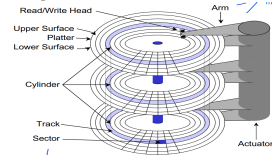


L1 - Data Storage

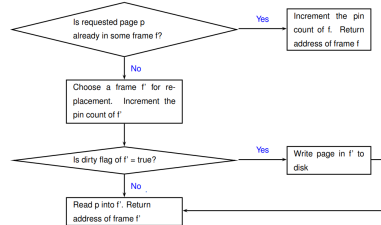
Magnetic Disks



- Disk Access Time** Seek time + Rotational Latency + Transfer time
- Response time** Queueing delay + Disk access time
- Rotational Delay** $\frac{1}{2} \frac{60s}{RPM}$
- Transfer Time** sectors on the same track * $\frac{TimePerRevolution}{SectorsPerTrack}$

Buffer Manager

- Buffer pool** Main memory allocated for DBMS
- pin count** is incremented upon pinning
- dirty bit** is updated when the page is unpinned (if modified)
- Replacement is only possible if pin count == 0

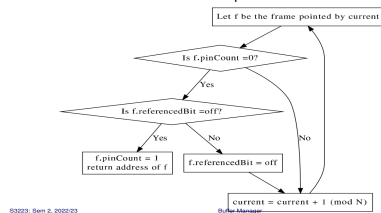


Replacement Policies LRU Policy

- Maintains a queue of pointers to frames with pin count = 0

Clock Replacement Policy

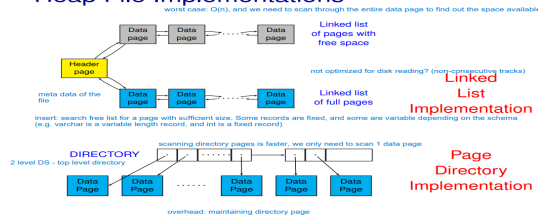
N = number of frames in buffer pool



- Simplifies LRU with a second chance round robin system
- Each frame has a **reference bit** that is turned on when pin count reaches 0
- Replaces a page when referenced bit is off and pin count is 0

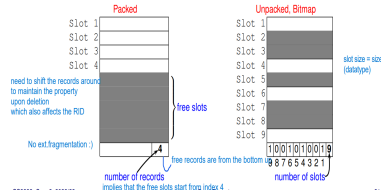
File Organisation

Heap File Implementations



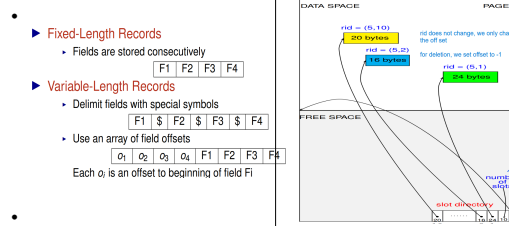
Page Formats: Fixed Length Records

- Packed Organisation** Store records in contiguous slots
- Unpacked Organisation** Uses a bit array to maintain free slots



Page Formats: Slotted Page (variable length record)

- Store records in slots of (record offset, record length)
- Record Offset: Offset of the record from the start of the page



L2 And L3 - Indexing

- A search key is a sequence of k attributes. If $k \geq 1$, composite key
- A search key is an unique index if it is a candidate key
- An index is stored as a file

Format of data entries

- Format 1: k^* is an actual data record with search value k
- Format 2: k^* is the form (k, rid)
- Format 3: k^* is the form (k, rid-list*)
- Note: Different formats affects the number of data entries stored in a page

Clustered Vs Unclustered

- Clustered:** Order of data entries is the same as the order of data records. Can only be built on ordered field (e.g. primary key)
- Unclustered:** Order of data entries does not correspond to the order of data records
- The implication is that we can read an entire clustered page with 1 I/O
- B+ Tree: Format 1 is clustered, Format 2 and 3 can be clustered if data records are sorted on the search key
- Hash: Only format 1 is clustered since hashing do not store data entries in search key order

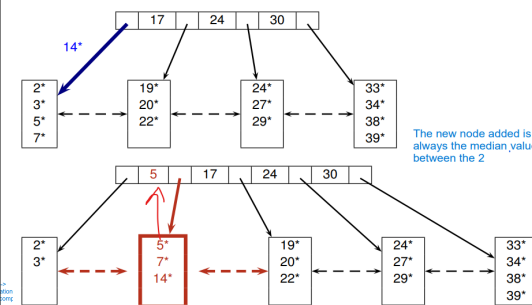
Tree Based Index - B+ Tree

- Leaf nodes are doubly linked and store Data Entries

- Internal nodes store index entries (p0, k, p1 ... pk, k, pk+1)
- Internal nodes contains m entries, $m \in [d, 2d]$ → space utilisation $\geq 50\%$
- Root contains m entries, $m \in [1, 2d]$

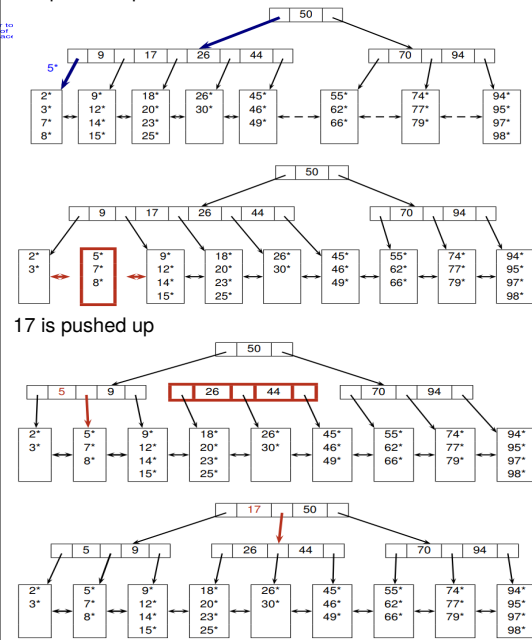
B+ Tree - Split Overflow Nodes

- Distribute d+1 entries to the new leaf node
- Create new entry index using smallest key in the new node (middle key)
- Insert new entry into parent node of overflowed node



B+ Tree - Overflow Propagation

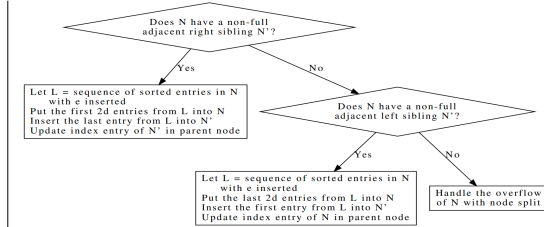
5 is pushed up



Excess middle node is pushed updated to parent node

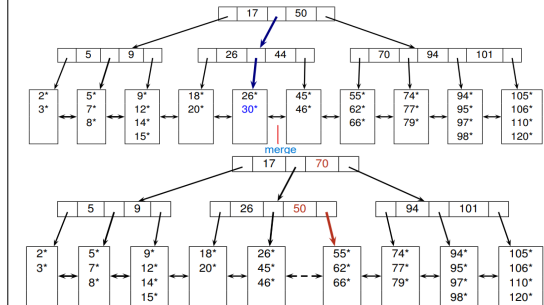
B+ Tree - Redistribution of data entries

Two nodes are siblings if they have the same parent node



B+ Tree - Underflow

- Underflow occurs when a node has less than d entries
- Underflow is resolved by redistributing entries between siblings
- An underflow node is merged if each of its adjacent siblings have exactly d entries



B+ Tree - Bulk Loading

- Initializing a B+ tree by insertion is expensive (need to traverse tree n times)
- 1. Sort all data entries by search key
- 2. Initialise B+ tree with an empty root page
- 3. Load data entries into leaf pages
- 4. In asc order, insert the index entry of each leaf page into the rightmost parent node

Hash Based Index

- Does not support range search, only equality queries

Static Hashing

- N buckets, each bucket has 1 primary page and ≥ 0 overflow pages
- To maintain performance, we need to routinely construct bigger hash tables and redistribute data entries

Dynamic Hashing - Extendible Hashing

- No overflow pages! A bucket can be thought of as a page
- At most 2 Disk I/Os for equality search (at most 1 if directory and bucket fits in memory)
- Instead of maintaining data entries, we maintain pointers to data entries in buckets
- Instead of maintaining buckets, maintain a directory of pointers to buckets
- The directory has 2^d buckets, where d is the global depth — large overhead if hashing is uniform
- Each director entry differs by a unique d-bit address
- Two directories are corresponding iff their addresses differ only in the dth bit
- All entries with the same local depth (l) have the same last l bits in h(k)

- Partition overflow:** Hash table $\pi_L^*(R_i)$ is larger than available memory buffers.
- Solution:** Recursively apply hash-based partitioning to overflowed partitions.
- Analysis:** Effective (no overflow) when B
- $$> \frac{|\pi_L^*(R)|}{B-1} * f \approx \sqrt{f * \pi_L^*(R)}$$
- If no partition overflow: $|R| + \pi_L^*(R)| + (\text{duplicate elimination})|\pi_L^*(R)|$
- Index based projection:** Do index scan if the wanted attributes \subseteq search key
- Join** $R \bowtie_{\theta} S$, where R is the outer relation and S is the inner relation

- Tuple-based**
 - Cost: $|R| + ||R|| * |S|$
 - for each tuple r in R
 - for each tuple s in S
 - if (r matches s) then output (r, s) 4 to result
- Page-based**
 - Load P_R and P_S to main memory
 - Cost: $|R| + |R| * |S|$
 - for each page P_R in R
 - for each page P_S in S
 - for each tuple $r \in P_R$
 - for each tuple $s \in P_S$
 - if (r matches s) then output (r, s) 4 to result
- Block nested-loop**
 - Allocate 1 page for S , 1 for output, B-2 for R
 - $|R| \leq |S|$
 - Cost: $|R| + (\lceil \frac{|R|}{B-2} \rceil * |S|)$
 - $|R| \leq |S|$
 - while Scanning R
 - read next (B-2) pages of R to buffer
 - for P_S in S
 - read P_S into Buffer
 - for $r \in buffer \wedge s \in P_S$
 - if (r matches s) then output (r, s) 4 to result
 - Without materialisation: $\lceil \frac{|R|}{B-2} \rceil * |T|$
- Index Nested Loop Join**
 - There is an index on the join attributes of S
 - Uniform distribution: r joins $\lceil \frac{|S|}{||\pi_{B_j}(S)||} \rceil$ tuples in S
 - format 1
 - B+Tree: $|R| + ||R|| * J$
- $J = \log_F(\lceil \frac{|S|}{b_d} \rceil)$ (tree traversal) + $\lceil \frac{|S|}{b_d ||\pi_{B_j}(S)||} \rceil$ (search leaf nodes)
- for $r \in R$
- use r to probe S 's index to find matching tuples
- Sort-Merge Join**
 - Sort R and S on join attributes and merge
 - Cost: $2|\pi_L^*(R)|(\log_{B-1}(\sqrt{\frac{|R|}{B}}) + 1) + 2|\pi_L^*(S)|(\log_{B-1}(\sqrt{\frac{|S|}{B}}) + 1) + |\pi_L^*(R)| + |\pi_L^*(S)|$
 - merging cost is $|R| + ||R|| * |S|$ if each tuple of R requires a full scan of S
 - Optimisation: $B \geq N(R, i) + N(S, i) + 1 \geq \sqrt{|R| + |S|}$
 - We can choose which relation to partition again if this is not met
 - Cost: cost of getting R , S + $k(\text{write out } (|R| + |S|)) + m(\text{merge } (|R| + |S|))$
 - If sorted on join column: $|R| + |S|$
- Grace Hash Join**
 - Partition into $B - 1$ partitions
 - If no partition overflow ($B > \sqrt{f * |S|}$):
 - $k(\text{Cost to partition } R, S) + \text{Cost of probe phase}$
 - Partition cost = cost of getting R + cost to write partitions ($|R|$)
 - Probe cost = $|R| + |S|$

L6: Query Optimisation

- Search space: Queries considered**
- Search Place** Queries being considered
- Linear** if at least one operand for each join is a base relation, bushy otherwise
- Left-deep** if every right join operand is a base relation
- Right-deep** if every left join operand is a base relation

2. Plan enumeration - for joins between 2 tables

Input: A SPJ query q on relations R_1, R_2, \dots, R_n
Output: An optimal query plan for q

```

01.  for i = 1 to n do
02.      optPlan( $\{R_i\}$ ) = best access plan for  $R_i$ 
03.  for i = 2 to n do
04.      for each  $S \subseteq \{R_1, \dots, R_n\}, |S| = i$  do
05.          bestPlan = dummy plan with cost(bestPlan) =  $\infty$ 
06.          for each  $S_j, S_k, |S_j| \in [1, i], S = S_j \cup S_k$  do
07.               $p$  = best way to join optPlan( $S_j$ ) and optPlan( $S_k$ )
08.              if ( $cost(p) \leq cost(bestplan)$ ) then
09.                  bestPlan =  $p$ 
10.          optPlan( $S$ ) = bestPlan
11.  return optPlan( $\{R_1, \dots, R_n\}$ )

```

- Decide on the plan for the 2 operands

- Decide on the plan to join: Block nested loop, sort merge join, grace hash join
- 3. Cost Model**
 - Uniformity:** uniform distribution of all values
 - Independence:** Independent distribution of values in different attributes
 - Inclusion:** for $R \bowtie S, \text{if } ||\pi_A(R)|| \leq ||\pi_B(S)||$ then $\pi_A(R) \subseteq \pi_B(S)$

Plans-no join, 1 table

- Table scan** Scan the entire table. Cost: $|R|$
- Index scan** Scan the index. Cost: $2 + |\text{leaf pages satisfying the predicate}| + |\text{entries satisfying predicate}|$ (unclustered)
- Index intersection with $I_a I_b$** Cost to partition predicate1(R) + Cost to partition predicate2(R) + cost to intersect partitions 1,2 + cost to RID lookup
- cost to partition: Scan index for matching pages + cost to write partitions from matching entries

Histogram

- Equiwidth** Each bucket has equal number of values
- Estimate: $\frac{1}{|bucket|} * ||bucket||$
- Equidepth** Each bucket has equal number of tuples
- Sub-ranges can overlap, tuples of the same value can be in 2 adjacent buckets
- $\frac{1}{|bucket_A|} * ||bucket_A|| + \frac{1}{|bucket_B|} * ||bucket_B|| + \dots$
- MCV** Separately track the top-k MCV and exclude them from the bucket

Size of query

- Join** $||R|| * ||S|| * \frac{1}{max(||\pi_b(R)||, ||\pi_b(S)||)}$
- Select - OR** $(1 - (p(a = x) * p(b = y))) * ||R||$
- Select - AND** $p(a = x) * p(b = y) * ||R||$

L7: Transaction Management

View Equivalent

- If T_i reads A from T_j in S , then T_i must also read A from T_j in S'
- For each data object A , Xact (if any) that performs final write on A in S must also perform final write on A in S'

Conflicting actions - WW, WR

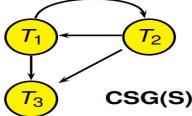
- Dirty Read** T_2 read uncommitted write from T_1
- Unrepeatable Read** T_2 updates an object that T_1 has previously read and T_2 commits while T_1 is still in progress $\rightarrow T_1$ can get a different value from read
- Lost Update** T_2 overwrites the value of an object that has been modified by T_1 while T_1 is still in progress

- View serializable prevents These
- Blind write** $R_1(X), W_2(Y), W_1(X)$ Blind write: $W_2(Y)$
- Conflict Serializable** Conflict equivalent to serial schedule, view serializable and not blind write
- Non Conflict Serializable** find conflicting action pairs($R_1(x) W_1(x)$), ($R_2(x) W_1(x)$)

Conflict Serializability Graph

- V contains a node for each committed Xact in S
- E contains (T_i, T_j) if an action precedes and conflicts with one of T_j 's actions

$R_1(A), W_2(A), Commit_2, W_1(A), Commit_1, W_3(A), Commit_3$



Schedules

- Cascading aborts** T_i read from $T_j \rightarrow T_j$ aborts $\rightarrow T_i$ aborts
- Recoverable** $\forall T \in S$ T_2 must commit after T_1 if T_2 reads from T_1
- Cascadeless** Whenever T_i reads from T_j in S , Commit must precede this action
- Theorem 4: Cascadeless \rightarrow Recoverable (not iff)
- Strict** to use before-images, $\forall W_i(O) \in S$, O is not read or written by another Xact until T_i either aborts or commits
- Theorem 5: Strict \rightarrow Cascadless (not iff)

2PL

- To read an object O , a Xact must hold a S-lock or X-lock on O
- To write to an object O , a Xact must hold a X-lock on O
- Once a Xact releases a lock, the Xact can't request any more locks
- Theorem 1: 2PL is conflict serializable

Strict 2PL

- To read an object O , a Xact must hold a S-lock or X-lock on O
- To write to an object O , a Xact must hold a X-lock on O
- A Xact must hold on to locks until Xact commits or aborts
- Theorem 2: Strict 2PL is strict and conflict serializable
- Strict 2PL prevents cascading rollback and deadlock and ensures recoverability

Detect deadlocks

- Waits-for graph (WFG) \rightarrow Deadlock is detected if WFG has a cycle
- Breaks a deadlock by aborting a Xact in cycle

Deadlock Prevention

- Each Xact is assigned a timestamp when it starts
- Assume older (smaller time stamp) Xacts have higher priority than younger Xacts
 - Suppose T_i requests for a lock that conflicts with a lock held by T_j
 - Two possible deadlock prevention policies:
 - Wait-die policy:** lower-priority Xacts never wait for higher-priority Xacts
 - Wound-wait policy:** higher-priority Xacts never wait for lower-priority Xacts

Prevention Policy	T_i has higher priority	T_i has lower priority
Wait-die	T_i waits for T_j	T_i aborts
Wound-wait	T_j aborts	T_i waits for T_j

L8: MVCC

Multi Version Serializable Schedle (MVSS)

- multiversion view equivalent** if S and S' have the same set of read-from relationships
- i.e. $R_i(x_j)$ occurs in S iff $R_i(x_j)$ occurs in S'
- Monoversion Schedule** each read action returns the most recently created object version
- MVSS** if there exists a serial Monoversion schedule that is multiversion view equivalent to S
- Note that a MVSS is not necessarily conflict serializable schedule if it is not a valid monoversion schedule
- E.g. $W_1(x_1), R_2(x_0), R_2(y_0), W_1(y_1), C_1, C_2$ is MVSS with (T_2, T_1) but contains conflicting actions $W_1(x_1)$ and $R_2(x_0)$

Snapshot Isolation (SI)

- Each Xact has a snapshot of the database at the start of the Xact and sees only versions from that snapshot and **its own writes**

- FUW** T needs to acquire X-lock on O (if not - wait), and if O has been updated by a concurrent T' then T aborts
- FCW** (no locks) before committing T checks if O has been updated, abort if it has been updated
- Write-skew anomaly**, not MVSS: $R_1(X_0), R_2(X_0), R_1(Y_1), R_2(Y_2), W_1(X_1), Commit_1, W_2(Y_2), Commit_2$
- Read-only anomaly**,not MVSS: $R_1(b), R_2(a), W_1(b), C_1, R_2(b), W_2(a), R_3(a), R_3(b), C_3$

Transaction Dependencies

- WW** from T_1 to T_2 : T_1 commits some version of X and T_2 writes the immediate successor
- WR** from T_1 to T_2 : T_1 commits some version of X which is read by T_2
- RW** from T_1 to T_2 : T_1 reads some version of X and T_2 commits the immediate successor
- DSG** V = xacts, E = Dependencies, use \rightarrow for concurrent transactions and \rightarrow for non-concurrent

L10-Recovery

Policies

- Steal:** Allows dirty pages to be written to disk before commit
- Force:** Requires all dirty pages to be written to disk when commit
- No-steal: no undo, Force: no redo. Pgsql uses steal and no-force

Restart: analysis, redo, undo

- Analysis:** identifies dirtied pool pages and active Xacts at time of crush
- Redo:** redo actions to restore db to pre-crush
- Undo:** undo actions of Xacts that did not commit

Analysis: Xact table

- When the first log record is created, create a new entry T with status U
- Update lastLSN for T to be r 's LSN
- Remove T if end log is seen

Analysis: Dirty Page Table

- New dirtied page will be added to the DPT with recLSN= r .LSN
- Remove entry when it is flushed to disk

	prevLSN	XactID	type	pageID	length	offset	before image	after image
10	-	T_1	update	P500	3	21	ABC	DEF
20	20	T_2	update	P600	3	41	HJ	KLM
30	10	T_2	update	P500	3	20	GDE	QRS
40	10	T_1	update	P505	3	21	TUV	WXY

DIRTY PAGE TABLE			XACT TABLE		
pageID	recLSN		XactID	lastLSN	status
P500	10		T_1	40	U
P600	20		T_2	30	U
P505	40				

Redo Phase

- Redo LSN = min(LSNs), then fetch page LSN,
- if $r.LSN > pageLSN$ and Page is in DPT, redo
- Update pageLSN to $r.LSN$

Undo Phase

- Start from largest LSN from L
- if update, create CLR with undoNextLSN= r 's prevLSN, update-L-TT(r .prevLSN)
- if CLR, update-L-TT(r .undoNextLSN)
- update-L-TT(lsn): add lsn to L if its not null, else add end log record for T and remove it from TT