

01. DBMS STORAGE

Data entry formats

1. k^* is an actual **data record** (with search key k)
2. k^* is of the form **(k , RID)** - fixed length (k , \bullet)
3. k^* is of the form **(k , RID-list)** - e.g. (k , {RID11, RID12})

02. TREE-BASED INDEXING

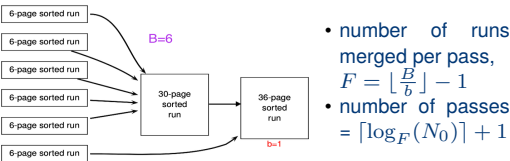
04.1 SORTING

External Merge Sort

- **sorted run** \rightarrow sorted data records written to a file on disk
- divide and conquer
 1. create temporary file R_i for each B pages of R sorted
 2. merge: use $B - 1$ pages for input, 1 page for output
- total I/O = $2N(\lceil \log_{B-1}(N_0) \rceil + 1)$
 - $2N$ to create $\lceil N/B \rceil$ sorted runs of B pages each
 - merging sorted runs: $2N \times \lceil \log_{B-1} N_0 \rceil$

optimisation with blocked I/O

- sequential I/O - read/write in *buffer blocks* of b pages
- one block (b pages) for output, remaining blocks for input



- number of runs merged per pass, $F = \lfloor \frac{B}{b} \rfloor - 1$
- number of passes = $\lceil \log_F(N_0) \rceil + 1$

Sorting with B^+ -trees

- when *sort key is a prefix of the index key* of the B^+ -tree
- sequentially scan leaf pages of B^+ -tree
 - for Format-2/3, use RID to retrieve data records

04.2 SELECTION: $\sigma_p(R)$

- $\sigma_p(R)$: selects rows from relation R satisfying predicate p
- **access path**: a way of accessing data records/entries
 - **table scan** \rightarrow scan all data pages
 - **index scan** \rightarrow scan index pages
 - **index intersection** \rightarrow combine results from index scans
- **selectivity** of an access path \rightarrow number of index & data pages retrieved to access data records/entries
 - more selective = fewer pages retrieved
- index I is a **covering index** for query $Q \rightarrow$ if all attributes referenced in Q are part of the key of I
 - Q can be evaluated using I without any RID lookup (**index-only plan**)

Matching Predicates

- **term** \rightarrow of form $R.A \text{ op } c$ or $R.A_i \text{ op } R.A_j$
- **conjunct** \rightarrow one or more terms connected by \vee
 - **disjunctive** conjunct \rightarrow contains \vee

- conjunctive normal form, **CNF predicate** \rightarrow comprises one or more conjuncts connected by \wedge

$$\underbrace{(\text{rating} \geq 8)}_{\text{term/conjunct}} \vee \underbrace{(\text{director} = \text{"Coen"})}_{\text{term/conjunct}} \wedge \underbrace{(\text{year} > 2003)}_{\text{term/conjunct}} \wedge \underbrace{(\text{language} = \text{"English"})}_{\text{term/conjunct}}$$

B^+ -tree matching predicates

- for index $I = (K_1, K_2, \dots, K_n)$ and non-disjunctive CNF predicate p , I matches p if p is of the form $(K_1 = c_1) \wedge \dots \wedge (K_{i-1} = c_{i-1}) \wedge (K_i \text{ op}_i c_i)$, $i \in [1, n]$
 - zero or more equality predicates
 - *at most one* non-equality comparison operator which must be on the last attribute of the prefix (K_i)
- matching index: matching records are in contiguous pages
 - non-matching index: not contiguous \Rightarrow less efficient

Hash index matching predicates

- for hash index $I = (K_1, K_2, \dots, K_n)$ and non-disjunctive CNF predicate p , I matches p if p is of form

$$(K_1 = c_1) \wedge (K_2 = c_2) \wedge \dots \wedge (K_n = c_n)$$

Primary/Covered Conjuncts

- **primary conjuncts** \rightarrow subset of conjuncts that I matches
 - e.g. $p = (\text{age} \geq 18) \wedge (\text{age} \leq 20) \wedge (\text{weight} = 65)$ for $I = (\text{age}, \text{weight}, \text{height})$
- **covered conjuncts** \rightarrow subset of conjuncts covered by I
 - each attribute in covered conjuncts appears in key of I
- primary conjuncts \subseteq covered conjuncts

Cost of Evaluation

let p' = primary conjuncts of p , p_c = covered conjuncts of p

B^+ -tree index evaluation of p

1. navigate internal nodes to find first leaf page
- $$\text{cost}_{\text{internal}} = \begin{cases} \lceil \log_F(\lceil \frac{\lfloor |R| \rfloor}{b_d} \rceil) \rceil & \text{if } I \text{ is a format-1 index} \\ \lceil \log_F(\lceil \frac{\lfloor |R| \rfloor}{b_i} \rceil) \rceil & \text{otherwise} \end{cases}$$
2. scan leaf pages to access all qualifying data entries
- $$\text{cost}_{\text{leaf}} = \begin{cases} \lceil \frac{\|\sigma_{p'}(R)\|}{b_d} \rceil & \text{if } I \text{ is a format-1 index} \\ \lceil \frac{\|\sigma_{p'}(R)\|}{b_i} \rceil & \text{otherwise} \end{cases}$$
3. retrieve qualified data records via RID lookups
- $$\text{cost}_{\text{RID}} = \begin{cases} 0 & \text{if } I \text{ is a covering format-1 index,} \\ \|\sigma_{p_c}(R)\| & \text{otherwise} \end{cases}$$
- reduce cost with **clustered** data records (sort RIDs): $\lceil \frac{\|\sigma_{p_c}(R)\|}{b_d} \rceil \leq \text{cost}_{\text{RID}} \leq \min\{\|\sigma_{p_c}(R)\|, |R|\}$

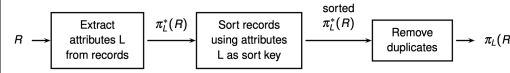
hash index evaluation of p

- **format-1**: cost to retrieve data records $\geq \lceil \frac{\|\sigma_{p'}(R)\|}{b_d} \rceil$
 - **format-2**: cost to retrieve data entries $\geq \lceil \frac{\|\sigma_{p'}(R)\|}{b_i} \rceil$
- cost to retrieve data records =
- $$\begin{cases} 0 & \text{if } I \text{ is a covering index,} \\ \|\sigma_{p'}(R)\| & \text{otherwise} \end{cases}$$

05.1 PROJECTION $\pi_{A_1, \dots, A_m}(R)$

- $\pi_L(R)$ eliminates duplicates, $\pi_L^*(R)$ preserves duplicates

Sort-based approach



cost analysis

1. extract attributes: $|R|$ scan + $|\pi_L^*(R)|$ output temp result
2. sort records: $2|\pi_L^*(R)|(\log_m(N_0) + 1)$
3. remove duplicates: $|\pi_L^*(R)|$ to scan records

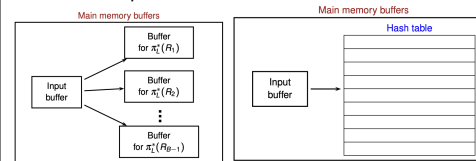
optimised sort-based approach



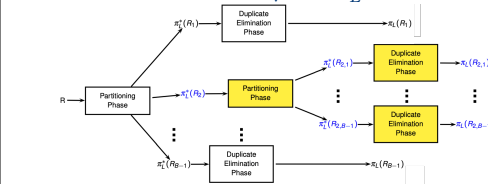
- if $B > \sqrt{|\pi_L^*(R)|}$, same I/O cost as hash-based approach
 - $N_0 = \lfloor \frac{|R|}{B} \rfloor \approx \sqrt{|\pi_L^*(R)|}$ initial sorted runs
 - $\log_{B-1}(N_0) \approx 1$ merge passes

Hash-based approach

1. **partitioning phase**: hash each tuple $t \in R$
 - $R = R_1 \cup R_2 \cup \dots \cup R_{B-1}$
 - for each R_i & R_j , $i \neq j$, $\pi_L^*(R_i) \cap \pi_L^*(R_j) = \emptyset$
 - for each t : project attributes to form t' , hash $h(t')$ to one output buffer, flush output buffer to disk when full
 - one buffer for input, $(B - 1)$ buffers for output
2. **duplicate elimination** from each $\pi_L^*(R_i)$
 - for each R_i : initialise in-mem hash table, hash each $t \in R_i$ to bucket B_j with $h' \neq h$, insert if $t \notin B_j$
 - write tuples in hash table to results



- **I/O cost** (no partition overflow): $|R| + 2|\pi_L^*(R)|$
 - partitioning cost: $|R| + |\pi_L^*(R)|$
 - duplicate elimination cost: $|\pi_L^*(R)|$
- partition overflow: recursively apply partitioning
 - to avoid, $B >$ size of hash table for $R_i = \frac{|\pi_L^*(R)|}{B_1} \times f$
 - approximately $B > \sqrt{f \times |\pi_L^*(R)|}$



Projection using Indexes

- if index search key contains all wanted attributes as a *prefix*
 - **index scan** data entries in order & eliminate duplicates

05.2 JOIN $R \bowtie_{\theta} S$

R = outer relation (smaller relation); S = inner relation

! for **format-2** index, add cost of retrieving record

nested loop joins

- **tuple-based** nested loop join: $|R| + |R| \times |S|$
- **page-based** nested loop join: $|R| + |R| \times |S|$
- **block nested loop join**: $|R| + (\lceil \frac{|R|}{B-2} \rceil \times |S|)$
 - $|R| \leq |S|$
 - 1 page output, 1 page input, $(B - 2)$ pages to read R
 - for each $(B - 2)$ pages of R : for each P_S of S : check r, s
- **index nested loop join**: $|R| \times \left(\log_F(\lceil \frac{\|S\|}{b_d} \rceil) + \lceil \frac{\|S\|}{b_d \lceil \pi_{B_j}(S) \rceil} \rceil \right)$
 - joining $R(A, B) \bowtie_{A=C} S(A, C)$ with B+tree index on $S.A$
 - for each tuple $r \in R$, use r to probe S 's index for match

sort-merge join

- sort R & S : $2|R|(\log_m(N_R) + 1) + 2|S|(\log_m(N_S) + 1)$
- merge cost: $|R| + |S|$ (worst case $|R| + |R| \times |S|$)
- **optimised sort-merge join**
 - merge sorted runs until $B > N(R, i) + N(S, j)$; then do merge and join at the same time
 - I/O cost: $3 \times (|R| + |S|)$
 - if $B > \sqrt{2|S|}$, one pass to merge initial sorted runs
 - $2(|R| + |S|)$ for initial sorted runs, $|R| + |S|$ for merging

hash join

1. partition R and S into k partitions on join column
 - $\pi_A(R_i) \cap \pi_B(S_j) = \emptyset \quad \forall R_i, S_j, i \neq j$
 - $R = R_1 \cup R_2 \cup \dots \cup R_k, \quad t \in R_i \iff h(t.A) = i$
 - $S = S_1 \cup S_2 \cup \dots \cup S_k, \quad t \in S_i \iff h(t.B) = i$
2. join corresponding partitions: $R \bowtie_{R.A=S.B} S = (R_1 \bowtie S_1) \cup \dots \cup (R_k \bowtie S_k)$

Grace hash join

for *build relation* R and *probe relation* S ,

1. **partition** R and S into k partitions each, $k = B - 1$
 2. **probing phase**: hash $r \in R_i$ with $h'(r.A)$ to table T
 - 2.1. $\forall s \in S_i, r \in \text{bucket } h'(s.B)$: output (r, s) if match
- I/O cost: $3(|R| + |S|)$ (no partition overflow)
 - $B > \frac{f \times |R|}{B-1} + 2$ (input & output buffer) $\approx B > \sqrt{f \times |R|}$
 - during probing, $B >$ size of each partition + 2
 - **partition overflow** if R_i cannot fit in memory
 - recursively apply partitioning to overflow partition

General join conditions

- **multiple equality-join** conditions: $(R.A = S.A) \wedge (R.B = S.B)$
 - index nested loop join: use index on some/all join attris
 - sort-merge join: sort on *combination* of attributes
 - other algos: no change
- **inequality-join** conditions: $(R.A < S.A)$
 - index nested loop join: requires B^+ -tree index
 - not applicable: sort-merge join (too much rewinding), hash-based joins
 - other algos: no change

07. Transaction

- An active Xact is a Xact still in progress
- Schedule** = A list of actions from a set of Xacts where the order of the actions within each Xact is preserved
- Serial Schedule** = A schedule where the actions of Xacts are not interleaved
- We say that T_j **reads O from T_i** in a schedule S if the last write on O before $R_j(O)$ in S is $W_i(O)$
- We say that T_j reads T_i if T_j has read some object from T_i
- We say that T_i performs the final write on O in a schedule S if the last write on O is $W_i(O)$
- An interleaved Xact schedule is **correct** if it is "equivalent" to some serial schedule over the set of Xacts

View Serializable Schedules (VS)

Two schedules S and S' (over the same set of Xacts) are **view equivalent** denoted by $S \equiv_v S'$ if they satisfy:

- 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 , the Xact (if any) that performs the final write on A in S must also perform the final write on A in S'

A schedule S is a **view serializable schedule (VS)** if S is view equivalent to some serial schedule over the same set of Xacts

Anomalies due to interleaved schedules

Two actions on the same object **conflict** if

- at least one is a write action,
- the actions are from different Xacts

The following are anomalies:

- Dirty read (due to WR conflicts)
 - T_2 reads O modified by T_1 and T_1 has not yet committed.
 - $W_1(x), R_2(x)$
 - T_2 can see an inconsistent DB state
- Unrepeatable read (due to RW conflicts)
 - T_2 updates O that T_1 reads and T_2 commits while T_1 is still in progress
 - $R_1(x), W_2(x), Commit_2, R_1(x)$
 - T_1 can get a different value if it reads O again
- Lost update (due to WW conflicts)
 - T_2 overwrites O that was modified by T_1 while T_1 is still in progress
 - $R_1(x), R_2(x), W_1(x), W_2(x)$
 - T_1 's update is lost
- Phantom read
 - T re-executes a query on a predicate and gets a different set of results due to a recently committed T.
 - Can be prevented by **predicate locking**: Grant T an S on p, another T' request for X on p is blocked. Also see **index locking**

Conflict Serializable (CS)

- conflict equivalent** denoted by $S \equiv_c S'$ if they order every pair of conflicting actions of two **committed** Xacts in the same way.
- conflict serializable schedule (CS)** if it is conflict equivalent to a serial schedule over the same set of Xacts.

- Conflict serializability graph** denoted as $CSG(S) = (V, E)$ s.t
 - V contains a node for each committed Xact in S
 - E contains T_i, T_j if an action in T_i precedes and conflicts with one of T_j 's action.
- Theorem:** A schedule is **CS** iff its CSG is acyclic.
- Theorem:** If S is **CS**, then S is also **VS**.
- Note:** CS3223 uses **serializable** to mean **CS**
- A write on O by T_i is a **blind write** if T_i did not read O prior to the write
- Theorem:** If S is **VS** and S has no blind writes, then S is also **CS**.

Recovery

- Cascading abort:** For correctness, if T_i read from T_j then T_i must abort if T_j aborts. We say that T_j 's abort is cascaded to T_i .
- Recoverable schedule:** For every Xact T that commits in S, T must commit after T' if T reads from T' .
- Recoverable schedules guarantee that committed Xacts will not be aborted but cascading aborts of active Xacts are possible.
- Cascadeless schedule:** If whenever T_i reads from T_j in S, $Commit_j$ must precede this read action.
- Theorem:** A cascadeless schedule is recoverable.
- Before-images:** Log the before-images of writes to undo the aborted Xacts. See Chap 10. But this does not always work.
 - $W_1(A), W_2(A), Abort_1$. Undoing $W_1(A)$ is incorrect.
- Strict:** For every $W_i(O)$ in S, O is not read or written by another Xact until T_i either aborts or commits.
 - Recovery using before-images is more efficient.
 - Concurrent executions are more restrictive.
- Theorem:** $Strict \subsetneq Cascadeless \subsetneq Recoverable$.

08. Concurrency Control

Transaction Scheduler

For each input action (R, W, C, A) to the scheduler,

- output action to S
- Postpone action by blocking Xact
- Reject action and abort Xact

Lock-Based concurrency control

Lock compatibility matrix						
Lock Requested	-	IS	IX	S	X	
IS		✓	✓	✓	✓	×
IX		✓	✓	✓	×	×
S		✓	✓	×	×	✓
X		✓	×	×	×	×

- If lock request not granted, T becomes **blocked**, execution is suspended and T is added to O's request queue.
- When a lock on O is released, lock manager checks request of first T in request queue for O. If can be granted, T acquires lock on O and resumes after popped from queue.
- When an Xact commits/aborts, all locks are released and T is removed from any request queue it's in.
- $S_i(O)$: Xact T_i requests S on O
- $X_i(O)$: Xact T_i requests X on O
- $U_i(O)$: Xact T_i releases lock on O

Two Phase Locking (2PL) Protocol

- To read O, T must hold S or X on O.
- To write O, T must hold X on O.
- Once T releases a lock, T cannot request anymore.
- 2PL = growing and shrinking phase.
- Theorem:** 2PL S are CS.

Strict 2PL

- Same as 2PL points 1 and 2.
- T must hold onto locks until T commits or aborts.
- Theorem:** Strict 2PL S are strict and CS.

Lock Management

Deadlock: cycle of T waiting for locks to be released by each other.

Deadlock Detection

- Waits-for graph (WFG)**
 - Nodes represent active T.
 - Add edge $T_i \rightarrow T_j$ if T_i waiting for T_j to release lock.
 - Remove edge when lock request is granted.
- Deadlock detected if WFG has a cycle.
- Break deadlock by aborting a T in the cycle.
- Alternative: timeout.

Deadlock Prevention

- Older T have higher priority than younger T
 - Each T is assigned timestamp when started.
 - Older T has a smaller timestamp.
- Suppose T_i requests for a lock that conflicts with a lock held by T_j

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

- Wait-die policy
 - non-preemptive: only requesting T can be aborted
 - younger T may abort repeatedly
 - T that has all locks is never aborted
- Wound-wait policy is preemptive
- To avoid starvation, a restarted T must use original timestamp

Lock Conversion

Increases concurrency by allowing lock conversions, previously serial only schedules can become interleaved.

- $UG_i(A)$: T_i upgrades S on A to X.
 - Blocked if another T is holding S on A.
 - Allowed if T_i has not released any lock.
- $DG_i(A)$: T_i downgrades X on A to S.
 - Allowed if T_i has not modified A and T_i has not released any lock.

Lock-based Isolation levels

Isolation Level	Dirty Read	Unrepeatable Read	Phantom Read
READ UNCOMMITTED	possible	possible	possible
READ COMMITTED	not possible	possible	possible
REPEATABLE READ	not possible	not possible	possible
SERIALIZABLE	not possible	not possible	not possible

Degree	Isolation level	Write Locks	Read Locks	Predicate Locking
0	Read Uncommitted	long duration	none	none
1	Read Committed	long duration	short duration	none
2	Repeatable Read	long duration	long duration	none
3	Serializable	long duration	long duration	yes

- Short duration** lock acquired for an operation can be released after operation ends before T commits/abort
- Long duration** lock acquired for an operation is held until T commits/abort

Lock Granularity

- Highest(coarsest) to lowest(finest): db, relation, page, tuple
- If T holds M on D, T implicitly holds M on granules finer than D.
- Protocol:** Before acquiring S/X on G, acquire IS/IX on granules coarser than G top-down. Release locks bottom-up.

NOTATION

Notation	Meaning
r	relational algebra expression
$ r $	number of tuples in output of r
$ r $	number of pages in output of r
b_d	number of data records that can fit on a page
b_l	number of data entries that can fit on a page
F	average fanout of B ⁺ -tree index (i.e., number of pointers to child nodes)
h	height of B ⁺ -tree index (i.e., number of levels of internal nodes)
	$h = \lceil \log_F(\lceil \frac{ R }{b_l} \rceil) \rceil$ if format-2 index on table R
B	number of available buffer pages