

Data Partitioning

Desirable properties of fragmentation

- completeness: each item in R can be found in ≥ 1 fragment (nothing is lost)
• reconstruction: R can be reconstructed from fragments (must be lossless join)
• disjointness: data items are not replicated

Fragmentation Techniques

- Range Partitioning
– Use predicates on ≥ 1 attribute (e.g. < 100, [100,500], ≥ 500)
– Use catch-all predicate to guarantee correctness OR use a binary tree (p1 ∨ ¬p1)
• Hash Partitioning
– Good hash function + not-skewed data => data will be evenly distributed
– Method 1: modulo method - when adding / removing nodes: need to rehash everything (less elastic)
* In Ri if h(...) mod i is i
– Method 2: consistent hashing - easier to add / remove nodes (more elastic)
* Partition h(...) using n values into v1 < ... < vn (replicated to all nodes OR just master)
· Ist: ≤ v1 and ≠ v1, 2nd: (v1, v2)
* In Ri if h(...) is in Ri's range, Ri is catch-all node
* Non-uniform data & load distribution (can be managed with virtual nodes -> let variable amount map to same physical node)

Derived Horizontal Fragmentation

- Partition a R based on S using semi-join Ri = R x SA Si
• For completeness, RA ⊆ SA
• For disjointness, SA must be key
• So for both, RA must be FK of S and NOT NULL

Vertical Fragmentation

- key(R) must be in all partitions (for disjointness: only keys are duplicated)
• heuristic Attribute Affinity Measure: if commonly referenced together, should be in same partition

Complete Partitioning wrt Query

A partitioning F is complete wrt to Q if for all fragments Ri ∈ F: either return WHOLE partition OR nothing.

Mintern Predicate Partitioning

Mintern predicate m = a combination (positive / negative) ∧ for a set of predicates (|m| = 2^n). Use boolean algebra to simplify

- Q = {Q1, ..., Qk}, Qi = σpi(R)
• P = {p1, ..., pk}
• F = {R1, ..., Rm} where Ri = σmi(R)

Mintern predicate partitioning F (of Q) is always a complete partitioning wrt to the Q. Watch out for ∨.

Query Processing

Make query plan that minimises total cost (CPU, IO, comm) OR response time. Try to parallelise queries.

- 1. Normalisation (rewrite query into normal form)
• more common CNF (p ∨ p) ∧ (p ∨ p)
• DNF (p ∧ p) ∨ (p ∧ p)
• p is simple predicate: single attribute Ai op v
2. Semantic Analysis (check against schema + type check)
3. Simplification & Restructuring

Localisation Program

- Rewrite distributed query into fragment query
• U for horizontal partition; ⊞ for vertical partition

Reduction Techniques

Identify & remove queries that do not contribute to result

- Reduction with Selection: σp(Ri) = ∅ if Ri = σFi(R) and Fi ∧ p = false
• Reduction with Join: Ri ⊞ Sj = ∅ if there's no "intersection" of predicates on join attributes (a)
• Reduction with Derived Frag.: Sj ⊞ Ri = ∅ if Sj is derived from R and i ≠ j
• Reduction with Vertical Frag.: if missing required attribute, drop the fragment

Distributed Join Strategies R ⊞A S

- 1. both R and S are partitioned on join key
2. only R (not S) is partitioned on join key
3. neither are partitioned on join key

Communication Cost

- Collocated: 0
– all servers perform local join -> send results to server
• Directed: size(R) if R is repartitioned (R is NOT previously partitioned)
– repartition -> if in wrong server, send to correct one
• Repartitioned: size(R) + size(S)
• Broadcast: (n - 1) * size(R) where R is smaller one
– broadcast smaller table to ALL servers

Storage

LSM (Log-Structured Merge) Storage

Writing to B+ is random I/O (+ splitting & propagate); use LSM instead (append-only updates)

- Memory Table (in memory hash table with in-place updates)
• After threshold: sorted + flush to disk (sequential I/O)
• SSTables (Sorted String Table): immutable; records are sorted by K; each SSTable associated with range of key values + timestamp
• Commit Log Files used for durability

Compaction of SSTables

- Why?
– improves read performance by defragmenting table records
– improves space utilisation by removing tombstones (must ensure all other versions are gone) & stale values
* can remove at bottom-most level (because guaranteed not in any other SSTable at this level, and no higher levels)
• Size-tiered Compaction Strategy (STCS)
– Each tier has approx. same size
– compaction triggered at tier L if number of SSTables = threshold
* All SSTables at L are merged into one SSTable at L+1
* Tier L becomes empty

- Each object has ≤ 1 version in every SSTable
• Levelled Compaction Strategy (LCS)
– SSTables at level 0 can have overlapping key ranges
– For level ≥ 1
* each SSTable is same size
* key ranges do not overlap within the tier
* SSTable at L overlaps with at most F SSTables at L+1
– Lower level (across tiers) + larger index (within tier) is more recent
– Compaction: new tables stored at L+1, old tables removed
* for L ≥ 1: choose a SSTable (round-robin style with wrap-around) -> merge with all overlapping SSTables at L+1
* for L = 0: merge all SSTables at level 0 with all overlapping SSTables at level 1
* if inserting into SSTable violates its F condition, make new table
* triggered when number of level 0 reaches threshold; for L ≥ 1: size(L) > F * L * MB
– Each level stores F times as much data as previous
* For n records of m MB each, in worst case: last level stores a version of each of the n records.
* Therefore, F^{L-1} < mn ≤ F^L => L = ⌈ logF(mn) ⌉
– Each object has ≤ 1 version in every SSTable (in level 0), has ≤ 1 version in every other level.

Searching LSM

- Start at MemTable; go to next level, start at right-most table
• Check key range first: if within, use binary search
– At each level ≥ 1: either search 0 or one tables

Optimising SSTable Search

SSTables are stored in blocks.

- 1. Sparse index
• (k1, k2, ..., kn) if N blocks
• ki is first key value in block i
• Binary search the sparse index in-memory
2. Bloom filter
• if match ALL hash functions, might be in block (false positive)
• else, definitely not

Indexing

- Local Indexing
– Each node stores index for its data
– Easy to update; still need to scatter-gather
– Have to check all nodes if not partition key
• Global Indexing
– Index the entire DB, partition index with hashing into all the nodes
– Hard to update (need another server); good for searching
– Have to check single node even if not partition key

DynanODB

- item in table has PK, otherwise it's schemaless
– Single PK = (partition key)
– Composite PK = (partition key, sort key)
• Tables are partitioned by partition key
• Items in same partition are sorted by sort key (if given)
• Global Index: index key is simple or composite; partition key can be different from table PK
• Local Index: index key is composite only (sort key is indexed key); partition key must be same as table PK

Distributed Commit Protocols

- ACID
– Atomicity: all or nothing
– Consistency: if each xact is consistent and DB starts consistent, it ends consistent
– Isolation: Executions of xacts are isolated from each other
– Durability: if a xact commits, its effects persist

Centralised DBMS Recovery Manager

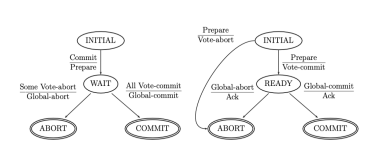
Recovery Manager ensures atomicity and durability.

- Implementing Abort
– Write-Ahead Logging Protocol: flush uncommitted update after before-image is flushed
– Undoes all updates by xact by restoring before-image
• Implementing Commit
– Force-At-Commit Protocol: commit a xact after after-images of all updated pages are flushed (write commit log, immediately flush)
• Implementing Restart
– Redo phase: redo all updates
– Undo phase: abort all active xacts

Failures in DDBMS

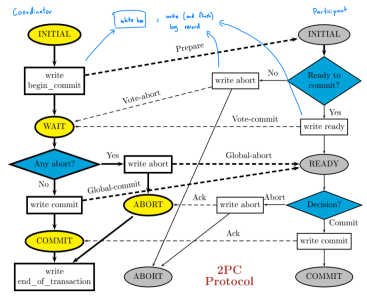
- Site failures
– fail-stop model: a site is either operational OR failed
– Partial site failure: some sites are operational, some have failed
– Total site failure: all sites have failed
• Communication failures (all sites operational)
– lost messages, network partitioning (split-brain problem)

Two Phase Commit (2PC): voting + decision



- (one state transition apart)
• ALL reach same global decision
• once voted, cannot change vote
• one abort, global = abort
• anyone can abort without voting
• commit => all voted commit
• no failure + all vote commit => commit
Log records are flushed / forced / synchronous writes
• Coordinator: Force write commit log record, don't force write abort log record
– (both) Recoverers in WAIT: after timeout, will abort
• Participant must force write ready log record
– Recoverers in INITIAL: will abort (but global might be commit)

- Participant: Force write commit / abort log record (if voted commit), don't need to force write (if voted abort)
– (vote commit): recovers in READY, will revoke commit, coordinator might not be able to inform global decision (because too long later)
– (vote abort): recovers in INITIAL, will abort -> OK!



Site Failures

- detected by timeouts
• Recovery Protocol (by server that failed)
– independent if can terminate without outside info
• Termination Protocol (by TC)
– non-blocking if can terminate without waiting for recovery

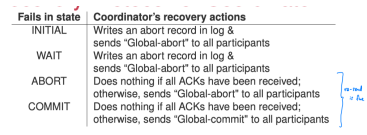


Table with 2 columns: 'Fails in state' and 'Participant's recovery actions'. Rows include INITIAL, WAIT, ABORT, and COMMIT with corresponding actions like 'Aborts transaction unilaterally' and 'Sends "/>

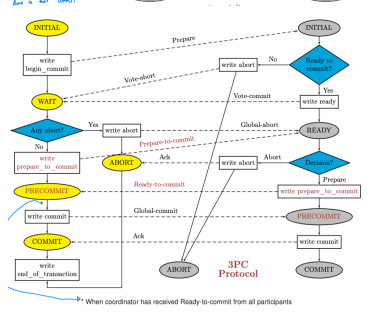
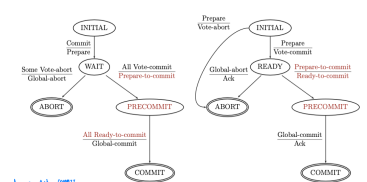
Table with 2 columns: 'Timeout in state' and 'Coordinator's termination actions'. Rows include WAIT, ABORT, and COMMIT with actions like 'Writes an abort record in log & sends "/>

Table with 2 columns: 'Timeout in state' and 'Participant's termination actions'. Rows include INITIAL and READY with actions like 'Aborts transaction unilaterally' and 'Participant is blocked!'.

Cooperative Termination Protocol

- Let participants communicate with each other
• If P timeouts in READY: P asks for decision
– if any node is COMMIT / ABORT: P does that & tells anyone who is READY
– if a node is INITIAL: it aborts, and replies ABORT
– if all are READY: blocking!

Three Phase Commit (3PC)



- Termination Protocol Changes
– Coordinator's pre-commit: write commit log record + send global-commit to operational participants
– Participants' ready + pre-commit: execute Termination Protocol X
• Recovery Protocol is same as 2PC

Termination Protocol 1

- Protocol
– Run leader-election to elect C
– C requests state from participants
* if any in COMMIT: Global-Commit to all
* if none in PRECOMMIT: Global-Abort to all
* else: Send Prepare-to-Commit to READY, receive Ready-to-commit from these, then Global-Commit to all
• 2 steps are needed: otherwise if it crashes again and none in PRECOMMIT, second if is triggered
– if any participant timeout (coordinator failed), elect new coordinator; any participant that fails is ignored; anyone that fails and recovers CANNOT participate
• Total Site Failure:
– Recovering TMs blocked until a TM P recovers:
* P can recover independently (state = INITIAL, ABORT, COMMIT)
• P notifies recovered TMs
* P was last TM to fail
• P executes termination protocol
• Without total site failure + comm. failure: non-blocking
• With total site failure: blocking
• With comm. failure: might be split-brain
Recover from earlier failure cannot rejoin: single node in PRE-

COMMIT, Coordinator & that node fail -> global will be Abort (but if he rejoins, it should be commit)

Termination Protocol 2 (handles comm. failure)

- Protocol
– Run leader-election to elect C
– C requests state from participants
* if any in COMMIT: Global-Commit to all
* if any in ABORT: Global-Abort to all
* if any in PRECOMMIT + no COMMIT / ABORT + majority in READY / PRECOMMIT: Prepare-to-Commit to -PRECOMMIT; receive Ready-to-commit
• if number of PRECOMMIT + Ready-to-commit is majority: Global-Commit
• else: blocked
* if no COMMIT / ABORT + majority in INITIAL / READY / PREABORT: Prepare-to-Abort to -PREABORT; receive Ready-to-abort
• if number of PREABORT + Ready-to-abort is majority: Global-Abort
• else: blocked
• Blocked TMs periodically re-tempt protocol; when failed TM recovers: executes protocol
• Non-blocking as long as majority are operational

Concurrency Control

- Concurrency Control Manager ensures isolation.
• VE if same read-froms & same final-writes
• VSS if VE to some serial schedule
– VSG - (Tj, Ti) if Tj reads from Ti, or both write to same variable and Ti does final-write
– VSG cyclic => not VSS
– VSG acyclic & (serial schedule from topo-sort is VE to S) => VSS
• conflicting actions if
– at least one of them is write action
– and actions are from different transactions
• CE: every pair of conflicting actions are ordered in the same way
• CSS: CE to some serial schedule (CSS => VSS)
– CSS => CSG is acyclic
• blind write: Xact no read before it writes
– VSS & no blind writes => CSS
• Recoverable Schedule (essential): for every Xact T that commits in S, T must commit after T' if T reads from T'

Lock-Based CC

Table with 4 columns: Requested, Held, S, X. Rows represent different lock types and states.

- if lock request not granted, Xact is blocked, Xact is added to O's request queue
• 2PL => CSS; once release a lock, no more request
• Strict 2PL => strict & CSS: Xact must hold onto lock until commit / abort
• Wait-For-Graph: Ti -> Tj if Ti waiting for Tj (must remove edge)
• Timeout mechanism: when Xact start, start timer, if timeout, assume deadlock
• Deadlock Prevention - older Xact has higher priority (not restarted on kill to avoid starvation)
– suppose Ti requests a lock held by Tj (Higher; Lower)
– wait-die (non-preemptive): Ti wait for Tj; Ti suicide => may starve
– wound-wait (preemptive): kill Tj; Ti wait for Tj
– if Tj dies, Ti still waits

MVCC (multiple ver.)

- read-only are never blocked / aborted
• update xacts not blocked by read-only xacts
• MVE if same read-from
• MVSS if MVE to some serial monoversion schedule
– monoversion: each read action returns the most recently created object version
– VSS ⊆ MVSS (not other way round)

Snapshot Isolation (MVCC protocol)

- SI: Xact T takes snapshot of committed state of DB at start of T
– can't read from concurrent Xacts
– Concurrent if overlap start & commits
– Oj is more recent than Oi if Ti commit after Tj
– Concurrent Update Property: if multiple concurrency Xact update same object, only one can commit (if not, may not be serialisable)
• First Committer Win (FCW): check at point of commit
• First Updater Win (FUW) - locks only used for checking (NOT lock-based)
– to update O: request X-lock on O; when commit / abort, release locks
– if not held by anyone:
* if O has been updated by concurrent Xact: abort
* else: grant lock
– else: wait for T' to abort / commit
* if T' commit: abort
* else: use (if not held by anyone) case
• Write Skew Anomaly (not MVSS)
– Both Xact read from initial value
• Read-Only Xact Anomaly (not MVSS)
– A Read-Only Xact reads values that shouldn't be possible
• SSI (Serialisable SI): produced by SI and is MVSS
• Garbage Collection: delete version Oi if starts a newer version Oj st for every active Xact Tj that started after commit of Tj, Tj commits before Tk starts (aka all active Xact can refer to Oj)

Distributed CC

- Global schedule S for T and {S1, ..., Sm} is VSS / CSS if
– Each local Si is VSS / CSS
– and local serialisation orders are compatible
• Centralised 2PL (C2PL)
– all locks are managed by central TM's lock manager
• Distributed 2PL (D2PL)
– each site manages locks for their own stuff

Distributed Deadlock Detection

- Centralised approach
– Each site maintains local Wait-For-Graph

