

# RDBMS and SQL

Madhavan Mukund

<https://www.cmi.ac.in/~madhavan>

Lecture 10, 29 October 2024

# Storing data

- RAM vs Hard disk vs SSD

- Blocks and latency

$10^9$  ops/sec

read  $10^9$  items/sec

?

$10^3$

$10^6$

delay  
transfer  
seek

large volume  $\rightarrow$  disk

1.4 billion people  
 $\times$  100 bytes

$140 \times 10^9$  bytes

140 GB

# Storing data

- RAM vs Hard disk vs SSD
- Blocks and latency



transfer a chunk of data at a time

Block is 4kb or more

"Accounting" – analysis

Ignore memory operations

Only count block seeks & transfers

# Storing data

- RAM vs Hard disk vs SSD
- Blocks and latency

Typical table will be  
split across blocks

File organisation

—how are these multiple  
blocks of one table  
arranged

# Fixed length records

- Blocks and block boundaries

*row in table* record 0  
record 1  
record 2  
delete? → record 3  
record 4  
record 5  
record 6  
record 7  
record 8  
record 9  
record 10  
record 11

10101	Srinivasan	Comp. Sci.	65000
12121	Wu	Finance	90000
15151	Mozart	Music	40000
22222	Einstein	Physics	95000
32343	El Said	History	60000
33456	Gold	Physics	87000
45565	Katz	Comp. Sci.	75000
58583	Califieri	History	62000
76543	Singh	Finance	80000
76766	Crick	Biology	72000
83821	Brandt	Comp. Sci.	92000
98345	Kim	Elec. Eng.	80000

*~~~~~*

# Deleting a record

- Compress

record 0	10101	Srinivasan	Comp. Sci.	65000
record 1	12121	Wu	Finance	90000
record 2	15151	Mozart	Music	40000
record 4	32343	El Said	History	60000
record 5	33456	Gold	Physics	87000
record 6	45565	Katz	Comp. Sci.	75000
record 7	58583	Califieri	History	62000
record 8	76543	Singh	Finance	80000
record 9	76766	Crick	Biology	72000
record 10	83821	Brandt	Comp. Sci.	92000
record 11	98345	Kim	Elec. Eng.	80000

# Deleting a record

- Compress
- Move last record

record 0	10101	Srinivasan	Comp. Sci.	65000
record 1	12121	Wu	Finance	90000
record 2	15151	Mozart	Music	40000
record 11	98345	Kim	Elec. Eng.	80000
record 4	32343	El Said	History	60000
record 5	33456	Gold	Physics	87000
record 6	45565	Katz	Comp. Sci.	75000
record 7	58583	Califieri	History	62000
record 8	76543	Singh	Finance	80000
record 9	76766	Crick	Biology	72000
record 10	83821	Brandt	Comp. Sci.	92000

# Deleting a record

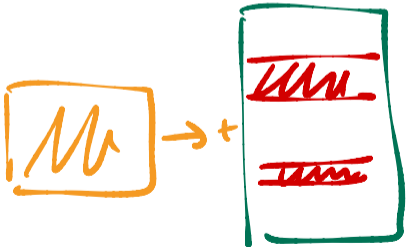
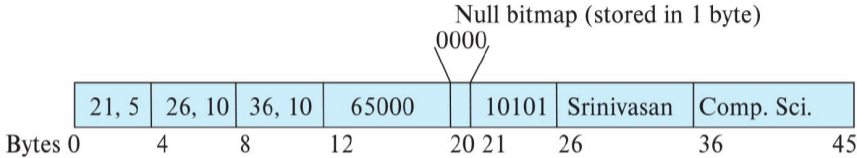
- Compress
- Move last record
- Maintain **free list** of empty slots

header				
record 0	10101	Srinivasan	Comp. Sci.	65000
record 1				
record 2	15151	Mozart	Music	40000
record 3	22222	Einstein	Physics	95000
record 4				
record 5	33456	Gold	Physics	87000
record 6				
record 7	58583	Califieri	History	62000
record 8	76543	Singh	Finance	80000
record 9	76766	Crick	Biology	72000
record 10	83821	Brandt	Comp. Sci.	92000
record 11	98345	Kim	Elec. Eng.	80000



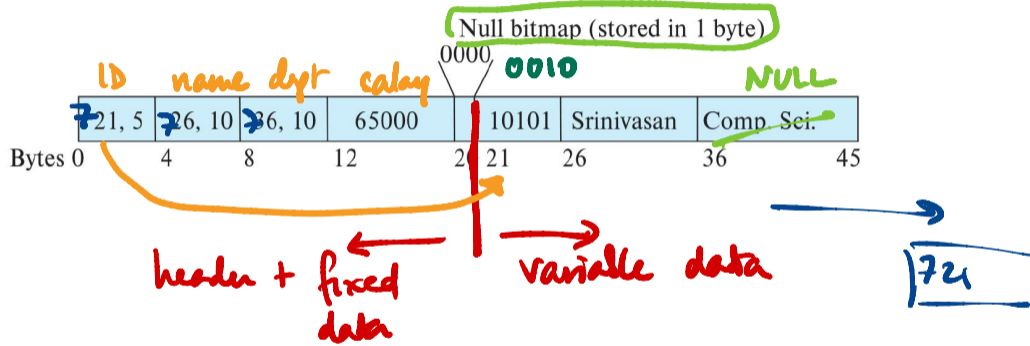
# Variable length records

- Single record structure



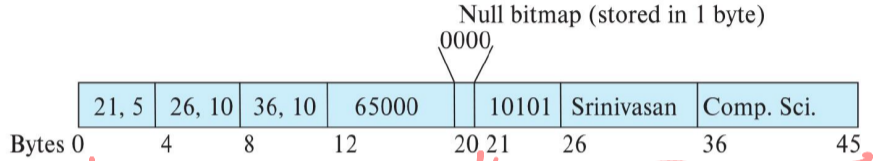
# Variable length records

- Single record structure



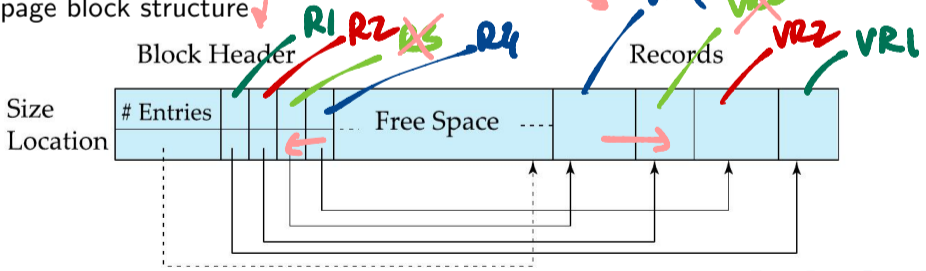
# Variable length records

- Single record structure



- Slotted page block structure

**Delete R3**



# Storing tables — heap file organization

- Use first available free slot

*Arbitrary order*

# Storing tables — heap file organization

- Use first available free slot
- Maintain free space map

Insert row with 6 units of space



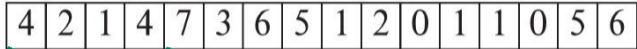
Block 1    Block 2

---

Block 16

# Storing tables — heap file organization

- Use first available free slot
- Maintain free space map



4 blocks

- Second level free space map



largest space in range of blocks

should fit in 1 block

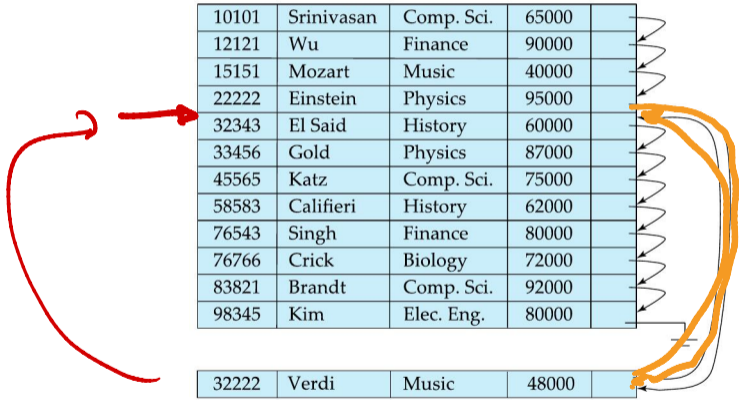
# Storing tables — sequential file organization

Sorted ↴

10101	Srinivasan	Comp. Sci.	65000	
12121	Wu	Finance	90000	
15151	Mozart	Music	40000	
22222	Einstein	Physics	95000	
32343	El Said	History	60000	
33456	Gold	Physics	87000	
45565	Katz	Comp. Sci.	75000	
58583	Califieri	History	62000	
76543	Singh	Finance	80000	
76766	Crick	Biology	72000	
83821	Brandt	Comp. Sci.	92000	
98345	Kim	Elec. Eng.	80000	

# Storing tables — sequential file organization

- Overflow block





- Why build an index?

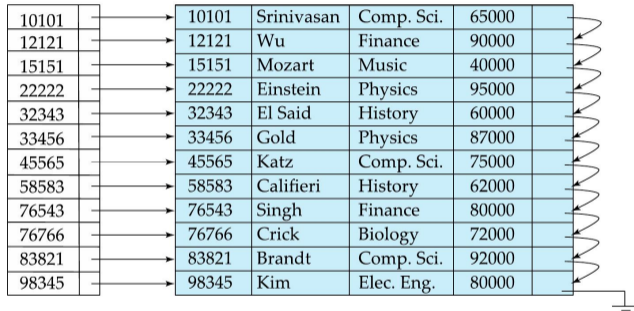
- Why build an index?
- Search key
  - As opposed to superkey, candidate key, . . .
  - May need multiple search keys for a table

- Why build an index?
- Search key
  - As opposed to superkey, candidate key, . . .
  - May need multiple search keys for a table
- Types of queries — point vs range
  - `ID = "10102"`
  - `salary > 75000`

- Why build an index?
- Search key
  - As opposed to superkey, candidate key, . . .
  - May need multiple search keys for a table
- Types of queries — point vs range
  - `ID = "10102"`
  - `salary > 75000`
- Maintaining an index
  - Inserts, deletes
  - Space

# Clustering index

- File is ordered with respect to index values
- **Index sequential file**
- Dense index — every value is present in the index



# Clustering index

- File is ordered with respect to index values
- **Index sequential file**
- Dense index — every value is present in the index
  - Index value may match multiple records

Biology		76766	Crick	Biology	72000
Comp. Sci.		10101	Srinivasan	Comp. Sci.	65000
Elec. Eng.		45565	Katz	Comp. Sci.	75000
Finance		83821	Brandt	Comp. Sci.	92000
History		98345	Kim	Elec. Eng.	80000
Music		12121	Wu	Finance	90000
Physics		76543	Singh	Finance	80000
		32343	El Said	History	60000
		58583	Califieri	History	62000
		15151	Mozart	Music	40000
		22222	Einstein	Physics	95000
		33465	Gold	Physics	87000

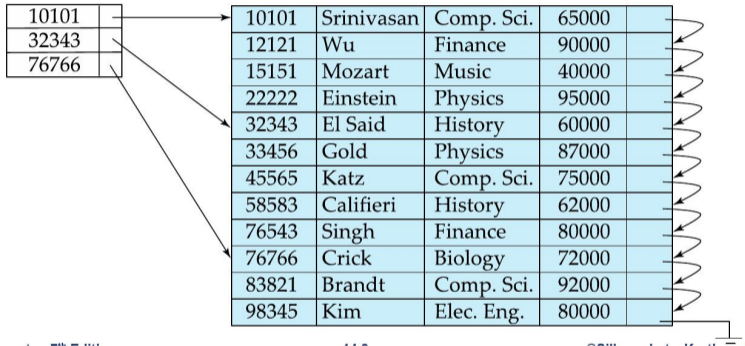
# Indexing — sparse indices

- Maintain indices for a subset of values
- Page headers in a dictionary



# Indexing — sparse indices

- Maintain indices for a subset of values
  - Page headers in a dictionary
- Align to block boundaries
  - Records are still sequential with respect to index
  - Sparse index identifies first record in each block

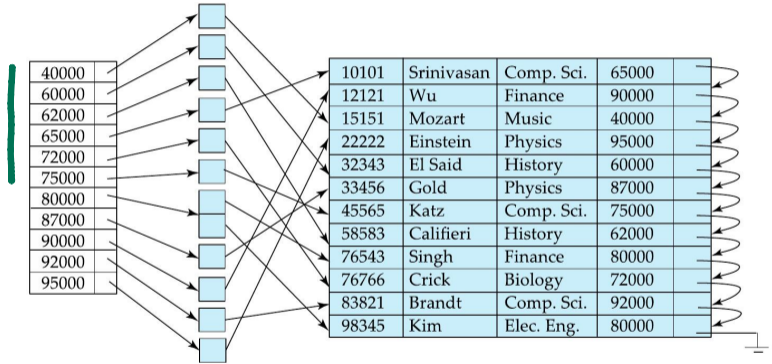




# Indexing — secondary index

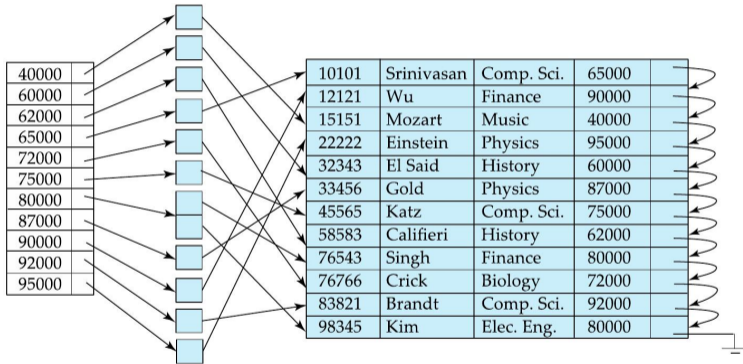
- Index for an attribute that does not match sequence in which table is stored

Salaries  $\leq$  75000



# Indexing — secondary index

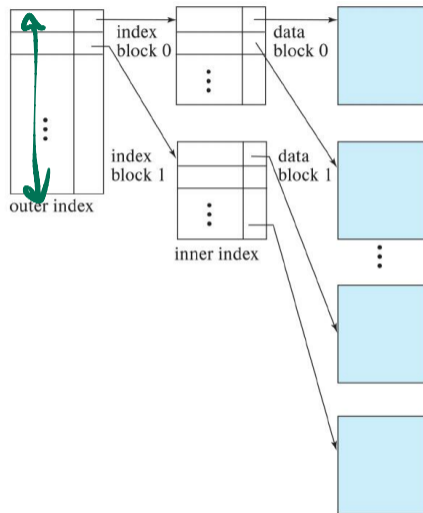
- Index for an attribute that does not match sequence in which table is stored
- Key points to block that contains pointers to matching records
  - Can have multiple records for same search key



- Typically, index will not fit in RAM

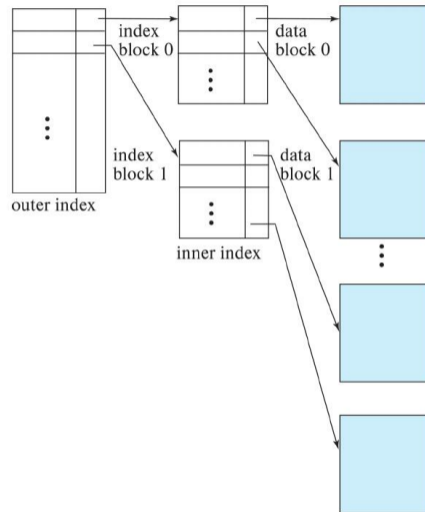
# Storage

- Typically, index will not fit in RAM
- Store index as a sequential file
  - Build a sparse index for the index file
  - Multi-level, till sparse index fits in one block



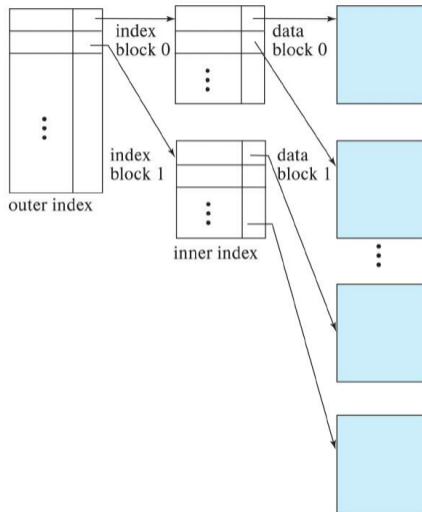
# Storage

- Typically, index will not fit in RAM
- Store index as a sequential file
  - Build a sparse index for the index file
  - Multi-level, till sparse index fits in one block
- Binary search to find required key



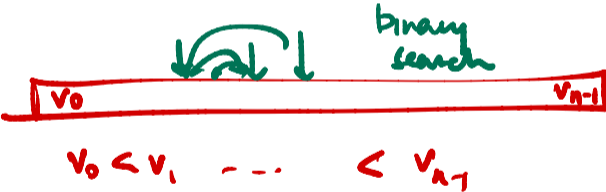
# Storage

- Typically, index will not fit in RAM
- Store index as a sequential file
  - Build a sparse index for the index file
  - Multi-level, till sparse index fits in one block
- Binary search to find required key
- Idea leads to a more efficient structure

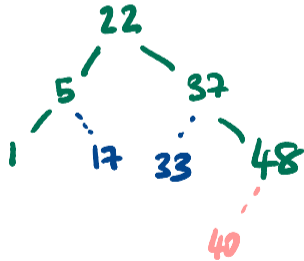
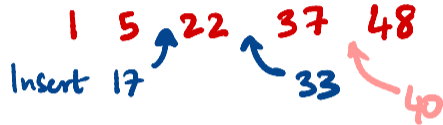


# Search trees

- Binary search trees
  - Binary search on dynamic data
  - Balanced tree has logarithmic height



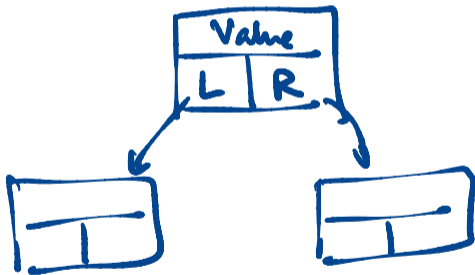
## Dynamic binary search



Binary search tree

# Search trees

- Binary search trees
  - Binary search on dynamic data
  - Balanced tree has logarithmic height





# Search trees

- Binary search trees
  - Binary search on dynamic data
  - Balanced tree has logarithmic height

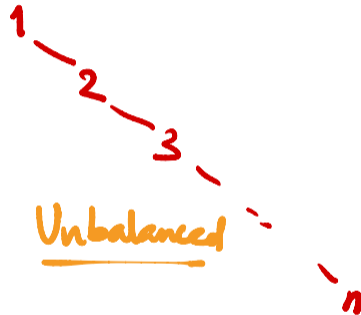
Bad case

Insert 1

Insert 2

⋮

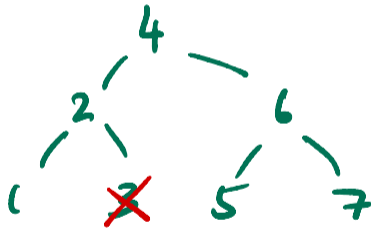
Insert n



# Search trees

- Binary search trees
  - Binary search on dynamic data
  - Balanced tree has logarithmic height

Perfectly balanced  
 $2^n - 1$  element



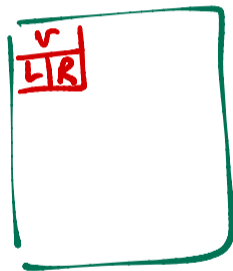
- Binary search trees
  - Binary search on dynamic data
  - Balanced tree has logarithmic height

Different notions of  
"approximate" balance

→  $n$  elements  
height is  $O(\log n)$   
↓  
longest path

- Binary search trees
  - Binary search on dynamic data
  - Balanced tree has logarithmic height

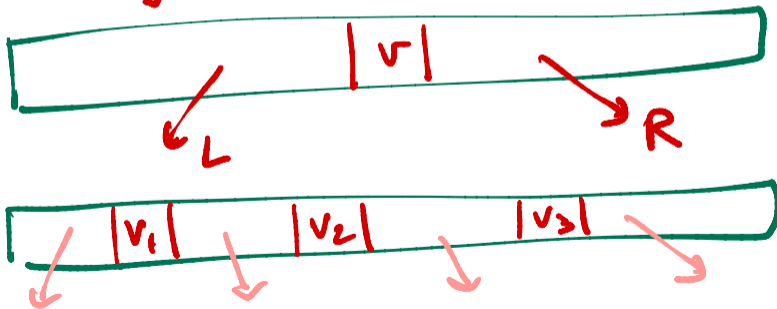
"Binary tree" on disk  
each node = 1 block



# Search trees

- Binary search trees
  - Binary search on dynamic data
  - Balanced tree has logarithmic height

Instead, multway branch



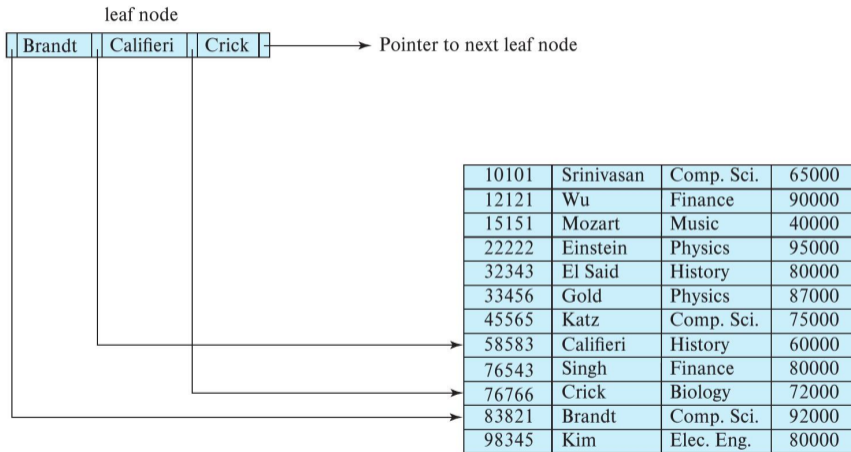
- Binary search trees
  - Binary search on dynamic data
  - Balanced tree has logarithmic height
- Block-based access
  - Binary tree node has one search key value, two pointers
  - Block can hold much more

- Binary search trees
  - Binary search on dynamic data
  - Balanced tree has logarithmic height
- Block-based access
  - Binary tree node has one search key value, two pointers
  - Block can hold much more
- Generalize to multiple key values, multiple pointers



# B+ trees

- Leaf nodes form a dense index — linked list of leaves, each one block



*instructor file*



# B+ trees

- Leaf nodes form a dense index — linked list of leaves
- Non-Leaf nodes form a sparse index

