



CUBRID Architecture Overview

Processes, Layered Stack, Query Pipeline, Distribution

2026-05 · Code Analysis Seminar

Agenda — seven axes

1. **Process model** — `cub_master`, `cub_server`, `cub_pl`, `cub_broker` + `cub_cas`
2. **Layered storage stack** — workspace → locator → catalog → heap/btree/ehash → page buffer + DWB → disk → volumes
3. **Query pipeline** — parser → semantic check → rewrite → optimizer → XASL → executor → scan manager → list-file → cursor
4. **Concurrency / logging / recovery** — MVCC + lock + WAL + checkpoint + recovery + vacuum
5. **Distribution** — heartbeat + HA replication + CDC + 2PC + flashback + backup
6. **PL family** — `cub_pl` JVM, JavaSP, PL/CSQL
7. **Cross-cutting infrastructure** — boot, sessions, thread pools, network, broker, errors, parameters, monitoring

This deck is a **router**, not a deep-dive. Each axis points at 2–5 detail docs in knowledge/code-analysis/cubrid/.

Who this is for · how CUBRID got its shape

For the engineer who has read zero lines of CUBRID source and needs the high-level shape before opening any of the ~70 detail docs. Also for the engineer who has read several detail docs and now wants to fit them together along a single axis.

CUBRID's design is **not generic** — every shape on the next twenty slides is a deliberate choice. The lineage shows:

Choice	Why it looks this way
Object-relational, OODB lineage	UniSQL → CUBRID. Schema is an in-memory object graph (<code>SM_CLASS</code>), not a flat relation; DDL manipulates the graph then persists.
Separate <code>cub_master</code> supervisor	HA story without embedded consensus — every node decides locally from its own peer table.
Broker tier with pre-forked CAS	Pool of long-lived workers + <code>SCM_RIGHTS</code> file-descriptor handoff — TCP termination decoupled from engine state.
PL family in a sibling JVM	Server immune to JVM stalls / GC pauses / user-code crashes; PL is a privileged JDBC client to its own server.
One source tree, three builds	<code>cub_server</code> / <code>libcubridsa</code> / <code>libcubridcs</code> from the same code, switched at runtime via <code>dlopen</code> .

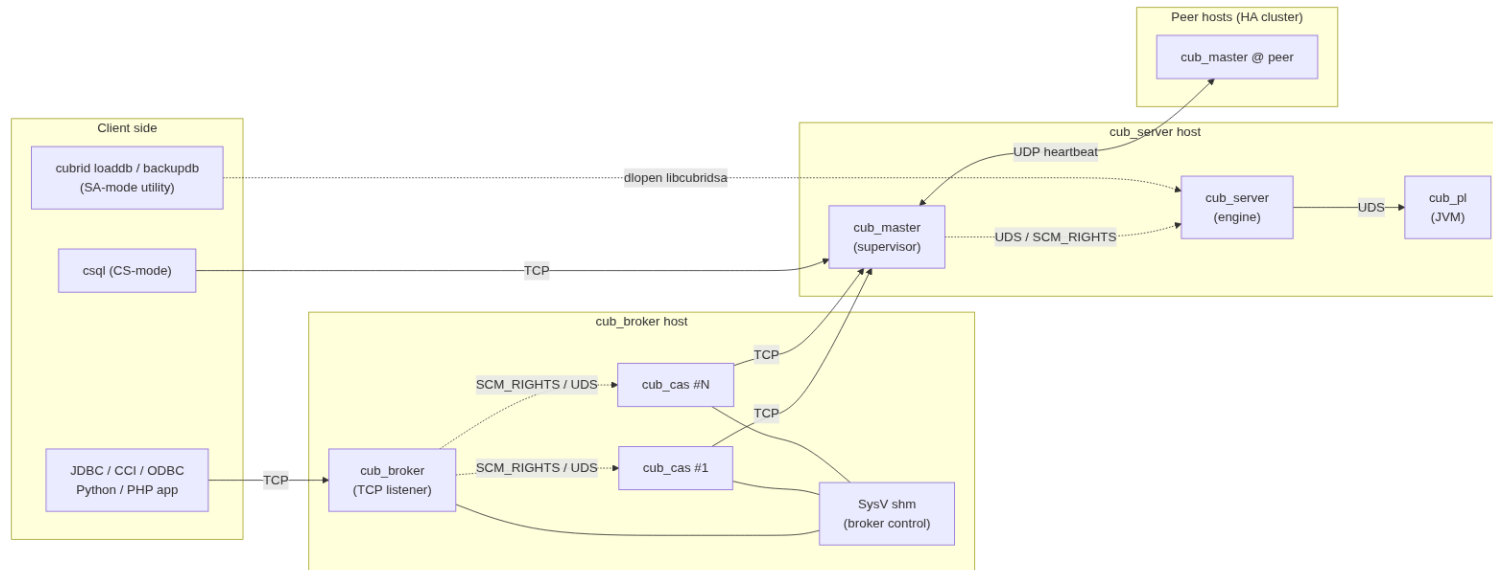
Detail: `cubrid-design-philosophy.md` collects the full rationale.

The four long-lived processes

- `cub_master` — per-host supervisor. Listens on TCP, **hands the file descriptor** to the right `cub_server`. UDP heartbeat between peers. Does **not** own DB state.
- `cub_server` — database engine. One per database. Owns volumes, page buffer, log, lock table, catalog, optimizer, MVCC, vacuum.
- `cub_pl` — PL JVM. One per `cub_server`. Java SPs + PL/CSQL; calls back over per-session UDS.
- `cub_broker` + `cub_cas` — broker tier. Pre-forked CAS pool; **rendezvous-passes** sockets via `SCM_RIGHTS`.

Plus utility processes (`loaddb` , `backupdb` , ...) that `dlopen libcubridsa.so` directly.

Process model — diagram

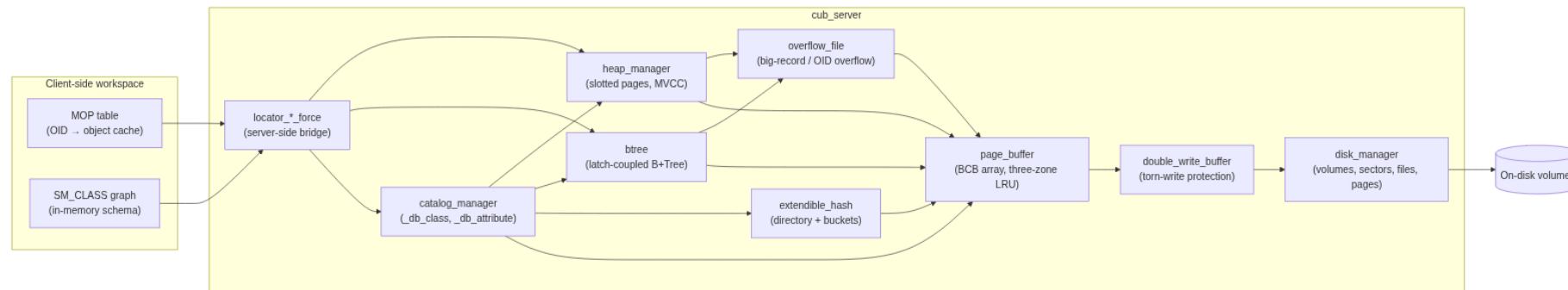


- **TCP listener on `cub_master`, not `cub_server`**. Master routes by DB name, hands over the file descriptor.
- **Broker is a separate tier** with its own TCP hop. `cub_pl` is a sibling JVM to `cub_server`.

Part II

Inside `cub_server`

Layered storage stack — diagram

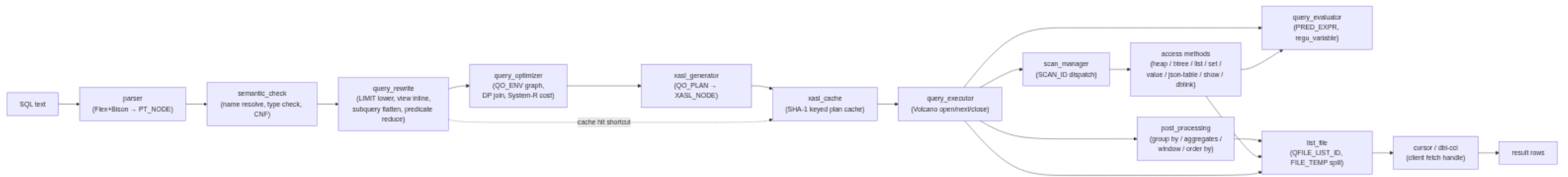


Rule: every record layer talks to the page buffer, never to the disk manager. That is what lets the buffer enforce WAL ordering.

Each storage layer in one sentence

Layer	One sentence	Detail doc
<code>disk_manager</code>	Volumes, sectors (64 pages), files, pages (I/O unit). Permanent / temporary cache split.	<code>cubrid-disk-manager.md</code>
<code>page_buffer</code>	<code>(VPID → BCB → frame)</code> hash, three-zone LRU, per-BCB read/write/flush latch.	<code>cubrid-page-buffer-manager.md</code>
<code>double_write_buffer</code>	Sequential staging volume <code>fsync</code> 'd before home-page write — torn-write protection.	<code>cubrid-double-write-buffer.md</code>
<code>heap_manager</code>	Slotted pages, MVCC headers, nine record types, big-record overflow spill.	<code>cubrid-heap-manager.md</code>
<code>btree</code> / <code>extendible_hash</code>	Latch-coupled B+Tree (<code>key OID</code>); Fagin-style ehash for class-name / repr-id lookup.	<code>cubrid-btree.md</code> · <code>cubrid-extendible-hash.md</code>
<code>catalog_manager</code>	Per-class disk repr + statistics (CTID); <code>_db_class</code> / <code>_db_attribute</code> system classes.	<code>cubrid-catalog-manager.md</code>
<code>SM_CLASS</code> / <code>locator</code>	In-memory schema graph + client → server bridge (<code>locator_*_force</code>).	<code>cubrid-class-object.md</code> · <code>cubrid-locator.md</code>

Query pipeline — diagram

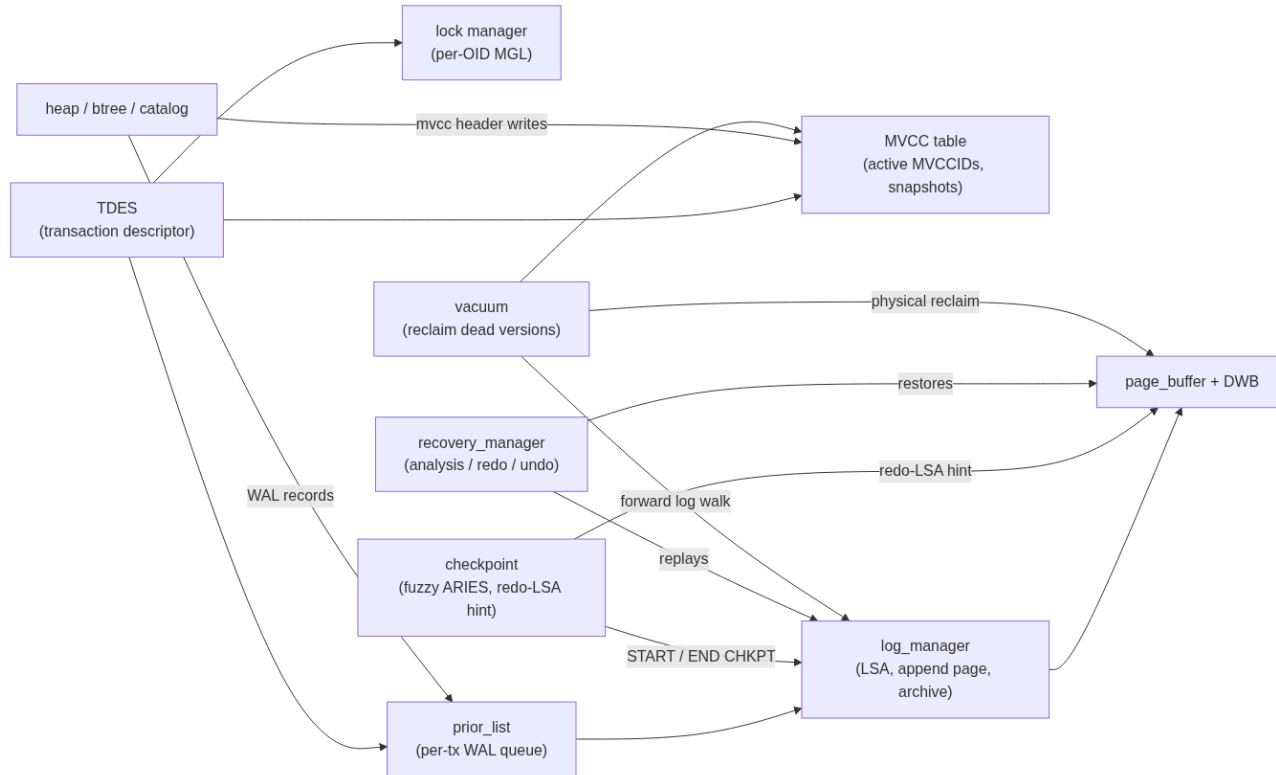


- **Compile once, execute many.** Stages 1–5 produce a serialised XASL tree cached on a SHA-1 of the rewritten SQL.
- **Optimizer IR \neq executor IR.** `QO_PLAN` (graph + cost) is lowered into `XASL_NODE` (recursive tree with `aptr/dptr/scan_ptr`).
- **Executor is Volcano-style.** Uniform open / next / close over polymorphic `SCAN_ID` operators.

Pipeline stages in one sentence

Stage	One sentence	Detail doc
parser	Flex/Bison → polymorphic <code>PT_NODE</code> tree, per- <code>PARSER_CONTEXT</code> block allocator.	<code>cubrid-parser.md</code>
semantic check	Name resolve, aggregate, host-variable, local; CNF for predicates.	<code>cubrid-semantic-check.md</code>
rewrite	LIMIT lower, view inline, subquery flatten, predicate reduce, auto-parameterize.	<code>cubrid-query-rewrite.md</code>
optimizer	<code>QO_ENV</code> graph, DP join enum, System-R cost.	<code>cubrid-query-optimizer.md</code>
xasl_generator	<code>QO_PLAN</code> → <code>XASL_NODE</code> tree with <code>REGU_VARIABLE</code> / <code>ACCESS_SPEC</code> sub-IRs.	<code>cubrid-xasl-generator.md</code>
query_executor	Volcano XASL interpreter; plan cached on SHA-1 of rewritten SQL.	<code>cubrid-query-executor.md</code> · <code>cubrid-xasl-cache.md</code>
scan_manager	Polymorphic <code>SCAN_ID</code> ; heap / btree / list / set / value / json-table / dblink / show.	<code>cubrid-scan-manager.md</code>
post-processing / cursor	GROUP BY, analytics, external sort → <code>QFILE_LIST_ID</code> ; client <code>CURSOR_ID</code> .	<code>cubrid-post-processing.md</code> · <code>cubrid-cursor.md</code>

Concurrency, logging, recovery — diagram

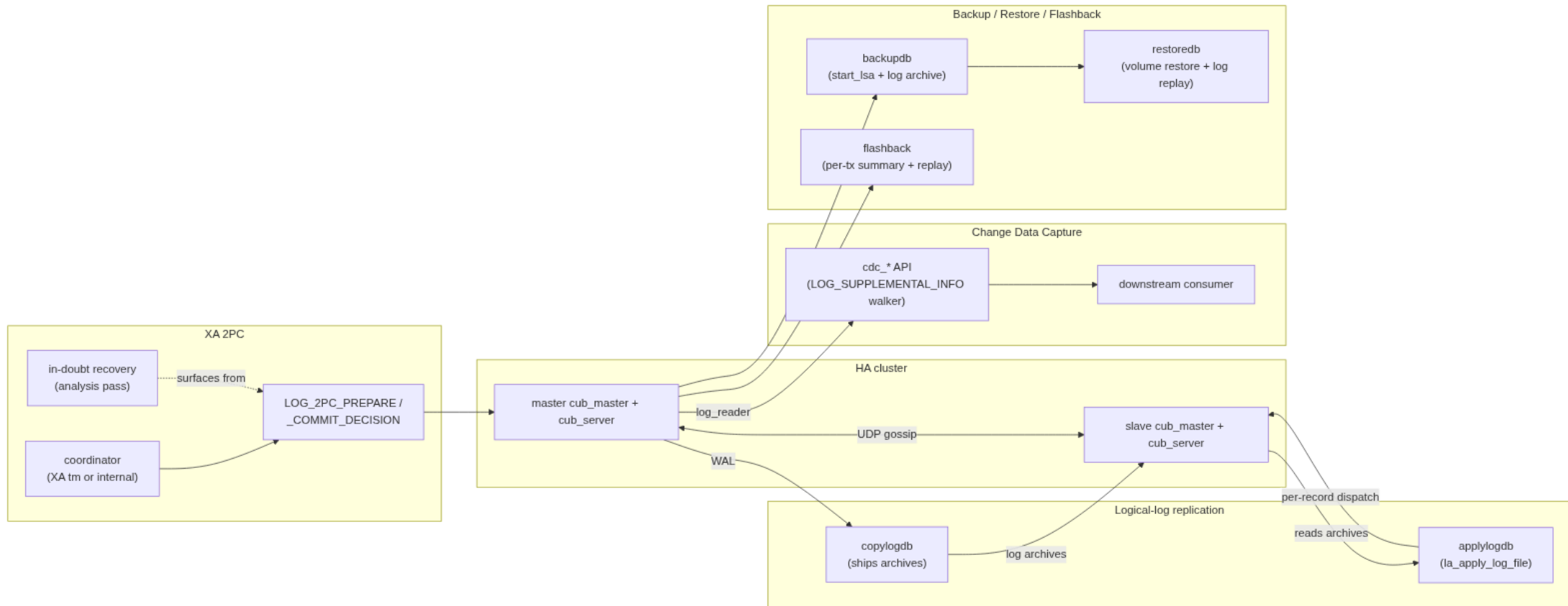


- Three timelines co-exist — **transactional** (MVCCIDs + locks), **physical** (WAL + LSAs), **page** (dirty BCBs).
- WAL is split via a per-tx **prior list**; group commit emerges from queue batching.

Each subsystem in one sentence

Subsystem	One sentence	Detail doc
transaction	Per-tx <code>TDES</code> ; isolation-level dispatch (SI / lock-based); nested savepoint rollback.	<code>cubrid-transaction.md</code>
MVCC	MVCCID assign, per-tx snapshot via <code>build_mvcc_info</code> ; read-only is transactionless .	<code>cubrid-mvcc.md</code>
lock_manager	Multi-granularity per-OID grant/convert/revoke, WFG cycle-scan deadlock detector.	<code>cubrid-lock-manager.md</code>
log_manager + prior_list	Per-tx prior-list queue, daemon drain → append-page pipeline → archive volumes.	<code>cubrid-log-manager.md</code> · <code>cubrid-prior-list.md</code>
checkpoint	Fuzzy-ARIES daemon; emits redo-LSA hint without forcing dirty pages.	<code>cubrid-checkpoint.md</code>
recovery_manager	Three-pass restart from <code>chkpt_lsa</code> ; per-page parallel redo.	<code>cubrid-recovery-manager.md</code>
vacuum	Forward WAL walk below oldest-visible MVCCID; master → worker dispatch.	<code>cubrid-vacuum.md</code>
double_write_buffer	Sequential staging volume <code>fsync</code> 'd before home write — torn-write protection.	<code>cubrid-double-write-buffer.md</code>

Distribution — diagram

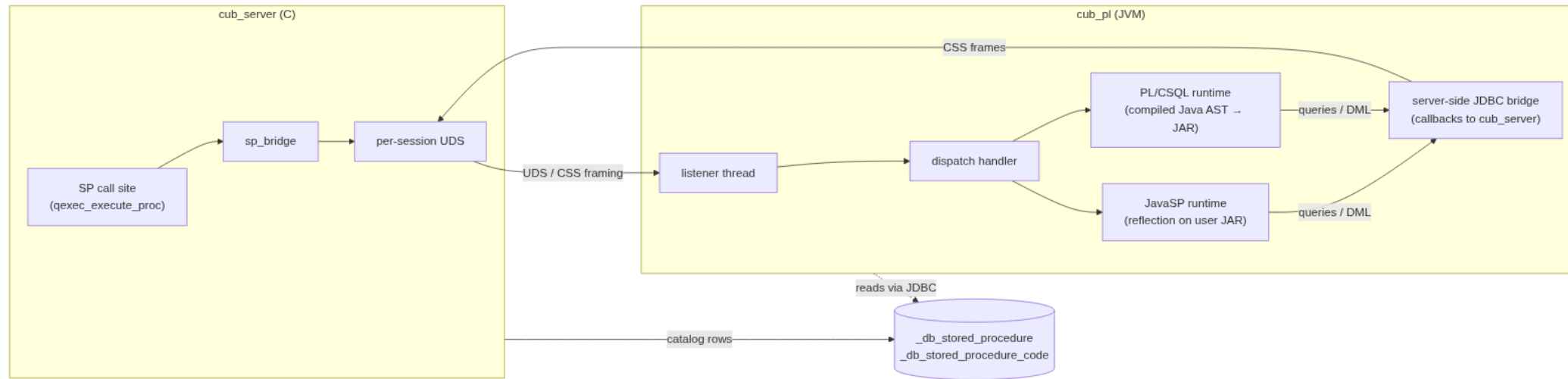


- **cub_master** is the gossip endpoint — heartbeat is independent per node; witness hosts guard split-brain.
- The WAL is CUBRID's single event log — recovery, vacuum, replication, CDC, flashback, backup all read the same stream; they differ only in **which record types** and **which direction**.

Distribution subsystems

Subsystem	One sentence	Detail doc
heartbeat	UDP gossip cluster liveness, per-node independent <code>calc_score</code> master election, slave → to-be-master → master FSM, witness-host (<code>ha_ping_hosts</code>) split-brain guard.	<code>cubrid-heartbeat.md</code>
HA replication	Master emits <code>LOG_REPLICATION_DATA</code> / <code>_STATEMENT</code> alongside physiological WAL; <code>copylogdb</code> ships archives; <code>applylogdb</code> walks them forward, per-record-type dispatch back into the storage layer.	<code>cubrid-ha-replication.md</code>
CDC	<code>cdc_*</code> API walks <code>LOG_SUPPLEMENTAL_INFO</code> records via <code>log_reader</code> ; legacy <code>la_*</code> HA applier shares the same record format internally.	<code>cubrid-cdc.md</code>
2PC	Coordinator + participant FSMs through <code>LOG_2PC_EXECUTE</code> ; prepared-state log records survive crash; in-doubt transactions surface during ARIES analysis pass.	<code>cubrid-2pc.md</code>
flashback	Two-phase forward log walk — per-tx summary then per-tx detailed log-info pull — shares CDC entry format, reads archived log volumes.	<code>cubrid-flashback.md</code>
backup / restore	Online physical backup — snapshot data volumes + log records bracketed by <code>start_lsa</code> and the next checkpoint; restore replays the log forward up to a stop time.	<code>cubrid-backup-restore.md</code>

PL family — `cub_pl` as a sibling JVM



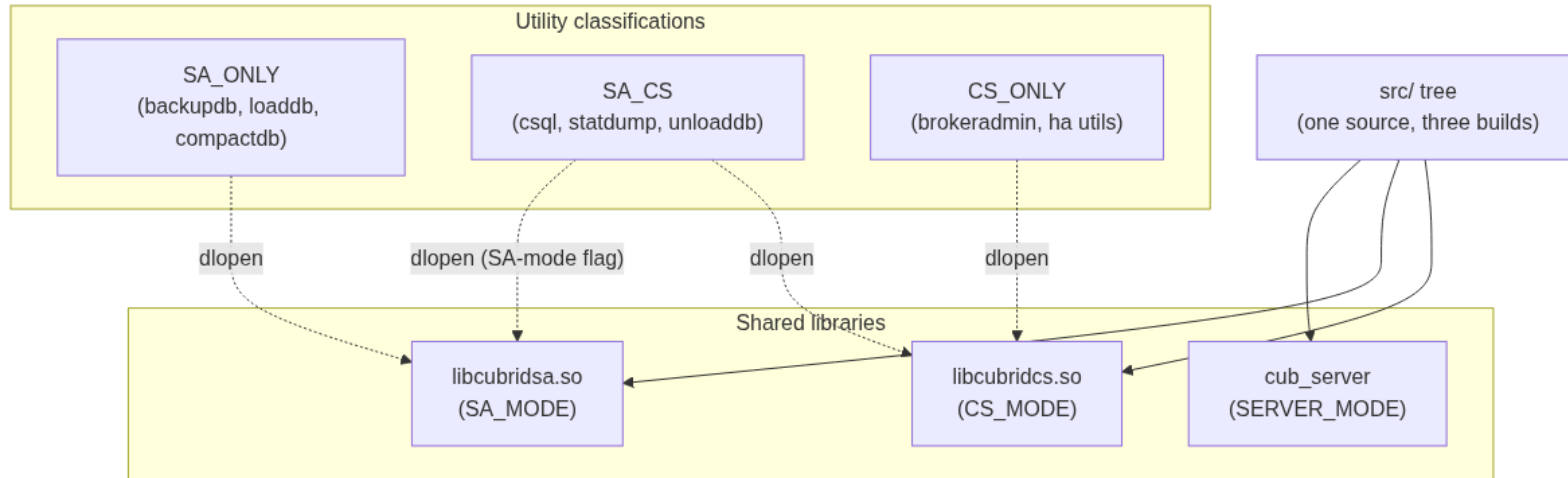
- **Two languages, one runtime.** Java SPs and PL/CSQL both execute in `cub_pl`. PL/CSQL is parsed by ANTLR 4 inside `pl_server`, lowered to a CUBRID-specific Java AST, emitted as Java source, compiled by `javax.tools.JavaCompiler`, packaged as a Base64 JAR.
- **JVM calls back via JDBC.** From the engine's point of view, PL is a privileged client — same CSS framing, same NRP dispatch, same prepared-statement cache.
- **Why a sibling process?** Server is immune to JVM stalls, GC pauses, user-code crashes.

Detail: [cubrid-pl-javasp.md](#), [cubrid-pl-plcsql.md](#), [cubrid-pl-server-bridge.md](#).

Cross-cutting infrastructure

Subsystem	One sentence	Detail doc
boot	Ordered subsystem init; <code>createdb</code> formats volumes + bootstraps catalog; restart hands off to <code>log_recovery</code> 's three-pass replay.	<code>cubrid-boot.md</code>
server session	<code>SESSION_STATE</code> lock-free hash by session id, cached on the connection entry, bound to per-thread TDES.	<code>cubrid-server-session.md</code>
thread + worker pool	<code>cubthread::entry</code> , <code>worker_pool</code> (cores → workers → task queue), <code>daemon</code> + <code>looper</code> pattern, <code>csect</code> RW primitive.	<code>cubrid-thread-worker-pool.md</code>
thread manager NG	CBRD-26177 — bounded epoll workers, coordinator-driven rebalancing, send/rcv budgets.	<code>cubrid-thread-manager-ng.md</code>
network protocol	<code>NET_SERVER_*</code> opcode table; <code>cub_master</code> hands off via <code>master::connector</code> ; epoll CSS framing; <code>or_pack_*</code> marshalling.	<code>cubrid-network-protocol.md</code>
broker	<code>cub_broker</code> forks fixed <code>cub_cas</code> pool; <code>SCM_RIGHTS</code> UDS rendezvous; SysV shared-memory control plane.	<code>cubrid-broker.md</code>
system parameters	<code>prm_Def[]</code> registry, <code>cubrid.conf</code> INI, env overrides, <code>db_set_system_parameters</code> SQL path, <code>SESSION_PARAM</code> .	<code>cubrid-system-parameters.md</code>
monitoring	<code>cubmonitor</code> + legacy <code>perf_monitor</code> / <code>pstat_Metadata</code> for <code>SHOW STATS</code> / <code>statdump</code> .	<code>cubrid-monitoring.md</code>

The three SA / CS / utility build variants



- **Why this exists.** Admin ops (`backupdb` , `loaddb` , `compactdb`) need to touch disk with no running `cub_server` — SA mode `dlopen` s the entire engine in-process.
- **csql is both.** `--CS-mode` links `libcubridcs.so` ; offline, links `libcubridsa.so` . A single flag flips the dispatch.
- Detail: `cubrid-sa-cs-runtime.md` , `cubrid-csql.md` .

Subcategory map — where to start

Subcategory	One-line description	Entry-point docs
server-architecture	Process-level shape — boot, broker, sessions, threads, network.	cubrid-boot.md · cubrid-broker.md · cubrid-network-protocol.md
storage-engine	The layered storage stack inside <code>cub_server</code> .	cubrid-disk-manager.md · cubrid-page-buffer-manager.md · cubrid-heap-manager.md · cubrid-btree.md
base-infra	Custom allocators + lock-free primitives that every layer composes with.	cubrid-overview-base-infra.md · cubrid-lockfree-overview.md · cubrid-private-allocator.md
query-processing	Parse → execute → return — the full SQL pipeline.	cubrid-parser.md · cubrid-query-optimizer.md · cubrid-query-executor.md · cubrid-scan-manager.md
txn-recovery	Concurrency, logging, ARIES recovery, vacuum.	cubrid-transaction.md · cubrid-mvcc.md · cubrid-log-manager.md · cubrid-recovery-manager.md
ddl-schema	Catalog, schema object graph, authorization, statistics.	cubrid-catalog-manager.md · cubrid-class-object.md · cubrid-ddl-execution.md · cubrid-authentication.md
replication-ha	Distribution, log streaming, change capture, flashback.	cubrid-heartbeat.md · cubrid-ha-replication.md · cubrid-cdc.md · cubrid-flashback.md
pl-language	Procedural extensions in the sibling JVM.	cubrid-pl-javasp.md · cubrid-pl-plcsql.md · cubrid-pl-server-bridge.md
i18n-specialty	Charset, collation, timezone, json-table, show.	cubrid-charset-collation.md · cubrid-json-table.md

Reading paths — three guided routes

Three vertical slices across subcategories. Pick the one that matches what you're debugging.

(a) Single query — "SELECT from broker to heap and back"

`cubrid-rpath-select.md` → parser → semantic-check → rewrite → optimizer → xasl-generator → query-executor → scan-manager → list-file → cursor

(b) Transaction commit — "what makes COMMIT durable"

`cubrid-rpath-write.md` → locator → heap-manager → mvcc → lock-manager → prior-list → log-manager → DWB → page-buffer → checkpoint → recovery-manager

(c) HA failover — "`cub_master` declares a new primary"

`cubrid-heartbeat.md` → master-process → ha-replication → cdc → 2pc → flashback → backup-restore



CUBRID™

Thank you

Q & A

- **Analysis (front door):** <knowledge/code-analysis/cubrid/cubrid-architecture-overview.md>
- **Detail docs:** ~70 files under <knowledge/code-analysis/cubrid/> across **eight subcategories**
- **Per-axis overviews (3 of 8):** <cubrid-overview-storage-engine.md> · <cubrid-overview-query-processing.md> · <cubrid-overview-txn-recovery.md> — and five more
- **Design rationale:** <cubrid-design-philosophy.md>