# 1. SQL

DDL: data definition language
 CREATE TABLE
  data types: text, float, integer, char
DSL: data manipulation language
 SELECT [DISTINCT] <column expression list>
   FROM <single table>
   [WHERE <predicate>]
   [GROUP BY <columns>
     [HAVING <predicate>]
   ]
   [ORDER BY <column list>]
   [LIMIT <integer>]

- when using group by, must have same regular expression in the select statement, or put it inside aggregate

ORDER of reading SQL logic:
 FROM, WHERE, SELECT, GROUP BY && HAVING, DISTINCT, ORDER BY, LIMIT
- if a WHERE clause evaluated to null, not in output
- aggregates can't be in WHERE clause

CROSS JOINS    - can't have aggregation in WHERE clause
FROM A, B
 • cross product. Every row in left with every row on Right

INNER JOIN
same as CROSS JOIN but need ON.

FULL OUTER JOIN
if either table no match, its still there

OUTER JOIN
- even if ON predicate doesn't match, still in output
  LEFT means
   - still if no match, left columns table are still there with null
  RIGHT
   - vice versa    - Duplicate key in Primary key not allowed by definition

Facts
- two same queries not deterministically the same if there's a tie

# Relational Algebra

• no duplicates
• all operators take in relation and output different relation

Projection (π)
 SELECT: select only columns specified

Selection (σ)
 WHERE clause:
  $\sigma_{age=12 \wedge name = 'SAM'}$

Union (∪)
 - must be same columns

(−) set diff

(∩) intersection

Joins (⋈)
no specification = natural join (join on tables w/ same names)

Rename (ρ)
 cats ⋈ name = dname  pname → dname (dogs)

Pattern Matching
 LIKE
  - if no % or _, then ' ' acts like equality
  - an underscore stands for matches single character
  - % matches any sequence of zero or more chars.

standard regex
 ~  • (.) represents single wildcard character
    (*) represents repetition of prev item zero or more times  a·c = abc valid

# 2. Disks, Files, Buffers

## Disk vs Flash (SSD)
 Disk: Accessing a page
  - seek time, rotational delay, transfer time
  • Random read much slower than sequential
  • random write faster than sequential write
 Flash: random write faster than sequential
  - faster than disk for lo-lox random IOs
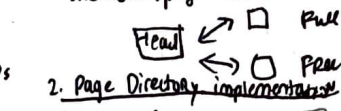 • Locality matters for both
 • Disk 10x/capacity/$ dollar
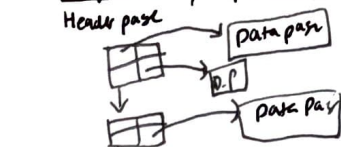
Files > Pages > Records

# File Structures

Heap File — no particular ordering
 1. Linked List implementation
  • Each data page has records, free space tracker, pointers
  • one header page as start


Head ⇄ □ Full
     ⇄ ○ Free

I/O
reading all data pages
in one header is only 1 I/O.

 2. Page Directory implementation
  Header page → data page
             → p.r
             → page page
  • pointer to data page and free space info about that page
  - faster insertion in terms of I/Os
    • read header, read data, write data, write header
  • fast insertion
  • slow search

Sorted File — pages are ordered and records on pages sorted by keys
 • implemented w/ Page Directories
 • LogN search
 • LogN+N insert because shift

Record Types
 Fixed length — only fixed types and all same length records
 Variable Length — variable length. Fixed length fields before variable length fields and header has pointers to the end of the variable length field
  Record ID: [page #, record # on page]

Page Formats
 1. Pages w/ FLR: page header to store # of records
   • Packed.
   • Unpacked: bitmap in header
 2. Pages w/ VLR:
   - page footer that maintains slot directory tracking
     (. slot count, free space pointer, entries
     4 bytes + 4 bytes + 8# # records    ↳ [record pointer, record length]

Heap vs Sorted
 Heap: good for frequent queries, inserts, updates. fine for frequent full scan
 Sorted: good for range searches, frequent lookup

# 3. B+ Trees/Indexes

Index: data structure that allows fast lookup, to certain key

B+ Tree Properties
 • d: order of tree. Each node except Root must have sorted $d \le x \le 2d$ entries assuming no delete.
 • Inner Node: 2d entries, 2d+1 children ptr. (tree fanout)
 • The keys in the children to the left of an entry must be < To Right ≥ than.

INSERTION
 1. Find leaf node to insert. Add key, record to leaf
 2. if overflow (> 2d):
   a) split into L, L2- D in L1 D+1 in L2
   b) if L is leaf, COPY L2's first entry into parent. else, move
   c) adjust pointer
 3. if parent overflow, recurse on it w/ step 2.

DELETE
 - just delete in leaf

Total capacity: $(2d)(2d+1)^h$ where h = # edges from R to L

> # ONLY FOR LEAF
> FillFactor is for preventing splits on new inserts, not IO.

Bulkload — build from scratch
 1. Sort data on key of index
 2. Fill leaf pages till ff. f.
 3. Add ptr from parent to leaf. If parent overflows, follow split parent
   a) keep d in L1. d+1 in L2
   b) move L2 first entry up.
 4) adjust ptrs

Storing Records
 Alt 1: leaf pages are records themselves
 Alt 2: leaf pages are pointers to corresponding record
 Alt 3: linked list of pointers to corresponding records

Clustering
 1. Unclustered


 2. Clustered

  - Better caching
  - when multiple query, we might have that page already since sequential.

Counting I/Os
 1. Read appropriate root to leaf. one I/O per node.
 2. Read appropriate data pages. IO per page. (account for clustering
 3. Write data page if modifying. if we want to write that spans multiple pages, IO per page
 4. update index page (node)

 1 I/O per linkedlist data page
  - when searching keys, there's not always a datapage associated so BEST CASE is NOT bringing I/O for data.

# 5. Sorting

Goal: sort pages in Disk

## External Merge Sort

Assume N data pages in Disk, B buffer pages available

1. Load B buffer pages fully and sort that group. do that for all $\lceil N/B \rceil$ sorted RUNS.

2. Merge B-1 pages recursively

$$O(2N \cdot \underbrace{(1 + \lceil \log_{B-1} \lceil N/B \rceil \rceil)}_{passes})$$

Two Way: $2N \cdot (1 + \lceil \log_2 N \rceil)$

# 5. Hashing

Goal: Group same values pages in Disk

## External Hashing

Assume B Buffer frames available

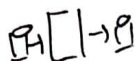### Divide Phase
- 1 input buffer
- B-1 output buffers / Partition
- group B-1 partitions into disk



### Conquer Phase
- Start when all partitions fit in B pages
- use fine grained hash function



Fact: 2 pass if $B(B-1)$ pages

If some of the pages/partitions are $\geq B$, first finish the pass and in the next pass recursively partition just that partition. So split that into B-1 as usual.

- Usually write more pages after reading *. Round up. So, 12 data pages, 6 buffer pages, each of 5 frames get 12/5 = 3 each. So write 15. and read 12
- Don't forget last pass of Reading/writing all the pages after fitting.

# 4. Buffer Management

- buffer manager works with disk space manager and RAM.
- between RAM and DISK
- B+ code makes read requests from Buffer Manager
- Metadata: FrameID, PageID, DirtyBit, PinCount

## Buffer Pool
Memory converted to buffer pool by partitioning into frames for pages.

1. FrameID: uniquely associated w/ memory address
2. PageID: determining which page a frame contains
3. Dirty Bit: verifying if modified
4. PinCount: tracking number of requestors using a page

## Handling Page Requests

Upon page request:
  if page exists in memory:
    1. pin count ++
    2. return address of page
  else:
    if space:
      1. read page into it
      2. pin count = 1
      3. return page
    else:
      REPLACEMENT POLICY!

## LRU
Pros: good for repeated access
      wins for random access/popular access
Cons: LRU page finding requires priority queue O(logN)
      sequential flooding

## Clock (LRU)
- extra metadata: Ref Bit

### Adding Page
1. if page exists: just return it. Dont move any clock pointers. set the reference bit

2. if doesn't exist:
   1. Move clock pointer around.
      1. if not reference bit, evict, put in, set ref bit, move clock
      2. if ref bit, unref bit, move clock

## MRU
- evict most recently used unpinned page
- Pincount reduced by whoever uses the page, NOT ALWAYS BUFFERMANAGER
- file/index management code/page requestor sets dirty bit.

# Conceptual Ideas

- B+ Tree:
    API: get(key, record ID)

- Buffer Manager:
    stores pages

- Remember a B+ Tree is just defined by a "page" Root, so this is something we need to fetch from disk, and keep fetching the other pages that pointers to into memory.

# Note 10. Parallel Query Processing

def: shared nothing - every CPU has own disk + memory
def: intra-operator parallelism - make one operator run as fast as possible
def: inter-operator parallelism - make query run as fast as possible by running operators in parallel

## Partitioning

1. **Range Partitioning**
- Each machine gets a certain range of values that it will store
- PRO: good for queries that lookup on a specific key
- PRO: used in parallel sorting and parallel sort merge join

2. **Hash Partitioning**
- Each record is hashed and is sent to a machine matches that hash value.

3. **Round Robin Partitioning**
- Distribute records evenly, one by one

def: Network cost - how much data we need to send over the network

## Parallel Sorting
1. Range partition the table
2. Perform local sort on each machine

## Parallel Hashing
1. Hash partition the table
2. Perform local hashing on each machine

## Parallel Sort Merge Join
1. Range partition each table using the same ranges on the join column
2. Perform local sort merge join on each machine

## Parallel Grace Hash Join
1. Hash partition each table using the same hash function on the join column
2. Perform local grace hash join on each machine

## Broadcast Join
When one big table and one small table, just send the small table to each machine and each machine does a 'local join.'

## Symmetric Hash Join (pipeline friendly)
1. Build two hash tables, one for each table in the join
2. When a record from R arrives, probe the hash table for S for all of the matches. When a record from S arrives, probe the hash table for R for all of the matches
3. whenever a record arrives add it to its corresponding hash table after probing the other hash table for matches

# Note 11. Transactions
def: inconsistent reads - a user reads only part of what was updated
def: lost update - two users try to update the same record so one of the updates gets lost
def: dirty reads - one user reads an update that was never committed, aborted
def: transactions - sequence of multiple actions that should be executed as a single, logical, atomic unit - Guarantee ACID properties

Atomicity - transaction ends in two ways: either commits or aborts. Atomicity means that either all actions in the Xact happen, or none
Consistency - if DB starts out consistent, ends consistent after Xact
Isolation - Execution of each Xact is isolated from that of others.
Durability - if Xact commits, effect persists, must survive failures

## Ensuring Isolation Property of Transactions
Serial Schedule - run all operations of one transaction to complete before beginning the next transaction
↳ not efficient, want to interleave transaction actions.

Equivalence - 1. involve same transactions
2. operations are ordered same way with same transactions
3. Each leave the database in the same state
↳
**Serializable**
we can ensure this #3 constraint by looking for conflicting operations
1. operations are from different operations transaction
2. both operate on same resource
3. at least one is a write

- conflict serializable if conflict equivalent to a serial schedule

## Dependency Graph
- one node per Xact
- Edge from $T_i$ to $T_j$ if:
  - an operation $O_i$ of $T_i$ conflicts w/ an operation $O_j$ of $T_j$
  - $O_i$ appears earlier than $O_j$ in the schedule
- Serializable if acyclic

## Deadlock
- Waiting for each other to release

## Avoidance
Set priority by its age: now - starttime
**Pipeline Breakers**
Wait-Die: if $T_i$ has higher priority, $T_i$ waits for $T_j$; else $T_i$ aborts
Wound-Wait: if $T_i$ has higher priority, $T_j$ aborts; else $T_i$ waits

## Detection of Deadlock
"Waits for Graph": one node per XACT and an edge from $T_i$ to $T_j$ if
- $T_j$ holds a lock on resource X
- $T_i$ tries to acquire a lock on resource X but $T_j$ must release its lock on resource X before $T_i$ can acquire its desired lock
- if cycle, shoot a XACT in cycle
- start backwards, and only draw arrow if they actually got it

**2PL** can't get locks after unlocking

**S2PL**: unlock after transaction

→ Topological Sort to get equivalent serial schedule

- if two schedules order every pair of conflicting operations the same way, then they are output equivalent (conflict equivalent)

mainly equijoin and natural join

## Notation
[R] = # pages in R
[S] = # pages in S
|R| = # Records in R

## Simple Nested Loop Join
- each Record R, get each page of S
- $O(\text{join}(R,S))$ where R is the outer loop: $[R] + |R|[S]$

## Page Oriented Nested Loop Join
- for every page in R, bring in one page S
$$[R] + [R][S]$$

## Block Nested Loop Join
- B-2 pages as one block of R
$$[R] + \left\lceil \frac{[R]}{B-2} \right\rceil [S]$$

## Index Nested Loop Join
- If we have an index on S, just look it up
$$[R] + |R| * (\text{cost to look it up matching records in S})$$

## Hash Join
### Naïve Hash Join
Fit R into B-2 pages memory and then read in each Record of S and look it up in the hashtable. $[R] + [S]$ I/Os

### Grace Hash Join
1. Repeatedly hash R and S into B-1 buffers so that we can get partitions of ≤ B-2 pages
2. If both R and S > B-2 pages keep partitioning
3. when either R or S is small enough i.e ≤ B-2 pages then load smaller one into mem and create hashtable, matching against other one.     — dont care about final write

cost
partition: $2([R] + [S])$ I/O or the partitioning cost
matching: $[R] + [S]$

### Memory Requirements
- partitioning phase divides R into (b-1) runs of size $\left\lceil \frac{[R]}{B-1} \right\rceil$
- matching phase requires each $[R]/B-1 ≤ B-2$
- $R < (B-1)(B-2)$
  no S constraint
- Naïve join better for R < memory
- GJ better for $R^2$ > mem > R

### Sort Merge Join
1. Sort R and sort S and then iteratively check. $[R] + [S]$ at the last step.

## Selectivity Estimation - how much a query plan costs

### Rules
capital letter = columns
lowercase = constants
- $X = a$ : 1/(unique vals in X)
- $X = Y$ : 1/(max(unique vals in X, unique vals in Y))
- $X > a$ : (max(X) - a / max(X) - min(X) + 1)
- cond1 AND cond2: selectivity (cond 1) * selectivity (cond 2)
- 1/10 default if dont have INFO
- cond1 OR cond2 = (cond1) + (cond2) - (cond1 · cond2)

### Selectivity of Join
join A and B on condition A.id = B.id
$$[A][B] / \max(\text{unique vals for A.id, unique vals for B.id})$$

### Common Heuristics
1. push down projects ($\pi$) and selects ($\sigma$) as far as they go
2. only consider left deep plans
3. Dont consider cross joins unless they are the only option

---

## Pass 1 of System R
- Single table/condition operations
  - considers Fullscan — O[P] I/os
  - consider Index Scan
    - alt 1 : (cost Reach level above leaf) + num leaf Read
    - can stop at condition that, so find the leaf to start at, and keep going right at leaf level because its sorted

    - alt 2/3 indies:
    (cost to reach level above leaf) + (# leaf nodes read) + (num data pages read)
    - Clustered index : # data pages read is the selectivity multiplied by total # data pages
    - unclustered index: I/O per each record so selectivity * # records

    example:
    clustered index = 2 + 0.5[L] + 0.5[B]
    unclustered = 2 + 0.5[L] + 0.5|B|

## Evaluating Query
- Either optimal I/O query on interesting queries
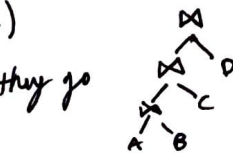- Interesting = sorted on a column used by GROUPBY or ORDERBY or used in a downstream join

## Pass 2...n
1. At each pass i, attempt to join i tables together, each from pass i-1 and pass 1
Advance optimal plan for each set and also the optimal plan for each interesting order for each set

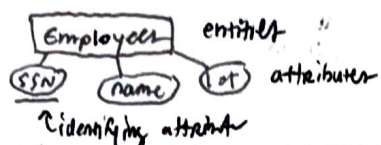\* Joined Tables must be on left
\* NO CROSS JOINS are considered
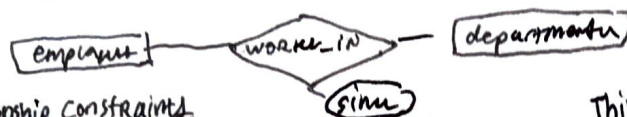
· when designing a database, we often use
Entity-Relational Models.

def: entity: a real world object described by a set of attribute values


Employees entities
SSN, name, lot — attributes
↳ identifying attribute

def: Relationship: association among two or more entities


employees — works_in — departments
time

**Relationship Constraints**

Thin black line — many-to-many Relationship

ex) many employees can work in many departments
many departments can have many employees

**Functional Dependencies and Normalization**

· X→Y means X column determines Y column in a table R.
i.e given any two tuples in table R, if their X values are
the same, then their Y values must be the same
· superkey: set of columns that determine all the columns in
the table
· candidate key: minimal set of columns that determine all the
columns in the table
· closure of F: F+: set of all FDs that are implied by F

Relation R w/ 4 FD F

**Decomposing a Relation — Boyce Codd Normal Form**
· Relation R with FDs F is in BNCF if for all X→A in F+,
A ⊆ X (called trivial form) or X is a superkey for R.

The procedure:
If X→Y violates BCNF, R becomes R-Y and XY

ex) Relation R = {C, S, J, D, P, Q, V} key C and F={JP→C, SD→P, J→S}
· to deal w/ SD→P, decompose into SDP, CSJDQV
· to deal w/ J→S, decompose CSJDQV into JS and CJDQV
end up w/ SDP, JS, and CJDQV

**Attribute Closures**
A+    Suppose if you have F= A→B, AB→AC, BC→BD, DA→C
then A+    closure = {AB, ABC, ABCD,

**key constraints** · each department has at most one manager but
employee can be manager for 1 or more dept

Thin Arrow — 1 to many — at most 1, or 0, and 0 or more
Thick Line — participation — at least one!
Thick Arrow · key + participation: at least and at most one
Department ⇒ Manager; each dep, one manager
exactly one

weak entity — entity that can be identified uniquely
only w/ key of another entity (owner entity)

primary key

X+: set of attributes closures, set of attributes
implied by X

**Decomposition Preservation**

IFF a dependency can be applied, i.e we have a set containing all the variables of a dependency