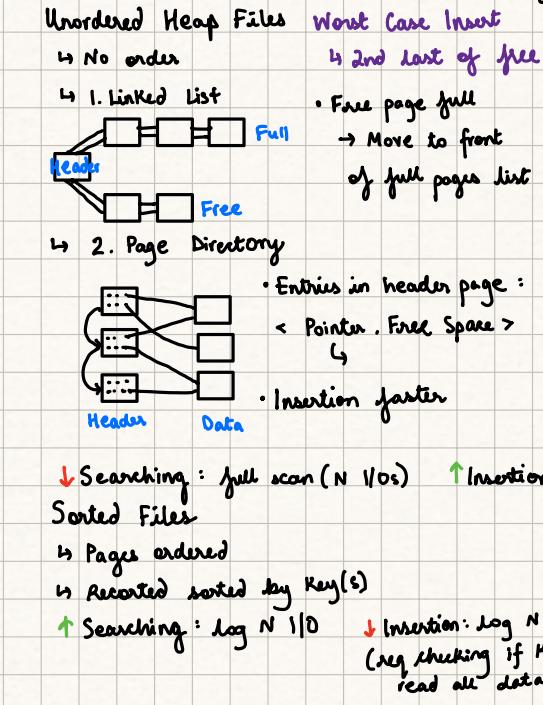


DISK, FILES

- Read: Disk $\xrightarrow{\text{Page}}$ RAM
- Write: RAM $\xrightarrow{\text{Page}}$ Disk Disk
- File: SSD (Flash) \downarrow Pages \downarrow Records
- \downarrow Cell \downarrow Sequential
- \downarrow Failure after erasure

FILES



SQL

- Tables make up relational databases
- NULL: falsy, any type
 - GROUPING
 - summarize cols [SUM, AVG, MAX, COUNT]
 - ignores NULL except COUNT(*)
 - check select statement w grouping?
 - where \rightarrow grouping \rightarrow having

SELECT [DISTINCT] <columns> ⑤

FROM <table> ①

- [INNER / LEFT / RIGHT / FULL] JOIN <table> FROM
ON <predicate>
- [WHERE <predicate> ②]
- [GROUP BY <columns> ③]
- [HAVING <predicate>] ④
- [ORDER BY <column list>] ⑥
- [LIMIT <# rows>]; ⑦
- LIKE $\sim \cdot : 0+ \text{chars}, - \text{any char}$

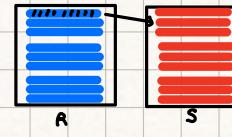
Subqueries:

WHERE

- WITH distinct in UNION- either common INTERSECT - in both ALL - all subquery meet condition ANY - any one

JOINS

SNLJ



Every record \rightarrow Every Record

for r_i in R:

for s_j in S:

if $\theta(r_i, s_j)$:

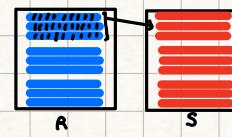
yield

$$[R][S] \rightarrow \# \text{ pages}$$

$$[R][S] \rightarrow \# \text{ records} / \text{page}$$

Do both ones, take the min one

PNLJ



for each P_R in R:

for each P_S in S:

for each row r in P_R :

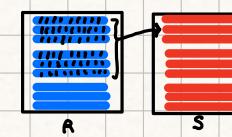
for each row s in P_S :

$$[R] + [R][S]$$

Keep smaller relation as outer loop

SNLJ, PNLJ use 3 buffers (2 input + 1 output)

BNLJ



for each block of $(B-2)$ pages C_R :

for each P_S in S:

for each row r in C_R :

for each row s in P_S :

$$[R] + \frac{[R]}{B-2} [S]$$

$(B-2)$ buffers for smaller one, 1 for other, 1 output

INLJ

for each record r_i in R:

for each record s_j in S where $\theta(s_j, r_i)$ is true:

<yield>

$$[R] + |R| \times (\text{cost to lookup in } S)$$

SORT-MERGE JOIN:

1. Sort both R and S

2. Merge matching tuples

while (True) {

while ($r < s$) { advance r }

while ($r > s$) { advance s }

mark

while ($r == s$) {

while ($r == s$) {

yield r, s

advance

3

reset s to mark

advance r

4

Don't need to materialize sorted relation

In final merge pass of sorting both, stream to SMJ

\rightarrow SAVES $2([R]+[S])$

Need:

runs in last merge + # runs in last merge of R of S

$$\leq (B-1)$$

Optimized: Total - $2[S] + [R]$

RELATIONAL ALGEBRA

Projection Π [select] takes columns

Selection σ [where] filters rows

Union \cup combine diff relations, remove duplicates] MUST have same attributes

Set Diff $-$ rows in table 1 not in table 2] same attributes

Group By λ [GROUP BY]

Intersection \cap rows in both tables

Cross Product \times 1 tuple for every possible pair

Alias ρ [AS] renames attributes, for join

Join \bowtie default natural join

\wedge (AND) \vee (OR)

can't have duplicates (tuples)

After Π , other cols DON'T exist

$$S_1 - (S_1 - S_2) = S_1 \cap S_2 = S_1 \times S_2$$

B+ TREES

- Order d , always balanced
- Must have sorted $d \leq x \leq 2d$ entries
- Inner Nodes have max $(2d+1)$ pointers
- Only leaf nodes have records
- Max values: $2df(2d+1)^h$

INSERTION

- Find leaf node, add $\langle K, R \rangle$ in order
- Overflow ($L > 2d$)
 - $L_1(d)$ $L_2(d+1)$
 - If L leaf, copy to parent
 - If L inner, PUSH to parent
 - Adjust pointers
 - Recurse

DELETION

- Just delete from leaf. NEVER from inner

STORAGE

- Alt 1 (By Val): Leaf contains data $\langle K, R \rangle$ can't support multiple indices, sorted
- Alt 2 (By Ref): Contains pointers $\langle K, \# \text{Page} \# \text{Record} \rangle$ \uparrow multiple indices
- Alt 3 (By List of Ref): List of pointers $\langle K, [\dots] \rangle$ \uparrow multiple records, same leaf node entry

CLUSTERING

Unclustered: $\sim 1/\text{O}$ per record

Clustered: $\sim 1/\text{O}$ per page

COUNTING I/Os

- Read root \rightarrow leaf
- Read data page (un vs clustered)
- Write data page
- Write index page

BULK LOADING

- Sort data
- Fill leaf to f
- Add pointer to parent
- Parent Overflow:

$L_1(d)$ $L_2(d+1)$

Move L_2 first entry to parent

- Adjust pointers

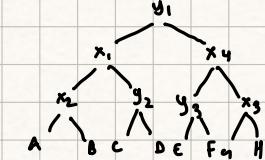
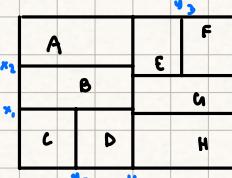
SELECT \rightarrow best case: doesn't exist, so only root \rightarrow leaf
 Check if pages overflow
 Can be built on any type of column
 full scan = # data pages
 Index Search: $h + l = \log_F(\# \text{leaves}) + 1$
 Bulk loading makes use of temporal locality

SPATIAL INDICES

Brute force search = # calc = # points

K-d tree search

K-d tree that organizes points in a K-dim space
 Balanced binary tree that splits on dimensions



EQUALITY SEARCH: $\log(BR) + 2$
 RANGE SEARCH: $\log(BR) + 1 + 2^m \# \text{pages}$

INSERTION: $\log BR + 4$ \rightarrow search, read leaf, write data, write index

DELETION: $\log BR + 4$

Searching:

- Use BST to find block it's in
- Backtrack if other neighbour block is closer

JOINS PT 2

GRACE HASH JOIN

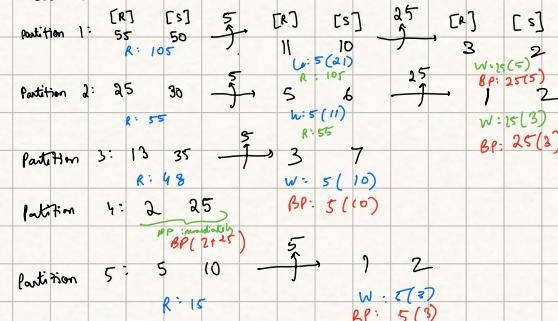
1. Partition Data

- Partition R,S into $(B-1)$ partitions
- Recurse until either $\leq (B-2)$

2. Build and Probe

- Build in memory hash table for smaller one

Pass 1: 500



COSTS

	Heap	Sorted	Clustered
Scan All	B	B	$\frac{1}{B} \cdot B$
Equality	$B/2$	$\log_2 B$	$\log_F(\frac{BR}{E}) + 2$
Range	B	$\log_2 B + \# \text{pages}$	$\log_F(\frac{BR}{E}) + 3 \cdot \# \text{pages}$
Insert	$2B$	$\log_2 B + B$	$\log_F(\frac{BR}{E}) + 4$
Delete	$\frac{B}{2} + 1$	$\log_2 B + B$	$\log_F(\frac{BR}{E}) + 4$

$B = \# \text{Data Blocks}$
 $A = \# \text{Records/Block}$
 $F = \text{Fanout}$
 $E = \# \text{Entries per leaf}$

BUFFER MANAGEMENT

4. Buffer frame holds 1 data page perfectly

5. Metadata Table: Frame ID, Page ID, Dirty Bit, Pin Count, memory addr

6. LRU: evict oldest set that is UNPINNED

7. CLOCK: LRU approx, uses ref bit \uparrow time

- skip pinned until unpinned + ref bit 0 found
- If ref=0, evict, write, read new page + set ref bit to 1
- Move clock hand to next frame

4. MRU: set ref bit to 1, DON'T move hand

5. MRU: evict most recently used \uparrow sequential scans

Requester sets dirty page

Page could have multiple pins

MRU > LRU sequential, = if $P \subset F$

SORTING

B buffer pages, $(B-1)$ input, 1 output

I/Os: $2N \cdot (1 + \log_{B-1} \frac{N}{B})$
 $\# \text{passes}$

$N=1000, B=11$

Pass 0: $1000/11 = 91$ runs, 11 each

Pass 1: $91/10 = 10$ runs, 110 each

Pass 2: $10/10 = 1$ run, 1100 each
 stop here

pages to sort in p passes: $B(B-1)^{p-1} \geq N$

I/Os does NOT depend on # B pages

HASHING

Hash into $(B-1)$ partitions, recurse until $\leq B$, 1 for streaming in

Group By + Duplicate Values

Hash func need to be distinct + 1 for conquer phase

Data skew \downarrow

Hash w/o recursive part: $B(B-1)$

B buffers, p_3 passes $\Rightarrow (B-1)^{p-1} B = \# \text{max pages}$

\uparrow part conquer

Stop when 1 run

QUERY OPTIMIZATION

Streaming operator: Work to produce each tuple

Blocking operator: Consume entire input THEN output generated

Selectivity Estimation

Pred	Selectivity
$C = V$	$\frac{1}{ C }, \frac{1}{ V }$
$C_1 = C_2$	$\frac{1}{\max(C_1, C_2)}$
$C \subseteq V$	$\frac{V - \min_c}{\max_c - \min_c} + \frac{1}{ C }$
$C > V$	$\frac{\max_c - V}{\max_c - \min_c} + \frac{1}{ C }$
$C \leq V$	$\frac{V - \min_c}{\max_c - \min_c}$
$C \geq V$	$\frac{\max_c - V}{\max_c - \min_c}$

Join Selectivity

- $|A| \cdot |B|$ (join set) $\min(A, B)$
- $|A| \cdot |B|$ (no. of joined tuples)
- tuple (no. of pages) records per page

Heuristic: Left deep, push down, NO cross join

System R:

Pass 1: Full scan/ Index scan

All 1: # h + sel x leaves

All 2/3 clustered: # h + sel x leaves + sel x data

All 2/3 unclustered: # h + sel x leaves + sel x rec

Pass 2 ... N: Join ($i-1$) + 1 tables

Advance lowest + interesting orders $\xrightarrow{\text{ORDER GROUP}} \xrightarrow{\text{JOINS}}$

↳ SMJ produces interesting order

↳ SNLJ, INLJ preserve left

TRANSACTIONS

- Inconsistent Read (WR): reads only a part of the update
- Lost (WW): multiple Xact write to same resource
- Dirty Read (WA): committed reads update never
- Unrepeatable Read (RW): different value

Conflict Serializability

If dependency graph is acyclic

Conflict: ≥ 1 write, dif Xact, same resource

$T_i \longrightarrow T_j$: T_i came earlier, conflict

Serial Equivalent:

Same Xact, order, state



View Equivalent: Same initial read, dependent read, winning writes

2PL: Ensures conflict serializable

- S before X
- No acquiring after releasing

Strict 2PL: Prevents cascading aborts

1. Release all locks together

Lock Management: Hash table { Resource : (Granted, Mode, Queue) }

Deadlock Avoidance (T_i waiting for T_j)

- Wait Die: T_i higher priority $\rightarrow T_i$ waits
 T_i lower priority $\rightarrow T_i$ aborts
- Wand Wait: T_i higher priority $\rightarrow T_j$ aborts
 T_i lower priority $\rightarrow T_i$ waits

LOCK GRANULARITY

Compatibility

	NL	IS	IX	S	SIX	X		Parent
NL	✓	✓	✓	✓	✓	✓		NL
IS	✓	✓	✓	✓	✓	✗		IS
IX	✓	✓	✗	✗	✗	✗		IX
S	✓	✓	✗	✗	✗	✗		S
SIX	✓	✓	✗	✗	✗	✗		SIX
X	✓	✗	✗	✗	✗	✗		X

RECOVERY

WAL: Log records to disk

- BEFORE page flush
- on commit: ACWAVS

Page flush if:

page LSN \leq flushed LSN

Xact commit if:

last LSN \leq flushed LSN

In DPT,

rec LSN \leq in mem LSN

rec LSN $>$ on disk LSN

Durability

- REDO: Reconstruct state

Start at smallest rec LSN

REDO All Update / CLR.

DON'T if

1. Page not in DPT

2. rec LSN $>$ LSN

3. page LSN \geq LSN

disk

DISTRIBUTED TRANSACTIONS

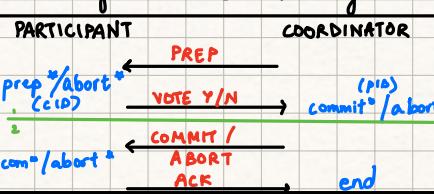
Distributed Locking

↳ Every node has own lock table (independent)

↳ UNION waits-for graph for distributed deadlock

2PC

↳ 2PC guarantees serializability



↳ Log + flush TMEN send

↳ Presumed Abort: DON'T flush abort

↳ Commit ONLY happens when no logs commit

2PC Recovery

- Participant Recovering

↳ No prepare: abort

↳ Prepare: ask coord

↳ Abort I: abort + vote NO

↳ Abort II: abort + send ack

↳ Commit: send ack

- Coordinator Recovering

↳ No commit: abort

↳ End: Nothing

↳ Abort: abort phase 2

PARALLEL QUERY PROCESSING

Architectures

Parallelism

1. Shared memory: CPUs share mem + disk

↳ Intra-query throughput

Each machine diff queries

Partitioning Schemes

One query across machines

↳ Sharding: 1 data page on 1 machine

↳ Replication: 1 data page on many machines

1. Range Partitioning ↗ Key lookup, range

Each machine gets range of values

2. Hash Partitioning ↗ Key lookup

Each machine gets hashed values

3. Round Robin ↗ parallelisation ↓ all req activation

Each machine gets same # pages

Network Cost

Data sent over network $\frac{pk(n-1)}{n}$

Parallel Sorting: Range partition + local sort

Broadcast Join: $(n-1) \cdot [p]$ smaller table.

Parallel Hashing: Hash partition + local hash

Parallel SMJ: Range partition BOTH + local SMJ

Passes: $1 + 1 + \# \text{passes}_R + \# \text{passes}_S + 1 + 1$

I/Os: $[R] + [S] + \left(1 + \log_{B-1} \frac{[R]}{B} \right) + \left(1 + \log_{B-1} \frac{[S]}{B} \right) + [R] + [S]$

Parallel GHJ: Hash partition BOTH + local GHJ

Hierarchical Aggregation: partitions should fit in memory

↳ COUNT $\rightarrow n: \text{count}$ c: sum

↳ AVG $\rightarrow c: \text{add} + \text{divide}$

No Steal Steal

REDO	UNDO
—	REDO

1. Steal / No-steal (Atomicity)

Can modified, uncommitted Xact be flushed?

↳ Steal: Yes (UNDO)

↳ No-steal: No (bad performance)

2. Force / No-force (Durability)

Are modified pages forced to disk on commit?

↳ Force: Yes

↳ No-Force: No (REDO)

Atomicity

3. UNDO:

Start at end, move up

UNDOes all updates for

Running/Aborting Xacts

↳ Write CLR UNDO Tx:LSN preLSN, undONEXTLSN

↳ Once all undone, write end record

↳ Remove Xact

Undo Logging: Steal Force

Dirty pages to disk B4 commit log

Flush pages B4 commit read

$\langle T, A, PREV \rangle$ if T com/ab → completed else write to disk

Redo Logging: NF/NS

commit B4 flush: redo if wrote after

commit but not flushed

$\langle T, A, nextVal \rangle$ read ↓ ignore uncom tx com: write new v

NOSQL

↳ Online Transaction Processing (OLTP):

- high no. of transactions ex. large no. of users
 - workload: frontend (social net, online stores)
 - queries: simple lookups, updates?
- ↳ Online Analytical Processing (OLAP):
- read only workloads
 - large amount of data
 - large no. of joins + aggregations

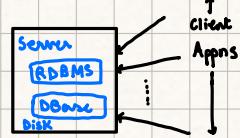
	OLTP	OLAP
main read	small no. of rec / query	aggregate over large no. of rec
main write	rand. access, low latency	bulk import or event stream
used	end user / customer via web app	internal analysis decision support
data	current state	history
size	GB - TB	TB - PB

Architectures:

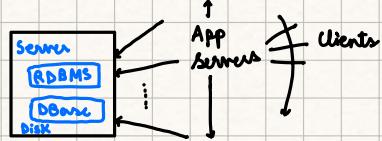
1. One-Tier Architecture



2. Two-Tier Architecture



3. Three-Tier Architecture



Partitioning / Sharding: Efficient for write-heavy workloads

Replication: Scale, resilient, extensive down time

CAP Theorem:

- ↳ Consistency: 2 clients simultaneous get same view
- ↳ Availability: Every request → Not a faulty response
- ↳ Partition Tolerance: Continue to operate despite drop/delay

Impossible for $> 2/3$ to be satisfied

↑ consistency \Rightarrow Error ↑ Availability \Rightarrow Scale