

- ### JOIN R(M) outer S(N) inner RMS
- ▷ Simple nested loop  $M + (M \times N)$
  - ▷ Block nested  $M + (\lceil \frac{M}{B} \rceil \times N)$
  - ▷ Index nested loop  $M + (M \times C)$  index lookup cost
  - ▷ Sort-Merge (good if R or S sorted / has cluster index)  
Sort R:  $2\{M, N\}(1 + \lceil \log_{B-1} \lceil \frac{M, N}{B} \rceil \rceil)$   
Merge:  $M+N$  if mostly unique,  $MN$  if all dup
  - ▷ Basic Hash (opt: bloom; good if know R fits)
    1. Build htab for R w h<sub>1</sub>
    2. Prob htab for every s ∈ S
  - ▷ Grace Hash (use disk, usually faster than sort-merge)
    1. Build htab for R & S with h<sub>1</sub>. Recursively partition with h<sub>2</sub>... if needed
    2. For each corresponding buckets do some loop join

- If no rec part & htab on disk:
- Part cost  $2 \times (M+N)$   
Prob cost  $M+N$
- Extern Merge Sort  
# passes =  $1 + \lceil \log_{B-1} \lceil \frac{M, N}{B} \rceil \rceil$   
IO cost =  $2N \cdot \# \text{ passes}$

F	N-F
WAL	flusher
S	
N-S	easy slow

- Half full: leaf  $\geq \lceil \frac{M-1}{2} \rceil$  keys internal  $\geq \lceil \frac{M}{2} \rceil$  ptrs
- Query Execution Models: top-down / bot.-up
- ▷ Iterator aka Volcano, Pipeline  
Next(), Open(), End()
  - ▷ Materialisation + OLTP → Can prop down LIMIT etc.  
Output() to get entire result
  - ▷ Vectorisation + OLAP + use vectorized instrs  
Next() returns batch

- ### Access Method
- ▷ Seq Scan → Prefetch → Bypass buf pool → Late materialise → Data Skipping: be lazy / keep zone map (page aggs)
  - ▷ Index Scan → histogram → multi-index (set eps)

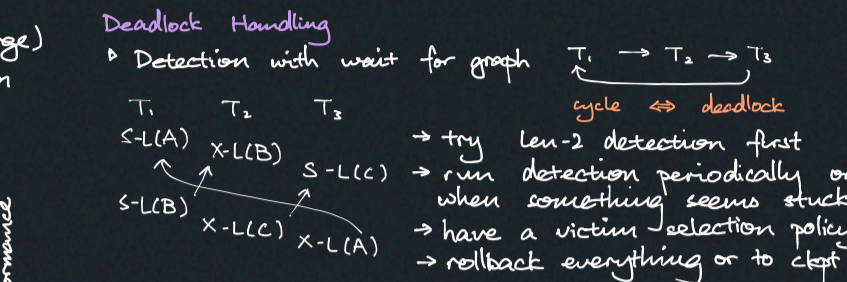
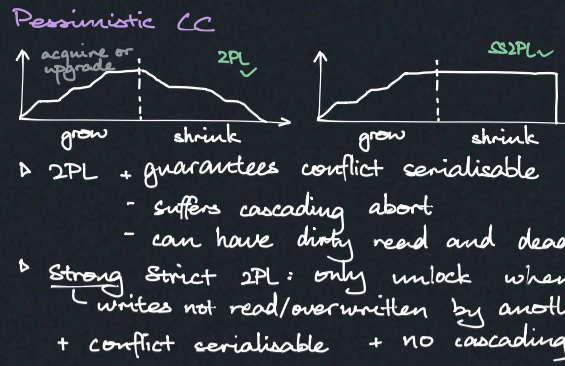
Del / update: Track modified records ids  
Halloween Problem: updated tuple moved and read again

Scheduling  
Node just calls children / put stuff in work queue  
walk dep. tree

- ### Intra-Query Par
- ▷ Horizontal gather distribute repartition - exchange ops
  - ▷ Vertical
  - ▷ Busy (do both)

- ### I/O Par
- ▷ RAID controller, multi disks behave like one to DBMS
  - ▷ Partition - shared log, separate directory

- ### Query Optim
- logical plan → physical plan, NP-hard w.r.t. # of joins
- ▷ Rules → push predicate σ down →  $\sigma_{A.id=B.id}(A \times B) \equiv \sigma_{A.id=B.id}(A) \times B$
  - push projection π down
  - ▷ Cost based, with cost model on  $\{CPU, Disk I/O, RAM, network\}$ 
    - ▷ Single-rel: find best access path for each col
    - ▷ Sys R: DP to find min cost among all left-deep join
    - ▷ Bot.-up: start from data source, reach logical query
    - ▷ Top-down: start from query, search down to find best physical

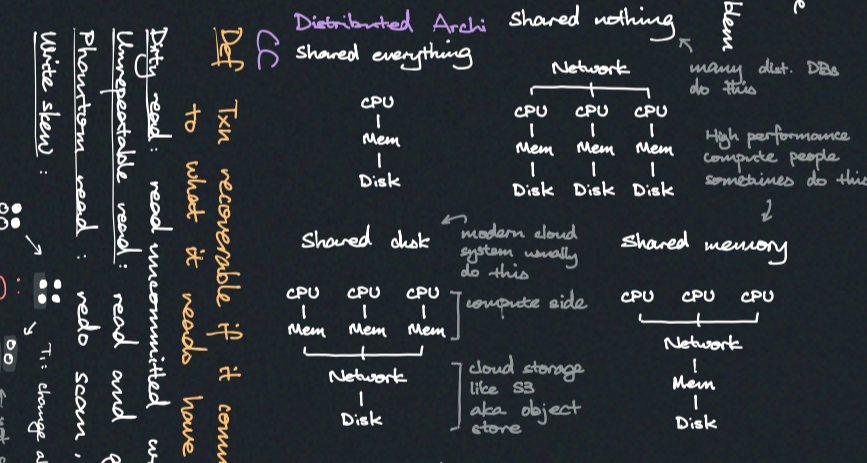
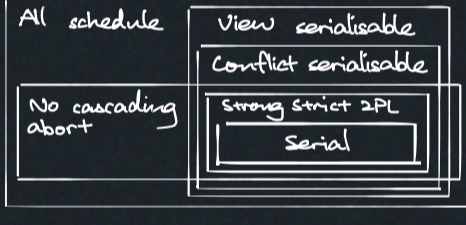


- Prevention: kill a txn if two try to grab same lock, based on some priority (say higher for older txns) (when dead txn restarts, reuse orig timestamp)
- requesting txn: if higher priority than holding txn then  
else  
wait + die  
Wound + wait

- ### Distributed Joins
- select \* from R join S on R.id = S.id, 2 nodes
1. If S replicated: do local join and concat result
  2. If S, R partitioned by same ranges: again do local join
  3. If S, R partitioned differently and R part. by join key  
broadcast to case 1 to have S copied on both
  4. Else - neither S, R part. by join key  
split up networking to go to case 2  
like hash join, just with hashed partitions on nodes
- Semi-join: project join key and send only that

### Dist Query Exec

- ▷ Push query to data send query to where the data is
- ▷ Pull data to query Bring needed data to processor running the query

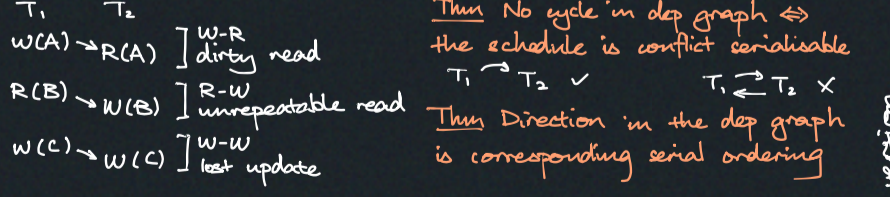


IS	IS	S	IX	S	IX	S	IX	S	IX
X	✓	✓	✓	✓	✓	✓	✓	✓	✓
IS	IX	S	IX	S	IX	S	IX	S	IX
X	✓	✓	✓	✓	✓	✓	✓	✓	✓
IS	IX	S	IX	S	IX	S	IX	S	IX
X	✓	✓	✓	✓	✓	✓	✓	✓	✓
IS	IX	S	IX	S	IX	S	IX	S	IX
X	✓	✓	✓	✓	✓	✓	✓	✓	✓

- ▷ Dirty read: read uncommitted write
- ▷ Dirty write: read uncommitted write
- ▷ Phantom read: read and get different result
- ▷ Phantom write: redo scan, different result
- ▷ Dirty write: not committable
- ▷ Dirty read: non-overlapping write sets
- ▷ Dirty write: dirty
- ▷ Dirty read: dirty
- ▷ Dirty write: dirty
- ▷ Dirty read: dirty
- ▷ Dirty write: dirty

- ### Cost Estimation
- Histogram → Independence assumption (?) → Take samples
  - Assume uniform dist. → Inclusion assumption: join key match across rels

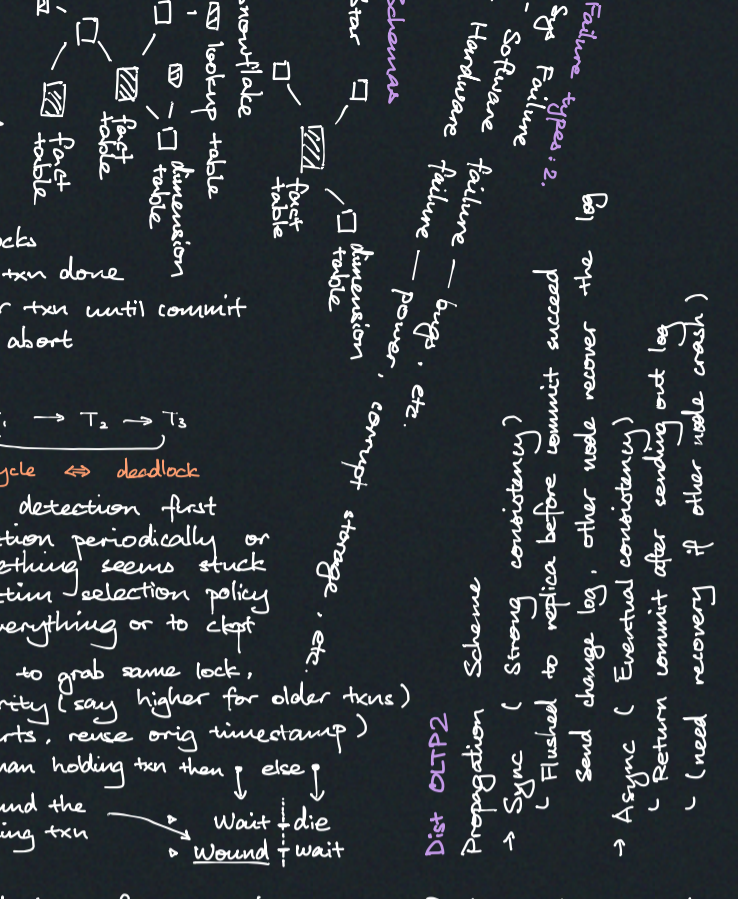
- ### Txn Manage, wanted:
- ▷ Atomicity All or nothing, no partial completion
  - ▷ Consistency Start consistent and end consistent
  - ▷ Isolation Different txns isolated
  - ▷ Durability Persist after txn commit



- ### Verify conflict serialisable: use dep. graph or keep swapping ops not in conflict to get serial schedule
- Def View serialisable allows blind write i.e. write without read (other constraints same as conflict serialisable)
- S<sub>1</sub>, S<sub>2</sub> view = p
  - T<sub>1</sub> reads initial val in S<sub>1</sub> ⇔ T<sub>1</sub> reads initial val in S<sub>2</sub>
  - T<sub>1</sub> reads T<sub>1</sub>'s write in S<sub>1</sub> ⇔ T<sub>1</sub> reads T<sub>1</sub>'s write in S<sub>2</sub>
  - T<sub>1</sub> does last write in S<sub>1</sub> ⇔ T<sub>1</sub> does last write in S<sub>2</sub>
- Write: make new version  
Read: find the correct version  
Benefits: W doesn't block R they can find the right version and time travel  
R doesn't block W (when W-W, kill second writer)

- ### MVCC (optimistic)
- ▷ Snapshot: like repo at commit, with copy of Txn Table
  - ▷ Storing versions
  - ▷ Delta storage
  - ▷ GC: Tuple level: chaining w by proc  
→ each thread identifies on old-to-new chain only  
▷ Txn level: each txn keeps write set  
→ re-apply committed txns  
→ revert aborted or unfinished txns
  - ▷ Buf Pool Policy  
STAL: uncommitted-txn can overwrite most recently committed value in volatile storage  
FORCE: changes must write to non-volatile before txn commits (and return success to user)
  - ▷ Isolation Levels

- ▷ Serializable - no phantoms, no dirty reads, all reads repeatable
- ▷ repeatable read - get all locks first, index lock, and strong strict 2PL
- ▷ same as above but without index lock
- ▷ read committed - allow phantoms and unrepeatable reads
- ▷ same as above but release S lock immediately
- ▷ read uncommitted - allow all three
- ▷ same as read committed but no S lock (snapshot isolation) - allows write skew
- ▷ allows write skew cuts through timeline
- ▷ all read consistent with a snapshot at txn's start
- ▷ commit only if writes don't conflict with other txns' writes since snapshot taken
- ▷ (Cursor stability) - prevents lost update



- ▷ Storage failure
- ▷ Disk failure
- ▷ Controller problem
- ▷ Failure types: 3

**Timestamp Ordering CC (Optimistic)**

Txn gets unique TS, atomic monotonic, need serialisability in this order

Basic: tuples get R & W TSs, conflict serialisable, may not be recoverable

Read - can't read sth written in the future

- If written by future txn: abort reading thread, restart with new TS

- Else update R-TS, make local copy for repeated reads

Write - can't write correctly if read or written in future

- If bad write: abort and restart self

- Else update W-TS, make local copy

Optimisation with Thomas Write Rule, view serialisable

Allow but ignore blind write: if  $TS(T_i) < W-TS(A)$

Issues: overhead, long txns can get starved

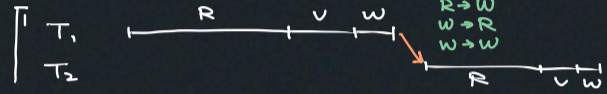
Optimistic CC (OCC), good if txns touch few things, no deadlocks, easy to abort, wast work if validation fail

1. Read Phase - all reads & local writes, work in local workspace track all reads & writes

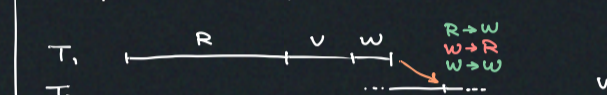
2. Validation - any conflict with other txns? (restart if failed) *declare timestamp at this point*

3. Write Phase - apply local changes to DB

Cases



T1 finishes W before T2 starts W



T1 finishes R before T2 finishes R



Validation Scope:

Backward - Check against those that committed

Forward - Check against those that could commit in future

Write phase:

Serial commit - one txn write at a time

Parallel commit - use latches to write but ensure ordering

Checkpoint

Naive: pause starting new txns, wait for active txns to finish, flush all pages

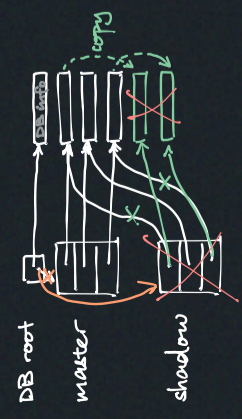
Slightly better: pause write txns, flush modified pages, save ATT and DPT at start of checkpoint

Fuzzy: put in CHECKPOINT-BEGIN, memcopy ATT and DPT, let execution continue while preparing checkpoint, when done, put in CHECKPOINT-END + ATT + DPT (start recovery at CHECKPOINT-BEGIN)

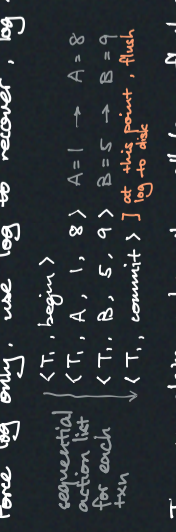
ATT txn modes: running, committing, undo candidate

Phantom Problem when not locking new DB objects eg. no lock on new tuple

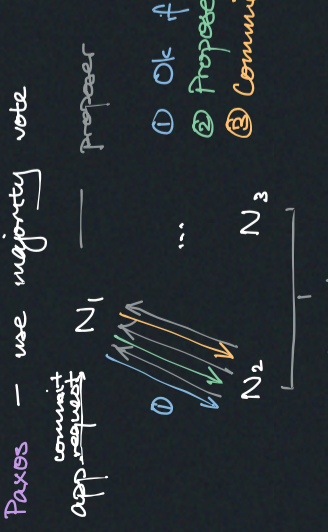
Re-run scan at commit to check for inconsistency  
 Predicate locking  
 Index locking  
 Shadow Paging (N-S,F) -> Virtual lock on non-existent key  
 Gap lock for non-existent ranges  
 Key range lock -> Hierarchical lock  
 Like file table after fork  
 Keep master page table & per-txn shadow page table



Journal file  
 Before modifying a page, make and log a copy  
 Write Ahead Log (WAL)  
 Force log only, use log to recover, log to log



Txn complete only when all logs flushed, flush log before pages  
 Optimisation: Group commits and batch flush  
 Logging schemes:  
 - Physical - byte changes like git diff  
 - Logical - "update all records satisfying"  
 - Physiological - physical across pages, logical within page  
 Paxos - use majority vote

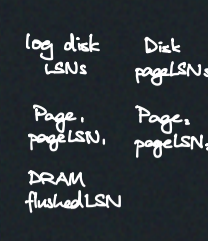
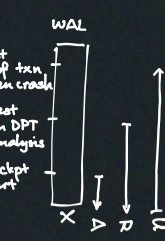


Proposer race - handled by proposer number  
 Multi-Paxos: Select leader to lead for a while  
 Raft: Only node with most up-to-date log becomes leader  
 Consistency Issues (CAP / PACELC)

Want:  
 - Consistency  
 - Always available  
 - Network partition tolerant  
 CAP Theorem: you can only get 2 of those  
 But for certain data struct, like those CRDT stuff, we can get all 3.

Dist OLTP!  
 K-safety: K replica must be avail. for each DB obj  
 Replica config  
 Primary - Replica  
 leader selection if primary down  
 Multi-Primary

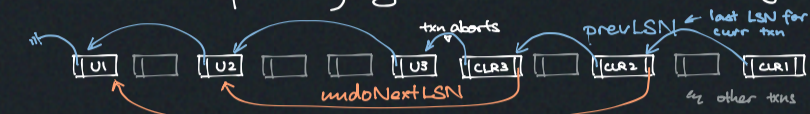
**ARIES**



But pool mgr can be modular  
 pageLSN1, pageLSN2  
 flushed to disk  
 flushed LSN  
 can directly evict this page  
 when evict, need to first flush log

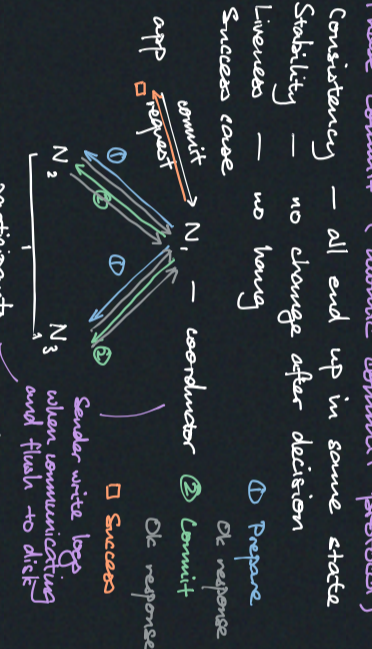
flushedLSN	mem	last logged LSN on disk
pageLSN	page*	newest update to page*
reclSN	DPT Dirty Page Table	oldest update to page* since last flush
lastLSN	ATT Active Txn Table	latest record of txn T <sub>i</sub>
MasterRecord	disk	LSN of last cplt

Commit: write log, ensure flushed, log END-TXN  
 Abort: write compensating log record (CLR) for each change, log END-TXN



Analysis: reconstruct ATT & DPT right before crash  
 - remove ended txns from ATT  
 - add record to ATT with Undo unless Committed for the txn  
 - add written pages to DPT  
 Redo: reapply logged changes to right before crash  
 when done, set committed txns to C and remove from ATT  
 Undo: in rev order undo changes of crashed/aborted txns

2 Phase Commit (atomic commit protocol)  
 - Consistency - all end up in same state  
 - Stability - no change after decision  
 - Liveovers - no hang  
 - Success case  
 - Failure: if any participant says abort, coordinator decides about coordinator tells everyone about  
 - Recovery: In in prepare state, tell coordinator, just say abort, If coordinator and is committing, tell participants commit  
 Optimisations:  
 - If last query, say it's last and ask if prepared  
 - Coordinator say commit success early, before phase 2  
 Distributed Partitioning  
 Physical (usually for shared nothing)  
 - Naive - by table  
 - Vertical - part by column  
 - Horizontal - by key hash, ranges, round robin, ...  
 Logical (usually for shared disk)  
 Middleware: routes queries to right part



Dist OLTP3  
 Prop timing  
 - Continuous log streaming  
 - On commit  
 Active vs passive  
 - Active - Passive: One send one receive  
 - Active - Active  
 Run query on all replica, if all agree then return  
 Ship queries btwn nodes  
 Consistency Issues (CAP / PACELC)  
 Want:  
 - Consistency  
 - Always available  
 - Network partition tolerant  
 CAP Theorem: you can only get 2 of those  
 But for certain data struct, like those CRDT stuff, we can get all 3.