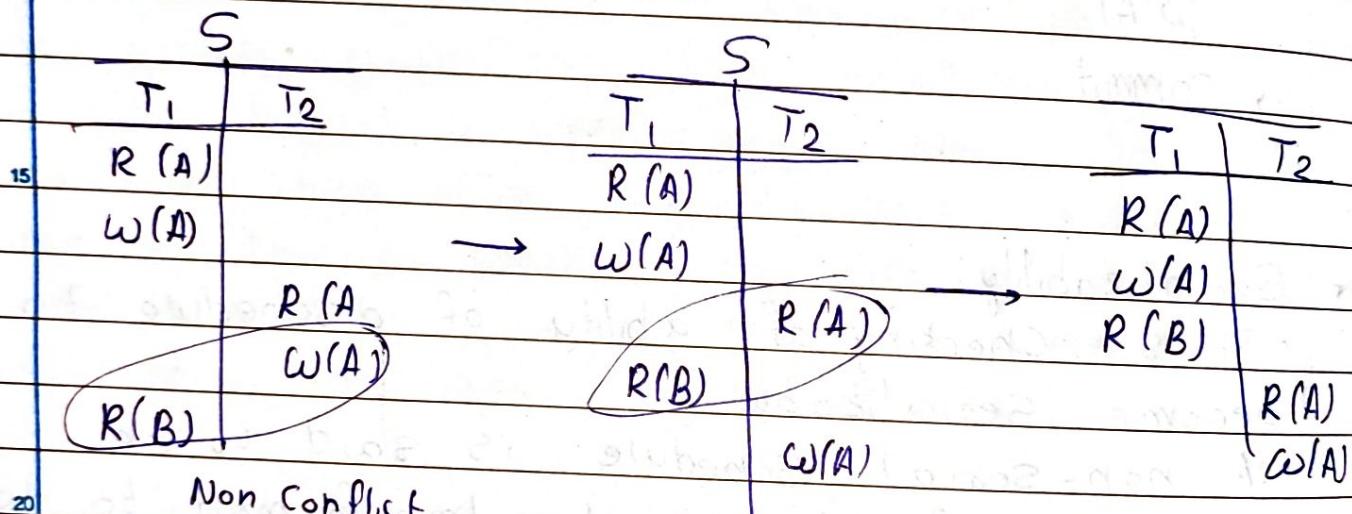
View

Camlin

There are two ways to check whether any non serial Schedule is serializable

1. Conflict Serializability

- In Conflict Serializability method if the adjacent pairs of transactions are non conflict then we can swap their position
- A Schedule is conflict serializable if it can be transformed into a serial schedule by swapping non-conflict operations.



- It uses a precedence graph (or conflict graph) to test serializability.
- The Schedule is ^{conflict}serializable, if precedence graph has no loops.
- We cannot say if the Schedule is serializable or not. We can just say it is not ^{conflict}serializable if loop is found.
- After loop is found we have to check by view serializability.

2. View Serializability

- It is a condition to ensure that the outcome of executing transactions concurrently is equivalent to some serial execution of those transactions.
- A Schedule will view serializable if it is view equivalent to a serial schedule.
- If a Schedule is conflict serializable, then it will also be view serializable.
- View Serializability focuses on the read and write operations of transactions and ensures that the views of the data seen by transactions remain consistent with a serial execution.

15

- * View Equivalence : Two Schedules are view equivalent if
 - Initial read match : for every transaction, if a read operation reads a value written by a previous transaction
 - Final Write Match : The final write of each data item in both schedules must be by the same transaction.
 - Intermediate Reads Match : If a transaction reads a value written by another transaction, the same write-read relationship must exist in both schedules.
- * View Serializable Schedule : A Schedule is view Serializable if it is view equivalent to a serial schedule.

* Initial Read :

S ₁		S ₂	
T ₁	T ₂	T ₁	T ₂
R(A)			W(A)
	W(A)	R(A)	

- 5 above two Schedules are view equivalent
as initial read by S₁ & S₂ is done by T₁. For A in both.

* Update Read :

S ₁			S ₂		
T ₁	T ₂	T ₃	T ₁	T ₂	T ₃
W(A)				W(A)	
	W(A)		W(A)		
		R(A)			R(A)

Here S₁, S₂ are not view serializable

- 15 as in S₁, T₃ is reading A updated by T₂
but in S₂, T₃ is reading A updated by T₁

* Final Write

- 20 A final write must be the same between both the schedules

S ₁			S ₂		
T ₁	T ₂	T ₃	I	T ₂	T ₃
W(A)					
	R(A)			R(A)	
		W(A)	W(A)		C(A)

S₁, S₂ are view serializable as final write for A is done by T₃ in both.