# SQL

(5) SELECT [Distinct] col1, col2, agg(col3)
(1) FROM table1, table2,
(2) WHERE <predicate> AND <predicate> OR
(3) GROUP BY <column list> ← Everything in SELECT must be in GROUP-BY or is an aggregate
(4) HAVING <predicate>
(6) ORDER BY <columns> [DESC] [ASC]
(7) LIMIT <integer>;

Must contain only GROUP-BY columns of aggregate functions

## SQL String Comparison

• Old School SQL     starts with B
  WHERE S.name LIKE 'B_ob' ← return Bob
  _ = any single char; % = 0+ chars
• Standard Regular Expressions
  WHERE S.name ~ 'B.*' ← returns Bob, McBob
  . = any char, * = repeat (0+ instances of previous)

## SQL Join Variants

From table1
  [ INNER | NATURAL | {LEFT | RIGHT | FULL}
  OUTER ] JOIN table2
  ON <qualification list>
• INNER: join tables where "ON" qualification
• NATURAL: join tables for pairs of attributes
              w/ same name

| FROM t1, t2 WHERE t1.id=t2.id | = | FROM t1 INNER JOIN t2 ON t1.id = t2.id |

FROM t1 NATURAL JOIN t2
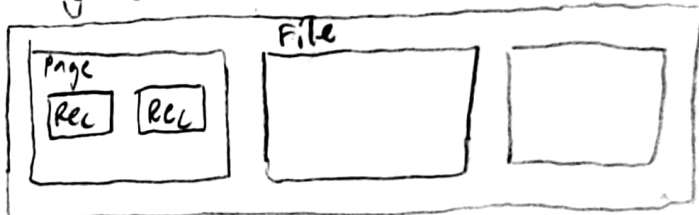
(assume only matching column names = "id")

• LEFT/RIGHT/FULL OUTER JOIN
  FROM t1 [ ] OUTER JOIN
  ON t1.id = t2.id
  - LEFT: if t1.id has no match, t2.id is NULL
        t1 ⬤) t2
  - Right: if t2.id has no match, t1.id is NULL
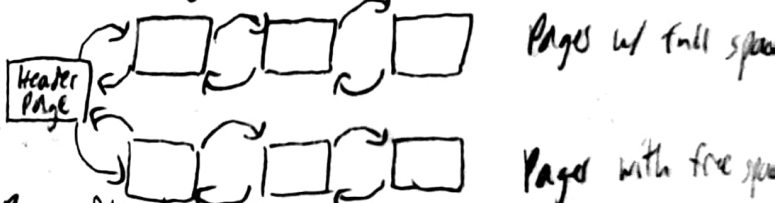        t1 (⬤ t2
  - FULL: LEFT and RIGHT both apply  t1 (⬤) t2

---

# Relational Table Properties   [No aggregates]

• Schema is fixed
  - unique attribute names
  - atomic (primitive) types
• Tables are not ordered (sets, multisets)
• Tables are flat (no nested tables)
  - First Normal Form

## Disk Representation     [pageID, location on page]

• DB File: collection of pages
• Pages: collection of records

File
Page
REC  REC

## Unordered Heap Files

Records placed arbitrarily across pages
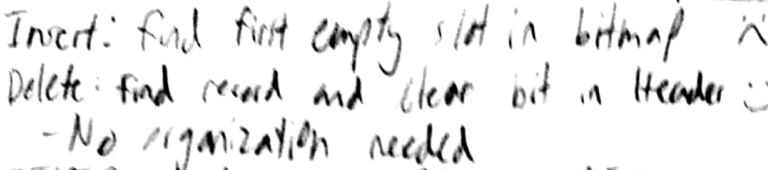1. File as Doubly Linked List

Header Page

Pages w/ full space

Pages with free space

2. Page Directory
Header Pages

<□,10>, <□, 3>          Data page w/ 3 free space   Data Page

<□, 13>...                Data Page

## Unordered Heap Files vs. Sorted File

| 2,5 | 1,6 | 4,7 | 3,10 | 8,9 | vs. | 1,2 | 3,4 | 5,6 | 7,8 | 9,10 |

| | Heap | Sorted | Assume |
|---|---|---|---|
| Scan all records | B·D | B·D | - append page to end of heap file |
| Equality search | 0.5·B·D | (log B)·D | |
| Range Search | B·D | (log B + pages)·D | |
| Insert | 2·D | (log B + B)·D | - Sorted files are packed |
| Delete | (0.5·B)·D | (log B + B)·D | |

# Page Layout

- FIXED record length, UNPACKED records

| Header | 0 | 1 | 0 | 0 | | Rec | X | Rec | X | Rec |
|---|---|---|---|---|---|---|---|---|---|---|

Insert: find first empty slot in bitmap ☆
Delete: find record and clear bit in Header ☺
  - No organization needed
- FIXED record length, PACKED RECORDS

| Header | Rec | Rec | Rec | Rec |
|---|---|---|---|---|

Insert: just append to end
Delete: Scan for record, delete, reorganize ☹
    all records (we pack)

VARIABLE record lengths, UNPACKED records
(slotted page)



| Rec 8 | Rec 6 | Rec 5 | Rec 7 |   →  footer   | 18 6 5 7 ☐ |

free space pointer

Delete: remove footer pointer (set record slot to null)
Insert: First empty pointer slot if free space pointer
NOTE: footer fixes one int and one pointer for EACH page
record length

- Fixed length:

|   | int | char(7) |   |
|---|---|---|---|
| 4 | 8 | 1 4 | 7 |

→ store 24 bytes

int double int

- Variable length:

(Primary Key)  (Not NULL)

| Header ☐☐ | 0 | 1 | Varchar(10) | int | double | int | text | Varchar(13) |

- 2 pointers to variable length attribute
- bitmap (2 bits) where int and text are NULL

# Hashing



Disk    B main memory buffers    Disk

---

# Index Files (B+ tree) — acts similar to B-tree

- Property: # entries per node. $d = \text{order} \leq \text{entries} \leq 2d$
- Alternative 1: record contents
    stored in the leaf node itself
- Alternative 2: Leaf nodes: ⟨k, rid of matching record⟩
- Alternative 3: Leaf nodes:
    ⟨k, list of rids of records⟩

- Property 1: leaf node entries: $d = \text{order} \leq \text{entries} \leq 2d$
- Property 2: all leaves same dist. from root
- Property 3: inner node w/ k items have k+1 children
- Leaf split    insert(8*)

| 4* 5* 6* 7* | → 



| 6 |

| 4* 5* |  | 6* 7* 8* |

not kept if inner node

- Bulk Load
  - Input sorted record keys
    (1*, 2*, 3*, ...)
  - Fill leaf page to fill factor > d
  - Update inner pages to full
  - Follow leaf and inner page split rules
    sorting

- General External Merge Sort
  - N pages to sort, B buffer pages



Disk    PASS 0    sorted runs (length B each, int variable)    Pass 1, 2,    sorted runs (length = B(B-1))

  - Number of passes: $1 + \lceil \log_{B-1} \lceil N/B \rceil \rceil$
  - Total I/Os = (I/Os per pass) × (# passes)
    $= 2N * (1 + \lceil \log_{B-1} \lceil N/B \rceil \rceil)$
  - Memory: B(B-1) after 2 passes
① Divide: split pages into B-1 partitions
    using hp hash function
② Conquer: Repeat for big partitions
    (> B pages), rehash w/ new
    hash function hp,
③ When partition fits in memory
    (<= B pages), read + write to
    check (Build Pass)