

SQL

No aggregates! Relational Table Properties

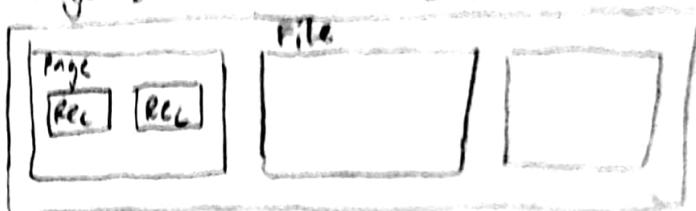
- (5) SELECT [Distinct] col1, col2, agg(col1)
 - (1) FROM table1, table2,
 - (2) WHERE <predicate> AND <predicate> OR
 - (3) GROUP BY <column list> ← Everything in SQL must be in GROUP-BY
 - (4) HAVING <predicate>
 - (6) ORDER BY <columns> [DESC] [ASC]
 - (7) LIMIT <integer>;
- Must contain only GROUP-BY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

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

Disk Representation

(page 30, location on page)

- DB File: collection of pages
- Page: collection of records



SQL String Comparison

- Old School SQL starts with %

WHERE S.name LIKE 'B_%o' ← returns Bob

- = any single char ; % = 0+ chars

- Standard Regular Expressions

WHERE S.name ~ 'B.*o' ← 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

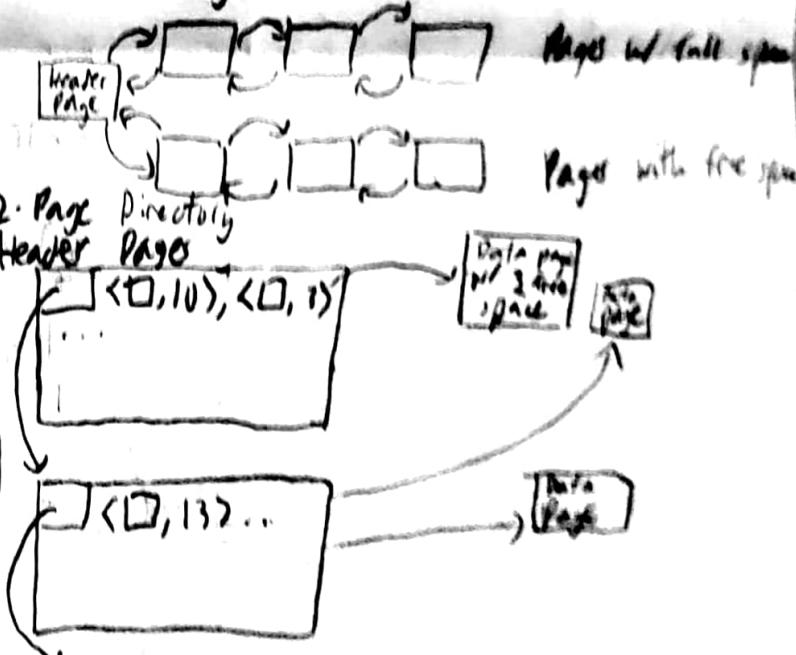
- Right: if t2.id has no match, t1.id is NULL
t1 () t2

- FULL: LEFT and RIGHT both apply to t1 () t2

Unordered Heap Files

Records placed arbitrarily across pages

- 1. File as Doubly Linked List



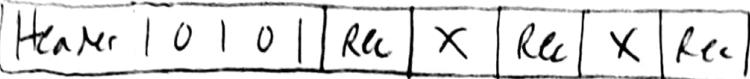
Unordered Heap File vs. Sorted File

2, 5	1, 6	4, 7	3, 6	7, 9	1, 2	3, 4	5, 6	7, 8	9, 10
------	------	------	------	------	------	------	------	------	-------

Scan all records	Index	Search	Average
Linear search	O(1)	O(log B) · D	Applies page to each of heap files
Range search	O(B · D)	O(log(B · page)) · D	Sorted files are packed
Insert	O(D)	O(log(B + B) · D)	
Delete	O(D · B)	O(log(B + B) · D)	

Page Layout

- FIXED record lengths, UNPACKED records

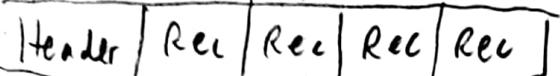


Insert: find first empty slot in bitmap X

Delete: find record and clear bit in header

- No organization needed

- FIXED record lengths, PACKED RECORDS



Insert: just append to end

Delete: Scan for record, delete, reorganize all records (w/c pack)

- VARIABLE record lengths, UNPACKED records (slotted page)

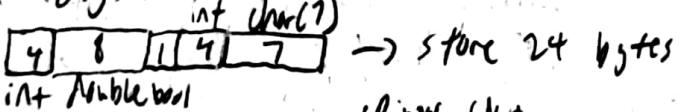


Delete: remove footer pointer (set record slot to null)

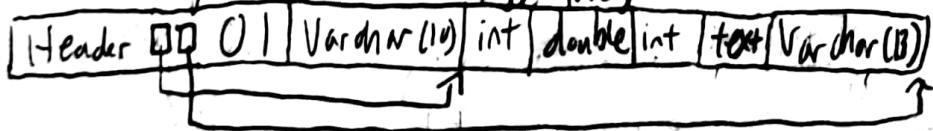
Insert: First empty pointer slot of free space pointer

NOTE: footer stores one int and one pointer for EACH PAGE

• Fixed length:



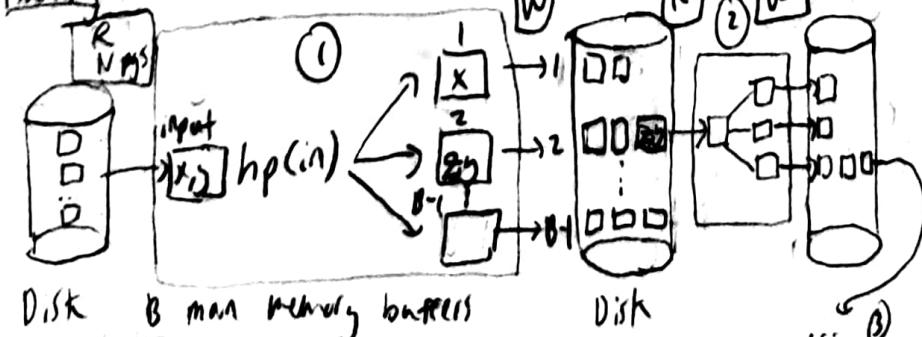
• Variable length:



- 2 pointers to variable length attributes

- bitmap (2 bits) where int and text are NULL

Hashing



Index Files (B+ trees) - acts similar to B-trees

• Property: # entries per node, $d = \text{order} \leq \text{entries}$

• Alternative 1: record contents stored in the leaf node itself

• Alternative 2: Leaf nodes: $\langle K, \text{rid of matching record} \rangle$

• Alternative 3: Leaf nodes: $\langle K, \text{list of rids of records} \rangle$

• Property 1: leaf node entries: $d = \text{order entries} \leq 2d$

• Property 2: all leaves same dist from root

• Property 3: inner node w/ K children have $K+1$ children

• Leaf split insert(f^*)



• Bulk Load

- Input sorted record keys (1*, 2*, 3*, ...)
not kept if inner node

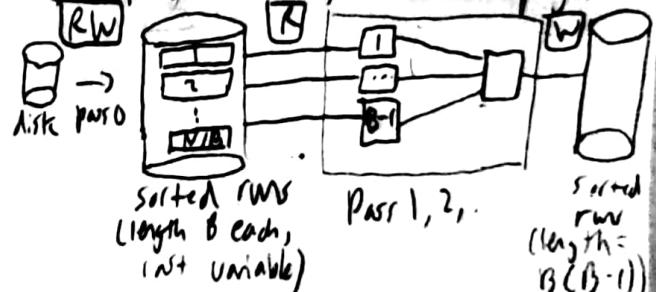
- Fill leaf page to fill factor $> d$

- Update inner pages to full

- Follow leaf and inner page split rules

• General External Merge Sort

- N pages to sort, B buffer pages



- Number of passes: $1 + \lceil \log_{B-1} \rceil N/B$

- Total I/Os = $(\text{I/Os per pass}) * (\# \text{ passes})$
 $= 2N * (1 + \lceil \log_{B-1} \rceil N/B)$

- 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 ($\leq B$ pages), read + write to check (Build Pass)

$\text{JOINS}[R] = \# \text{ pages in } R, p_r = \# \text{ records per page in } R, |R| = \# \text{ records}$

• Simple Nested Loop join on B:

for record r in R : { cost = sum R pages +
for record s in S : } scan S pages per record
if $\theta(r, s)$:
add $\langle r, s \rangle$ to result buffer

• Pages Nested Loop Join

for rpage in R : { cost = $|R| + [R][S]$
for spage in S : } cost = $[R] + [R][S]$
for rtuple in rpage
for stuple in spage.

if $\theta(r, s)$: add $\langle r, s \rangle$ result buffer

• Index Nested Loops Join

for each tuple r in R :
for each tuple s in $S, r == s$:
cost add $\langle r, s \rangle$ to res buffer
 $= |R| + |R|$ (cost to find matching tuples in S)

• Block Nested Loop Join

for rchunk of $B-2$ pages of R : { cost =
for spage of S : } $\{[R] + [R]/(B-2)][S]$
add matching tuples to res buffer

• Sort-Merge Join

(1) Sort R and S by join key

Mark = NULL

while r in range:

if not mark:

while $r < s$, advance r

while $r > s$, advance s

mark = s # start of block of s

if $r == s$:

result = $\langle r, s \rangle$

advance s

return result

reset s to mark

advance r

mark = NULL

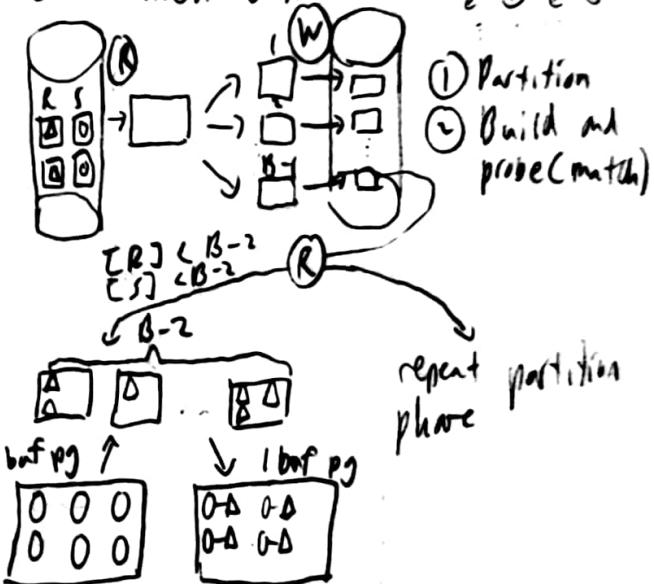
$r \rightarrow$

sid	shname
22	duran
28	allen
71	john
31	joe
44	mama
58	rusty

$s \rightarrow$

sid	bid
22	100001
28	101011
71	110110
31	10001
44	01001
58	00110

COST:
best: sort R + sort S + $|R| + |S|$
worst: sort R + sort S + $|R||S|$



Relational Algebra

Projection (Π) $\equiv \text{SELECT}$

Selection (σ_{cond}) $\equiv \text{WHERE}$

Renaming (ρ) $\text{Temp}_2(R_2.\text{rid} \rightarrow \text{sid}_2, S_1.\text{sid} \rightarrow \text{rid}_2)(R_1 \times S_1)$
output table old col \rightarrow new col old table
new

Union (\cup) \equiv tuples in R_1 OR R_2

Set-difference ($-$) \equiv tuples in R_1 but not in R_2

Cross-Product (\times) \equiv each R_1 row paired w/ each R_2 row

Intersection (\cap) $\equiv S_1 - (S_1 - S_2)$

Theta Join (Δ_θ) \equiv join on logical expression θ

Natural Join (Δ) \equiv equi-join on matching col names

Left Outer Join (Δ_L)

Right Outer Join (Δ_R)

Full Outer Join (Δ_F)

Group By ($\text{Yage}, \text{count}(*)) \gg 2$) \equiv group by age,
having count ≥ 2

Iterators

Materializing: save iterator as list on disk (cost I/O)

Streaming: don't need to save on disk (SELECT, π)

Pushdown: equivalent statements

• Selection (σ)

cascade: $\sigma_{c_1, c_2, \dots, c_n}(R) \equiv \sigma_{c_1}(\sigma_{c_2}(\dots(\sigma_{c_n}(R))\dots))$

commute: $\sigma_{c_1}(\sigma_{c_2}(R)) \equiv \sigma_{c_2}(\sigma_{c_1}(R))$

• Projections

cascade: $\Pi_{a_1}(R) \equiv \Pi_{a_1}(\Pi_{a_2}(\Pi_{a_3, a_4, \dots}(R)))$

• Cartesian Product

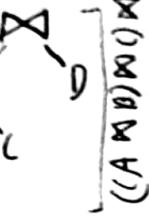
associative: $R \times (S \times T) \equiv (R \times S) \times T$

commutative: $R \times S \equiv S \times R$

Query Optimization

① Plan Space

- Try diff queries (equivalent relational algebra)
- Try different types of joins (brute hash, Block)
- Heuristics
 - Left-deep trees only
 - Avoid Cartesian products
 - Selection projection pushdown



② Cost Estimation (cost = #J/I + CPV factor * #tuples)

- Catalog: containing tuple statistics (updated periodically)
- Selectivity (sel) = |output| / |input| of σ term
 - σ term 1 \cap term 2 \cap term 3
 - reflects impact of σ term

$$\begin{array}{|c|c|c|} \hline & & \\ \hline \end{array} \xrightarrow{\sigma_{\text{sel}}=1} \boxed{1}, \text{ sel} = \frac{1}{5}$$

Predicate Sel Assume

$C = V$	$1/\text{distinct } C$	$ C $
$C = V$	$1/10$	
$C_1 = C_2$	$1/\max(\text{distinct } C_1, \text{distinct } C_2)$	$ C_1 , C_2 $
$C_1 = C_2$	$1/\text{distinct } C$	$ C_1 \text{ only}$
$C_1 = C_2$	$1/10$	
$C \leq V$	$(V - \min(C)) / (\max(C) - \min(C) + 1)$	$\max(C), \min(C)$
$C > V$	$(\max(C) - V) / (\max(C) - \min(C) + 1)$	$\max(C), \min(C)$
$C \leq V$	$1/10$	
$C > V$	$1/10$	
$C \leq V$	$(V - \min(C)) / (\max(C) - \min(C) + 1) + \frac{1}{10}$	$\max(C), \min(C)$
$C \geq V$	$(\max(C) - V) / (\max(C) - \min(C) + 1) + \frac{1}{10}$	$\max(C), \min(C)$
$C \leq V$	$1/10$	
$C \geq V$	$1/10$	
$C \geq V$	$(\max(C) - V) / (\max(C) - \min(C))$	$\max(C), \min(C)$
$C \geq V$	$1/10$	
$C \leq V$	$(V - \min(C)) / (\max(C) - \min(C))$	$\max(C), \min(C)$
$C \leq V$	$1/10$	

$$P_1 \wedge P_2 = S(P_1) * S(P_2) \quad \text{Independent term}$$

$$P_1 \vee P_2 = S(P_1) + S(P_2) - S(P_1) * S(P_2) \quad \text{Independent term}$$

$$\text{NOT } p \vdash S(p) \quad \text{Lemme: can we construct original relation?}$$

• Decompose R into A and B

Decomp is lossless iff F^+ contains:

$$X \cap Y \rightarrow X \quad \text{AND} \quad X \cap Y \rightarrow Y$$

③ Search Algorithm (best job)

pass 1: find minimum cost access method
(Index scan, full scan, etc.)
for each table, and interacting order
on interacting index

pass i:

- Consider only: left deep and not cartesian product (unless all are cartesian)
- advance only: cheapest cost plan every subset group of relations and interacting order

Interesting Orders

- ORDER BY attributes
- GROUP BY attributes
- downstream join attributes

Transactions and Concurrency

Transactions follow rules

- Atomicity: All operations in a transaction happens or no ops happen
- Consistency: Data starts and ends consistent
- Isolation: Execution of txns looks like running one at a time
- Durability: if txn commits, effects persist

Non-conflict swapping:

If two consecutive operations from two txns can be swapped if non-conflicting.

$T_1: R(A) \dots W(A)$ } ^{swaps} $\text{OR } T_2: W(A) \dots R(A)$

Conflicting

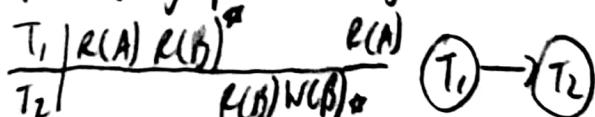
Conflict equivalent: schedules w/ same txn operations, different orderings (non-conflict swap)

Serializable: a schedule S is equivalent to a serial schedule

Conflict Serializable: schedule S is conflict serializable to some other serial schedule

Dependency graph: $(S \rightarrow \text{Serial})$

- One node per transaction
- If an operation in T_i conflicts w/ an operation in T_j and T_i comes first, add edge from T_i to T_j .
- Schedule is conflict serializable IFF dependency graph is acyclic



View Serializable: same as conflict serializable but calls $R(A)$, $W(A)$ non-conflicting, only $R(A)$, $W(A)$ is conflicting

- $(S \rightarrow V)$

Simple Locking

S	X	X
X	X	X

on a resource

- An S lock lets a txn read a resource
- many txns can hold S locks on a resource at once
- An X lock lets a txn modify a resource
- No other txn can have any type of lock while a txn has X lock
- If a txn can't get a lock it wants, it blocks and waits until another txn releases the conflicting lock

Deadlocks

$T_1: R(A), W(B) \rightarrow T_1$ waits for T_2 's lock on B, while T_2 waits for T_1 's lock on A

- Priority: older txns = higher priority
- wait-die: if T_i wants T_j 's lock
 - if T_i higher priority, wait for T_j to release
 - if T_i lower priority, T_i aborts
- wound-wait: if T_i wants T_j 's lock
 - if T_i higher priority, T_j aborts ("wound")
 - if T_i lower priority, T_j waits
- "waits-for" graph
 - One node for each txn
 - If T_j waits for T_i , edge from T_j to T_i
 - Cycle below deadlock

2-Phase Locking (2PL)

- A txn may not acquire a lock after releasing any lock held
- Conflict Serializable guaranteed



Strict 2-Phase Locking (Strict 2PL)

- Same as 2PL except all locks released at end of txn
- avoids cascading aborts



Multi-granularity Locking

NL	IS	IX	SIX	S	X
NL	✓	✓	✓	✓	✓
IS	✓	✓	✓	✓	X
IX	✓	✓	X	X	X
SIX	✓	X	X	X	X
S	✓	✓	X	X	X
X	✓	X	X	X	X

- To get S or IS lock on a node, need IS or IX on parent node
- To get X, IX, or SIX on a node, need IX or SIX on parent node

- Assume Strict 2PL unless interface or txn's new operation

Recovery

STEAL: txns can flush pages w/ uncommitted updates
 Pro: can maximize use of buffer pages
 Con: Must Undo if program crashes before commit

FORCE: force dirty pages to disk before commit
NO-STEAL: txns cannot flush pages w/ uncommited changes
NO-FORCE: dirty pages allowed to not be flushed on disk upon commit

ARIES (STEAL, NO FORCE)

Write ahead logging

- txns not committed until log records for changes written to disk
 - changes must be logged before data modified on disk

Undo Logging (Atomicity) (Dirty Pages)

$\langle \text{Start T} \rangle \langle \text{Commit T} \rangle \langle \text{Abort T} \rangle \langle T, X, V \rangle$

- If T commits, then $\text{FLUSH}(X)$ must be written to disk before $\langle \text{Commit T} \rangle$
- If T modifies X, then $\langle T, X, V \rangle$ written to log before $\text{FLUSH}(X)$

Redo Logging (Durability) (Updated Element)

$\langle \text{Start T} \rangle \langle \text{Commit T} \rangle \langle \text{Abort T} \rangle \langle T, X, V \rangle$

- If T modifies X, both $\langle T, X, V \rangle$ and $\langle \text{Commit T} \rangle$ must be logged before $\text{FLUSH}(X)$

Undo

$\text{flush}(X)$

$\text{Flush}(Y)$

$\langle \text{Commit}(txn) \rangle$

Redo

$\langle \text{Commit}(txn) \rangle$

$\text{flush}(Y)$

ATT

Txn	LastLSN	Status
T1	10	Running
T2	20	Committing
T3	40	Running
T4	50	Aborting

Page	reclsn
P1	0
P3	10

- Writes
- UPDATE log record for every write
- Commit:
- Write COMMIT log record, flush log to disk, "committed" is status, release txn locks, write END log record, remove from txn table
- Aabort (see UNDO)
- Write ABORT log record, write CLR record for each UPDATE the txn logged in reverse order, release txn locks, write END log record, remove from txn table
- Assume strict 2-phase locking and write-ahead logging
- Data structures
 - Transaction Table: information on currently active transactions
 - Dirty page table: dirty pages in memory (running, committing, aborting)

① Analysis: identify changes not yet written to memory at time of crash

• Start at last BEGIN-CHECKPOINT and scan all following log records

• For each record:

if record is UPDATE/CLR:

if txn not in ATT: add it (Txn ID, Status, LastLSN)

if txn in ATT: update LastLSN

if record is ABORT/COMMIT: update ATT (Status, LastLSN)

if record is END: remove transaction from ATT

② REDO: repeat all actions to restore database state to what it was at the time of the crash from the smallest recLSN

• For each UPDATE/CLR beginning from min(all recLSN in DPT):

REDO operation UNLESS any of the following is true:

- affected page not in DPT

- recLSN of page in DPT > LSN

- pageLSN of page (in DB NSK) > LSN

③ UNDO: undo the actions of transactions that did not commit

• For all "loser" (active) transactions:

undoLSN = LastLSN

+ while undoLSN != NULL:

if record at undoLSN is UPDATE:

undo = add record to CLR
 pageLSN = LSN
 of first modified page

ii) undoLSN = record's reclsn w/ undoNextLSN = record's pageLSN

Add END record for this txn to CLR