

Relational Algebra

* - requires both relations to have same schema

Operator	Type	SQL	Description
Π	unary	SELECT	<ul style="list-style-type: none"> vertical operator select attributes (cols) <p>set semantics { may lose tuples if we remove the only column two tuples differ in }</p>
σ	unary	WHERE	<ul style="list-style-type: none"> horizontal operator select subset of tuples satisfying some predicate
ρ	unary	AS	renames attributes
\cup^*	binary	UNION	<ul style="list-style-type: none"> returns relations containing tuples from both relations lose tuple if found in both relations
$-^*$	binary	EXCEPT	<ul style="list-style-type: none"> returns relation containing all tuples from L not found in R match each tuple $t \in L$ w/ each tuple $t \in R$
\cap^*	compound	INTERSECT	intersection of both relations
\bowtie	compound	INNER JOIN	inner join on θ
\bowtie	compound	NATURAL JOIN	natural join of L & R on every column w/ same name
γ	extension	GROUP BY	groups by column and aggregates by specified agg.
		HAVING	

External Hashing

- use hash function h_p
 - partition the data
 - stream partitions to disk
 - if $B \in$ buffer:
 - split data into $B-1$ partitions
 - if $n_{\text{partitions}} \leq B$
 - make in-memory
 - load them in \rightarrow hash table for each partition
- \rightarrow can apply duplicate removal, agg., etc.
- in-memory hash table must use separate hash h_r

Q): partitions too big after first pass?

A): recursive partitioning

- for each partition, apply new hash h_n
 - splits partitions into new smaller ones
- repeat until partitions fit in memory
 - \rightarrow every hash must be new

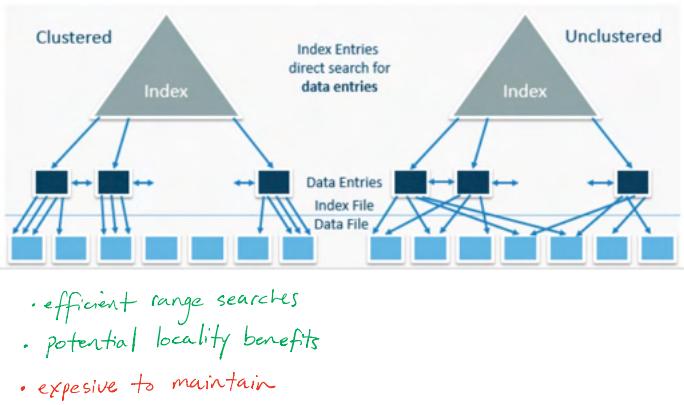
con): not good if more than B pages
of duplicates

How data is stored in the index

- k:** the search key, (a subset of the table's columns)
 - e.g. date of birth, or (lastname, firstname), etc.
- Three alternatives:
 - Alt 1. By value:** actual data record (with key value k)
 - Alt 2. By reference:** <k, rid of matching data record>
 - Alt 3. By list of refs.:** <k, list of rids of all matching data records>

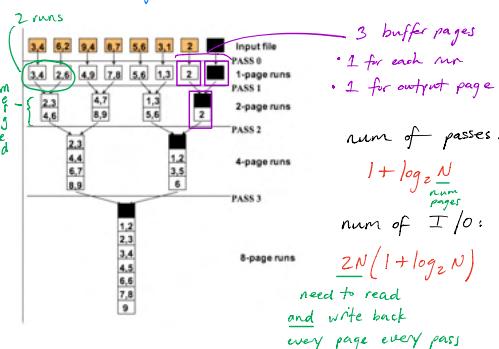
Clustering Overview

- Clustered:** data records are stored in approximate order by key
- Unclustered:** data records are in any order, not sorted by key



Sorting

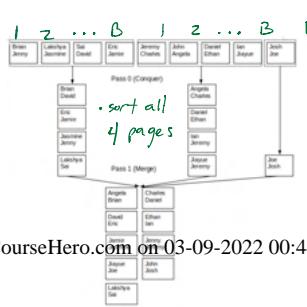
- 2-way external merge
- sort small amounts of data into runs of sorted tuples
- given 2 runs of sorted tuples
 - merge to form 1 larger run



full external merge sort

- assume B buffer pages
- load B pages, sort all at once
 - look at first tuple of each run, output tuple w/ lowest value
- merge $B-1$ runs at once
 - min PQ DS

$B = 4$



num of passes:

$$1 + \log_{B-1} (\lceil \frac{N}{B} \rceil)$$

num runs after pass 0

num of I/Os:

$$2N(1 + \log_{B-1} (\lceil \frac{N}{B} \rceil))$$

same principle as above

$SNLJ [R] + Pr [R][S]$

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

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

$INLJ [R] + |R| \cdot \frac{\text{cost} + \text{find matching tuples } \in S}{\text{differ based on attr.}}$

alt 1 : cost to traverse from root → leaf + read all leaves w/ matching tuples

alt 2/3 : cost of retrieving RIDs + cost of fetching actual records

$SMJ \underbrace{\text{cost}(\text{sort}(R))}_{\text{must satisfy cond. for optimization}} + \text{cost}(\text{sort}(S)) + [R] + [S]$

possible optimization iff $(\# \text{ runs in last pass for } R + \# \text{ runs in last pass for } S) \leq B-1$:
 $\rightarrow \text{reduces I/o cost by } 2([R]+[S])$

$GHTJ \underbrace{\text{cost(hashing)}}_{\text{cost of reading n each page in both partitions}} + \text{cost(Naive Hash Join)}$

$\text{cost}(\text{hash}(P)) : \text{I/o needed to read all pgs in } P + \text{I/o needed to write all resulting partitions}$

Selectivity Estimation - Connectives

Predicate	Selectivity	Assumption
p1 AND p2	$S(p1) * S(p2)$	Independent predicates
p1 OR p2	$S(p1) + S(p2) - S(p1) * S(p2)$	Independent predicates
NOT p	$1 - S(p)$	

Selectivity Estimation - Inequalities on Floats

Predicate	Selectivity	Assumption
$c < v$	$(\max(c) - v) / (\max(c) - \min(c))$	We know max(c) and min(c). c is a float.
$c \geq v$	$1 / 10$	We don't know max(c) and min(c). c is an integer.
$c \leq v$	$(v - \min(c)) / (\max(c) - \min(c))$	We know max(c) and min(c). c is a float.
$c \geq v$	$1 / 10$	We don't know max(c) and min(c). c is a float.

* We add 1 to the denominator in order for our [low, high] range to be inclusive.
E.g. range [2, 4] = 2, 3, 4 → $(4 - 2) + 1 = 3$

Selectivity Estimation - Equalities

Predicate	Selectivity	Assumption
$c = v$	$1 / (\text{number of distinct values of } c \text{ in index})$	We know c .
$c \neq v$	$1 / 10$	We don't know c .
$c1 = c2$	$1 / \text{MAX}(\text{number of distinct values of } c1, \text{number of distinct values of } c2)$	We know c1 and c2 .
$c1 \neq c2$	$1 / (\text{number of distinct values of } c)$	We know c but not (other column).
$c1 = c2$	$1 / 10$	We don't know c1 or c2 .

Transactions

- collections of operations
- a single logical unit
- would like all components to happen at once
- commit
 - save changes
 - successful tx
 - abort
 - reverts changes
 - unsuccessful tx
- should follow ACID properties

property	description	associated topic
Atomicity	all operations in tx happen or none	recovery
Consistency	database constraints (e.g. unique values) are maintained	
Isolation	should look like we only run one transaction at a time	concurrency locking
Durability	once a transaction commits, it should persist	recovery

Conflict Serializability

a): How to check if schedule is serializable?

b): you don't

► instead, check if schedules are conflict serializable

• stronger condition than serializability

$R(A) \vee W(A)$

• two operations in schedule conflict if:

- at least one operation is a WRITE
- they are on diff. TXs
- they work on same resource

a): Check for conflict equivalence/serializability?

A): Dependency graph (precedence graph)

- if all arrows are in same direction in dependency graph for a/s schedule
- two schedules are conflict equivalent

Theorem): All conflict serializable schedules have an acyclic dependency graph

► check for cycle in dependency graph

- if cycle → not conflict serializable
- else → conflict serializable

Dependancy Graphs (precedence)

transaction

transaction	1	2	3	4	5	6
schedule	$R(A)$	$R(A)$	$R(B)$	$R(B)$	$R(B)$	$W(B)$
	T ₁	T ₂	T ₃	T ₄	T ₅	T ₆

Concurrency at same time

- benefits of running many txs concurrently:

(a) throughput argument

- increases processor/disk utilization
work on another tx while one is busy waiting for data from disk

► more transactions per second

(b) latency argument

- run a lot of txs at once then a long running tx doesn't cause short tx to wait long time to start
- transactions may finish faster

note): still must maintain ISOLATION

• two schedules are conflict equivalent if every conflict is ordered the same way

• a schedule is conflict serializable if:

- it is conflict equivalent to a serial schedule

Serializability Hierarchy

$R(B) \in T_2$ comes first so T_2 has edge to T_1

one node per transaction

Deadlocks

method	description	approaches
Deadlock avoidance	avoid them at all costs	<ul style="list-style-type: none"> • wait-die • wait-and-wait
Deadlock detection	catch them and abort tx in deadlock	<ul style="list-style-type: none"> • waits-for graph • kill inactive txs

Deadlock Avoidance

Scenario): T_i wants lock but T_j has conflicting lock

- ① Wait-Die
 - T_i has higher priority \rightarrow waits for T_j to release lock
 - T_i has lower priority \rightarrow abort T_i (dies)
 - txs can only wait on lower priority txs
 - \blacktriangleright lowest priority txs cannot wait
- ② Wait-and-Wait
 - T_i has higher priority \rightarrow cause T_j to abort (wait)
 - T_i has lower priority \rightarrow waits for T_j to finish
 - txs can only wait on higher priority txs
 - \blacktriangleright highest priority txs cannot wait

2-Phase Locking (2PL)
idea): enforces conflict serializability

- method
 - tx may not acquire a lock after it has released any locks
 - two phases
 - Phase 1
 - start until some lock is released
 - tx is only acquiring locks
 - Phase 2
 - until end of tx
 - tx is only releasing locks

problem): doesn't protect from cascading aborts

- if some earlier tx T_1 aborts in the end, it will affect all other txs
 - \blacktriangleright all txs have to be rolled back

why? it allows a new tx to read new values before the transaction commits

Strict 2PL
idea): avoid cascading aborts

- method
 - same as 2PL except all locks get released together after a tx completes

2PL

Strict 2PL

This study source was downloaded by 100000806014767 from CourseHero.com on 03-09-2022 00:44:03

<https://www.coursehero.com/file/98747426/186-Midterm-1-Cheat-Sheet-1pdf/>

Deadlock Detection

- "waits-for" graph

no deadlock

deadlock

- cycle in graph \rightarrow deadlock
- abort one of $T_i, T_j \rightarrow$ end deadlock

Undo Logging

- Write log records to ensure **Atomicity** after a system crash:
 - <START T>: transaction T has begun
 - <COMMIT T>: T has committed
 - <ABORT T>: T has aborted
 - <TX, X, v>: T has updated element X, and its old value was v
- If T commits, then **FLUSH(X)** must be written to disk before <COMMIT T>
 - Force – we can UNDO any modifications if a Xact crashes before COMMIT
- If T modifies X, then <TX, X, v> log entry must be written to disk before **FLUSH(X)**
 - Steal – we can UNDO any modifications if a Xact crashes before FLUSH

Undo/Redo Logging

- UNDO logging provides **atomicity**
 - Undoes all updates for **running** transactions
- REDO logging provides **durability**
 - Redos all updates for **committed** transactions
- Requires us to recover from the beginning of the log!
 - Computationally expensive
- What if we could start mid-way through the log?
- What if we could get the best of both? **Steal / No-Force**

Write-Ahead Logging

- Two Rules:
 - A transaction is not committed until the log records for all its changes have been written to disk
 - Changes must be logged before the actual data is modified on disk
- With these two rules we can develop a system that recovers properly from failures

Redo Logging

- Write log records to ensure **Durability** after a system crash:
 - <START T>: transaction T has begun
 - <COMMIT T>: T has committed
 - <ABORT T>: T has aborted
 - <TX, X, v>: T has updated element X, and its new value was v
- If T modifies X, then both <TX, X, v> and <COMMIT T> must be written to disk before **FLUSH(X)**
 - No-Steal, No-Force – we can REDO any modifications if a Xact crashes before FLUSH

Force/No-force

Whether or not modified pages from a transaction are forced to disk on commit

- **Force** = Enforce
 - Easy to maintain durability
 - If we crash after committing, everything was already written to disk!
- **No-Force** = Don't Enforce
 - Faster, more difficult to maintain durability
 - Don't have to do unnecessary writes
 - If DB crashes mid-transaction, we'll need to redo any changes that didn't make it to disk

No-Steal, Force Example

T_1	Write(A)	Write(B)			*CRASH*	Commit
T_2			Write(B)	Commit		

- **No-Steal** = Don't Allow modified pages to be flushed to disk before Commit
 - Force** = Force modified pages to be flushed to disk on Commit
 - T_1 : page A will not be modified on disk. T_1 crashed before Commit (**No-steal**)
 - T_2 : page B is modified on disk. (**Force**)
- Durability of database maintained. Committed Transactions have their modifications flushed to disk.
- Atomicity is not maintained. T_1 didn't commit, but its changes to page B still made it to disk!

Steal, No-Force Example

T_1	Write(A)	Write(C)	*Flush to disk*		*CRASH*	Commit
T_2			Write(B)	Commit		

- **Steal** = Allow modified pages to be flushed to disk before Commit,
 - No-Force** = Do not force modified pages to be flushed to disk on Commit
 - T_1 : Pages A, C are modified on disk. (**Steal**)
 - T_2 : page B not modified on disk (**No-force**)
- Without additional policies in place, Atomicity and Durability are not maintained:
 - T_1 's changes are flushed to disk despite not committing,
 - T_2 's changes are not flushed to disk despite committing.

Multigranularity Lock Compatibility Matrix

	NL	IS	IX	SIX	S	X
NL	•	•	•	•	•	•
IS	•	•	•	•	no	
IX	•	•	•	no	no	no
SIX	•	•	no	no	no	no
X	.	no	no	no	no	no

Steal/No-steal

Whether or not modified pages from an uncommitted xact can be flushed to disk

- **No-Steal** = Don't Allow
 - Easy to maintain atomicity
 - If we crash mid-transaction, none of the changes were written to disk!
- **Steal** = Allow
 - Allows buffer manager to use frames optimally
 - Allows uncommitted changes to make it to disk so we may need to undo these

GMT -06:00 avoids some problems

- Allows buffer manager to use frames optimally
- Allows uncommitted changes to make it to disk so we may need to undo these

Parallel Query Processing

Types of Query Parallelism

- inter-query parallelism
 - works using queries as the unit of parallelism
- intra-query parallelism
 - works within a single query

Partitioning Data

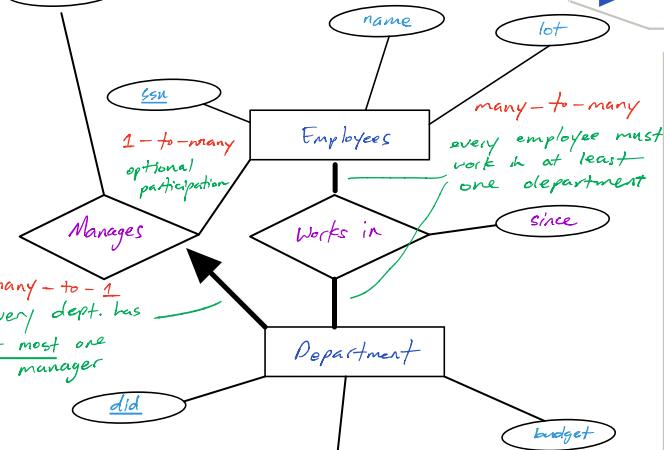
- range partitioning
 - on a key
 - divides data based on which range the key falls into
- hash partitioning
 - uses hash function
- round-robin partitioning
 - cycles through partitions in order as the data comes in
 - not ordered by key

Asymmetric Shuffles

- when data is already partitioned the desired way

Broadcast Joins

- sometimes one relation \Rightarrow huge other relation \Rightarrow small
- much cheaper to send all of tiny tables multiple times [than] partitioning huge table
- since



Weak Entities

- entity that can be identified uniquely only with the key of another entity
- will be a partial key dashed underline
- identifies the weak entity when combined w/ owner entity's key
- must be a 1-to-many relationship
 - 1 owner entity \rightarrow many weak entities
- total participation

Document Stores

- Purpose: when you don't want to allow too much flexibility
- Data consists of semi-structured data (JSON, XML)
- Advantages:
 - One of the most flexible (least-structured) data models
 - Values can be almost anything

- Operations:
 - persistency
 - No operations based on value because values are supposed to be incredibly flexible
- Distribution:
 - Key is sent to a machine based on its hash value
 - No replication: key exists on machine h(key)

JSON (Document Stores)

- Self-describing: each document can have its own schema
- Supports:
 - Objects: collection of (key, value) pairs (denoted by {}, {})
 - Key: string
 - Value: object, array, or atomic
 - Arrays: list of values (denoted by [], [])

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

•

Two-Phase Commit Logging

(1) Voting Phase

- coordinator → participants: PREPARE message
- participants
 - log & flush either a PREPARE or ABORT log record
 - log entry keeps track of coordinator ID
 - if participant votes NO
 - > transaction will abort
 - > begin cleaning up transaction
- participants → coordinator: vote YES or NO
- coordinator
 - log & flush after a COMMIT or ABORT log record
 - commit log entry contains all participant IDs

(2) Results Phase

- coordinator → participants: COMMIT or ABORT message
- participants
 - log & flush COMMIT or ABORT record to log
 - participants → coordinator: ACK message
- coordinator
 - add END record to log
 - remove txn from TxnTbl
 - no flush required

2PC : Blocking

Q): node goes down during 1st phase (voting)?

- all part. that voted YES
 - keep locks
 - wait for COMMIT/ABORT message

Q): a participant never recovers?

- coordinator can
 - reassign part.
 - recover from log
- continue
- if participant comes up again → tell it to recycle itself

Q): what if coordinator does not recover?

- checkmate!

Log Records

LSN	prevLSN	Txn Num	Log Payload
10	—	1	START
20	10	1	UPDATE
30	—	2	START
40	20	1	COMMIT

LSN of last log record for same transaction

Compensation Log Record (CLR)

idea): write a CLR for each undone operation

- write before actually undoing data
- undoNextLSN → next LSN to undo
- determined by prevLSN of current log record that is being undone

Data Structures: Transaction Table

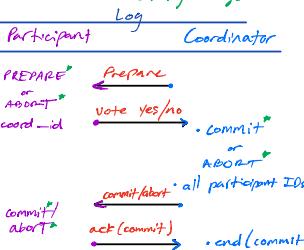
- tracks active transactions

(canceling, aborting, committing)
status of transaction

This study source was downloaded by 100008060114767 from CourseHero.com on 03-09-2022 00:44:03 GMT -06:00

1	R	33
2	C	42

2PC * → flush log before sending message



2PC : Handling Failures

- if coordinator thinks a participant went down:
 - abort txn if participant has not voted
 - if waiting for ACK
 - reperform recovery
 - must finish committing what we started
- if participant thinks coordinator is down:
 - if no PREPARE record in log
 - abort unilaterally (vote no)
 - if PREPARE record in log
 - all part. voted YES already
 - cannot proceed unilaterally

2PC : Recovery

- if commit/ABORT log record
 - do what usually do
 - if coordinator record contains all participant IDs
 - periodically send commit/abort messages until we get all ACKs back
- if PREPARE log record:
 - this is a participant node ID obtained from log record
 - participant → coordinator from log record
 - send message inquiring about txn status
- else
 - did not vote YES → abort the txn
 - if we receive VOTE YES/NO messages
 - this is the coordinator
 - respond with ABORT message

optimization): presumed abort

- assume a txn aborts if we have no log records
- part. that receive ABORT logs don't send ACKs

→ reduces messages sent/logging

enter a txn aborts:

- coordinator cleans up locally
 - remove txn from TxnTbl without ACKs

part. IPs do not need to be stored in ABORT records

→ not sending ACKs . ABORT records do not need to be flushed

Keeping Track of Page Updates

page LSN → LSN of last record that modified the page

stored in the page

updated in-memory when page is modified since from checkpoint

flushed to disk w/ pg

indicates how up-to-date the page on disk actually is

in-memory page LSN can differ from disk page LSN

Analysis Phase

(goal): Rebuild

TransactionTbl at time of crash

Dirty PgTbl

LSN after BEGIN CHECKPOINT record

ABORT Protocol

(1) Write ABORT record to log manager

(2) Undo all changes made by transaction.

last LSN → prevLSN → backwards
start backrecords via previous of each record

undo all log records that can be undone

→ adding CLR to log manager

(3) Write END record.

Redo Phase

(goal): redo everything from earliest recLsn < LSN

get back unflushed changes before crash

(a) pg not in DPT

→ pg on disk is up-to-date, no changes needed

(b) recLSN_pg > LSN

→ no need to undo

→ recLSN_pg <= 1st record to dirty pg

→ change is already flushed

(c) pg(LSN) = LSN

→ determines which changes have already been applied

Data Structures: Dirty Page Table

tracks dirty pages → modified in memory but not flushed to disk yet

Page num.	recLSN
46	11
63	24

recovery LSN indicating the first log record that caused the page to get dirty

Flushed LSN

max LSN that was successfully flushed to disk

$$\text{page LSN}_A \leq \text{flushed LSN}$$

before A can be flushed to disk, need to ensure that corresponding log records are on disk

ensures ATOMICITY

Checkpoints

idea): do not want to start from beginning each crash!

write BEGIN CHECKPOINT record

last record to start analysis at

records later than this may have changes that may or may not be reflected E checkpoint

write END CHECKPOINT record (sometimes later)

contains DPT
TransactionTable

after BEGIN checkpoint
do not know exactly when!

UNDO

TransactionTbl

- Any remaining transactions in the ATT didn't make it to completion, but may have some updates on disk (due to STEAL)
- We can't let their effects persist because of atomicity

UNDOLIST

UNDO

UNDOLIST

UNDOLIST