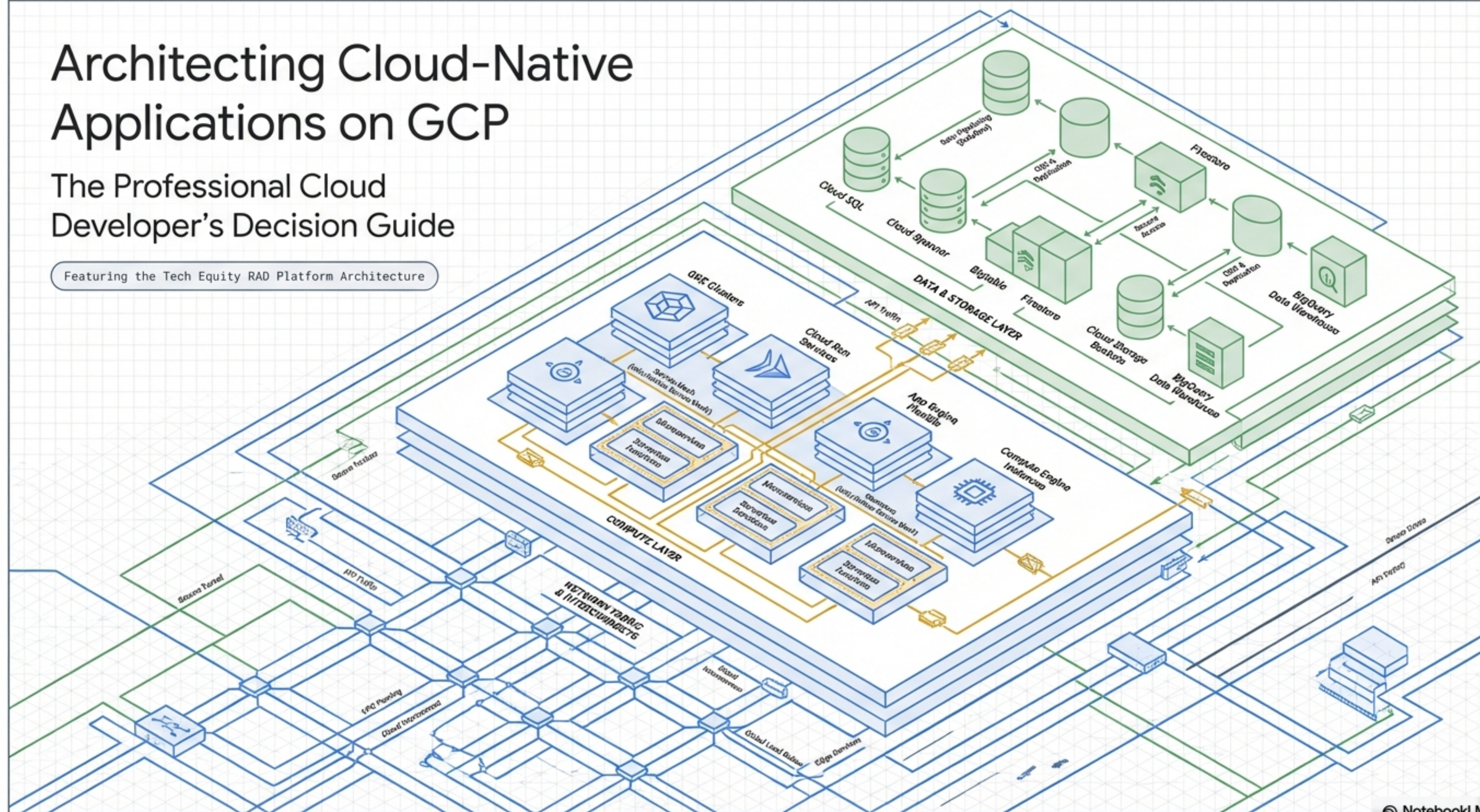


Architecting Cloud-Native Applications on GCP

The Professional Cloud Developer's Decision Guide

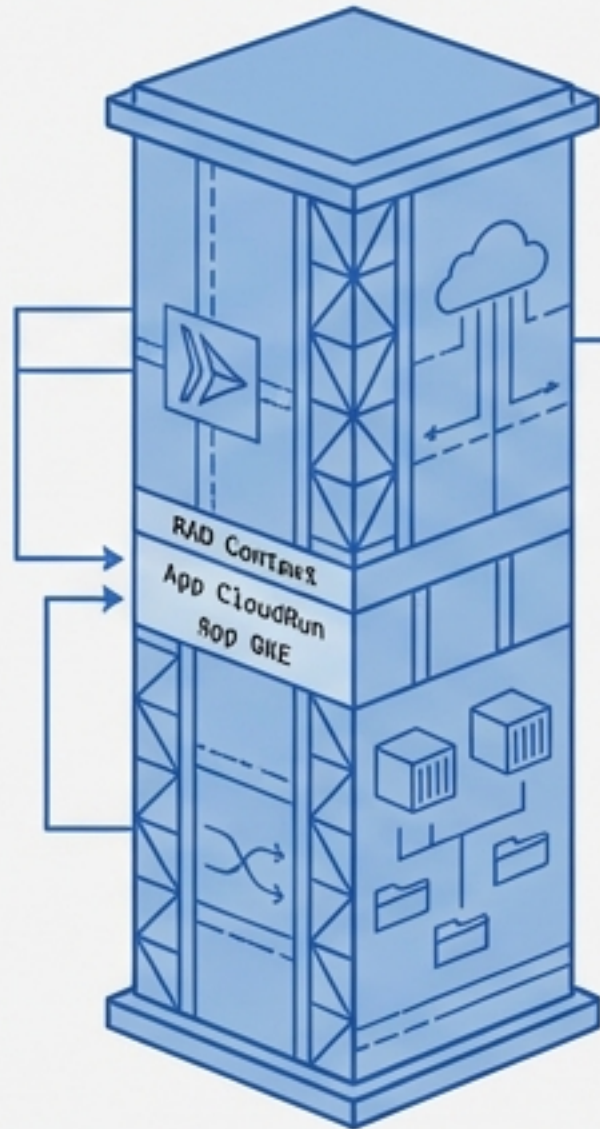
Featuring the Tech Equity RAD Platform Architecture



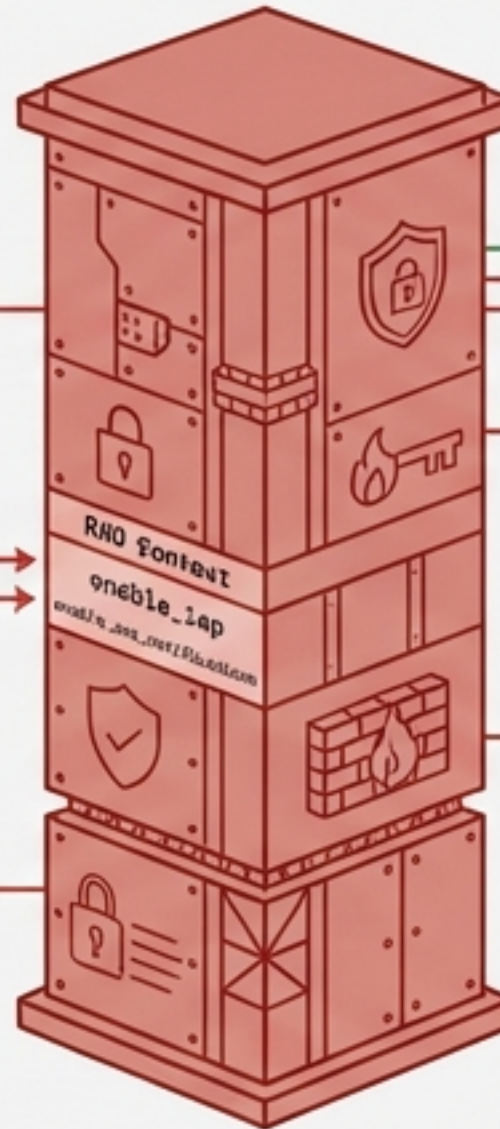
The Core Trinity of Cloud Architecture

Building cloud-native applications requires mastering three intersecting domains. This guide uses the Tech Equity RAD platform to map textbook theory to functional deployment variables.

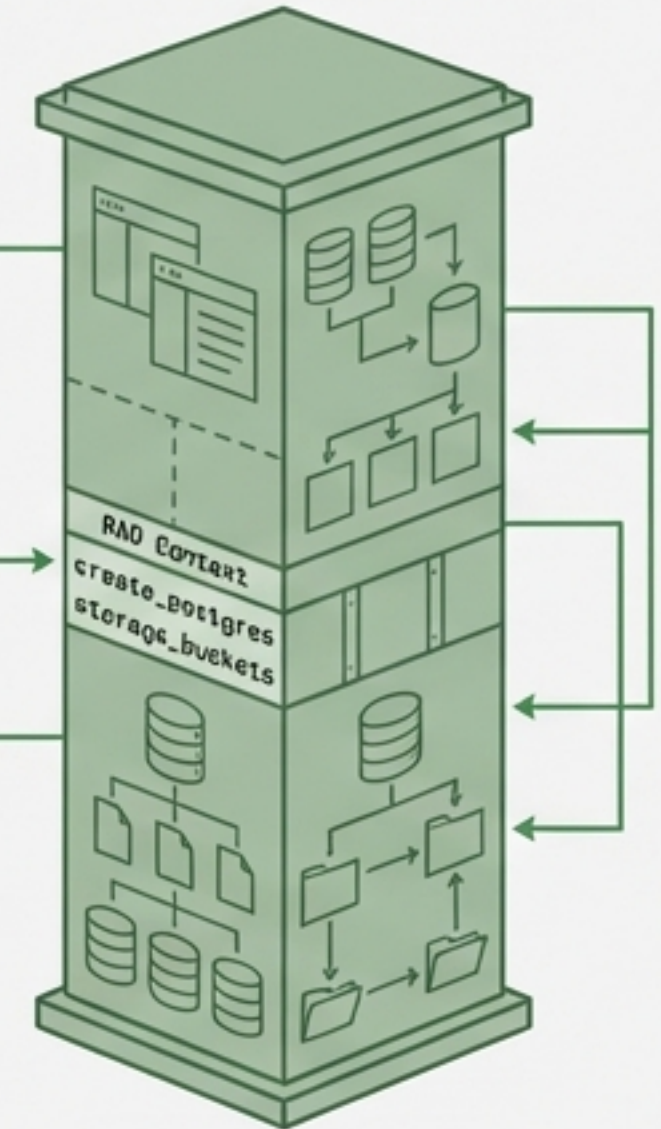
Performance & Compute
Routing traffic, scaling to zero, and orchestrating microservices.



Security & Access
Zero-trust identity, supply chain defense, and compliance.



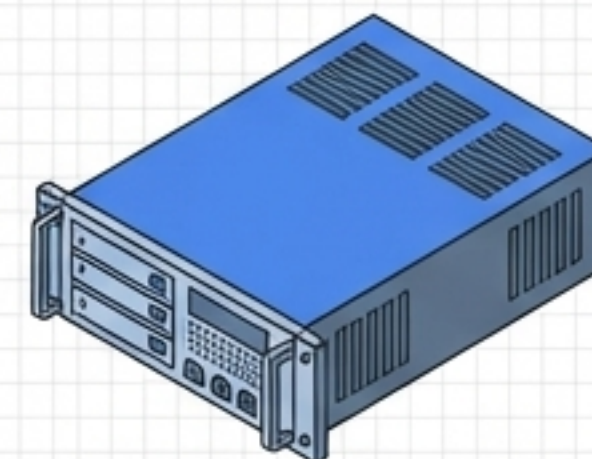
Data & State
Consistency models, schema design, and high-throughput analytics.



Choosing the Compute Foundation

Workload Requirements

Requires Custom OS / Kernel access / Legacy hardware?



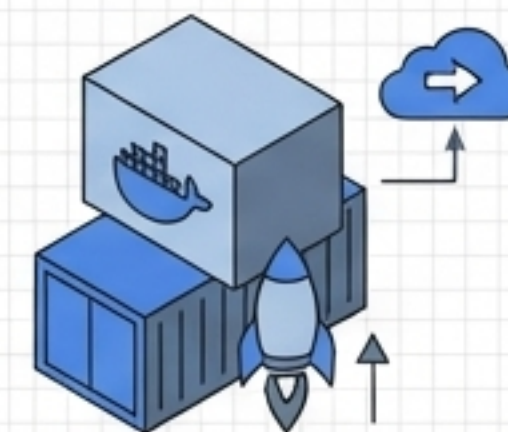
Compute Engine
Provisioned via Services GCP

Requires persistent connections, StatefulSets, or fine-grained networking?



GKE Autopilot
Provisioned via App GKE

Event-driven HTTP workload?
Needs to scale to zero?

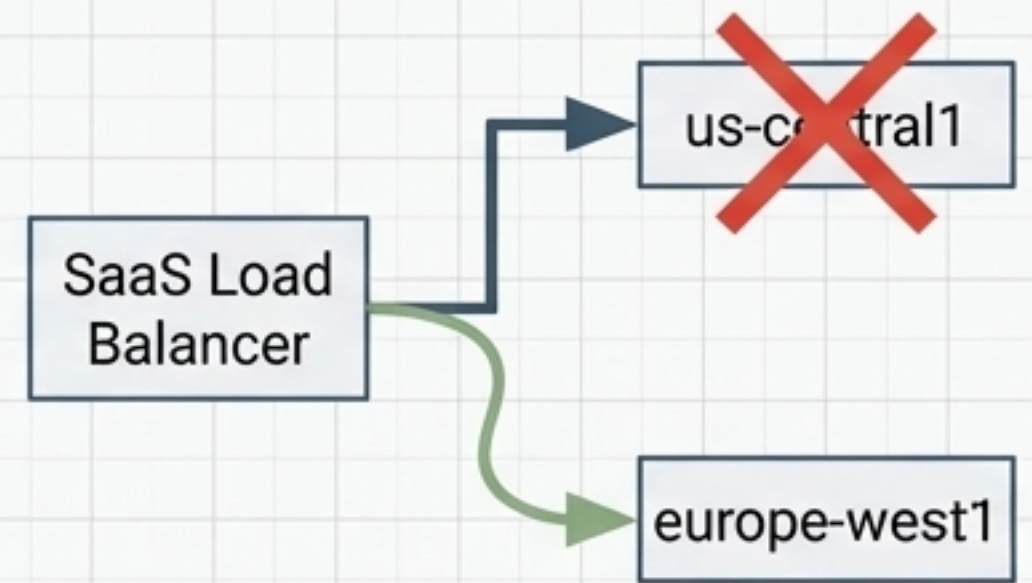
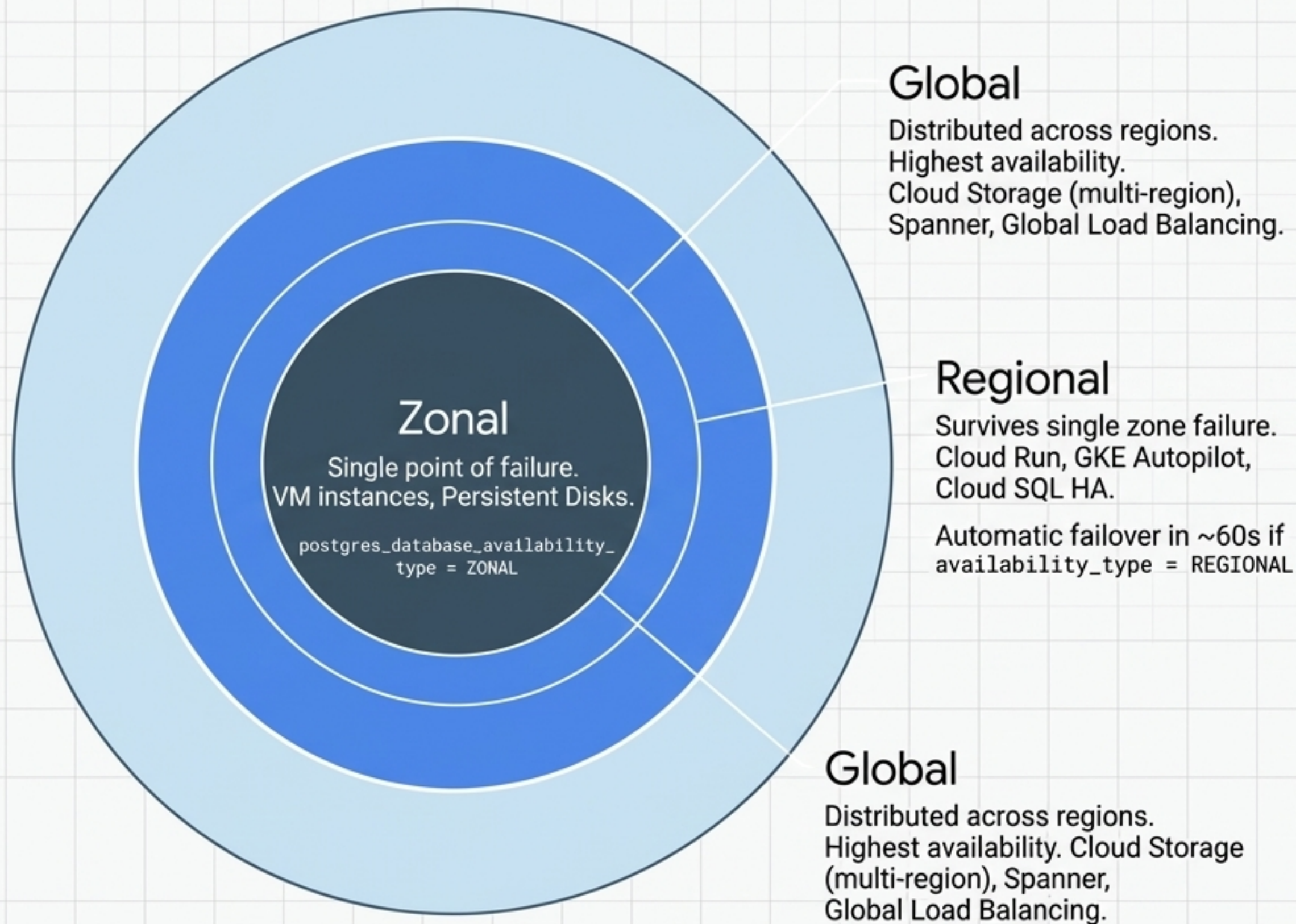


Cloud Run
Provisioned via App CloudRun

Cost Optimization

Right-sizing resources dictates cost. Set `min_instance_count = 0` to incur zero compute cost when idle. Use `max_instance_count` to bound scaling costs during spikes, and `container_resources` to prevent CPU/memory over-provisioning.

Availability and the Blast Radius

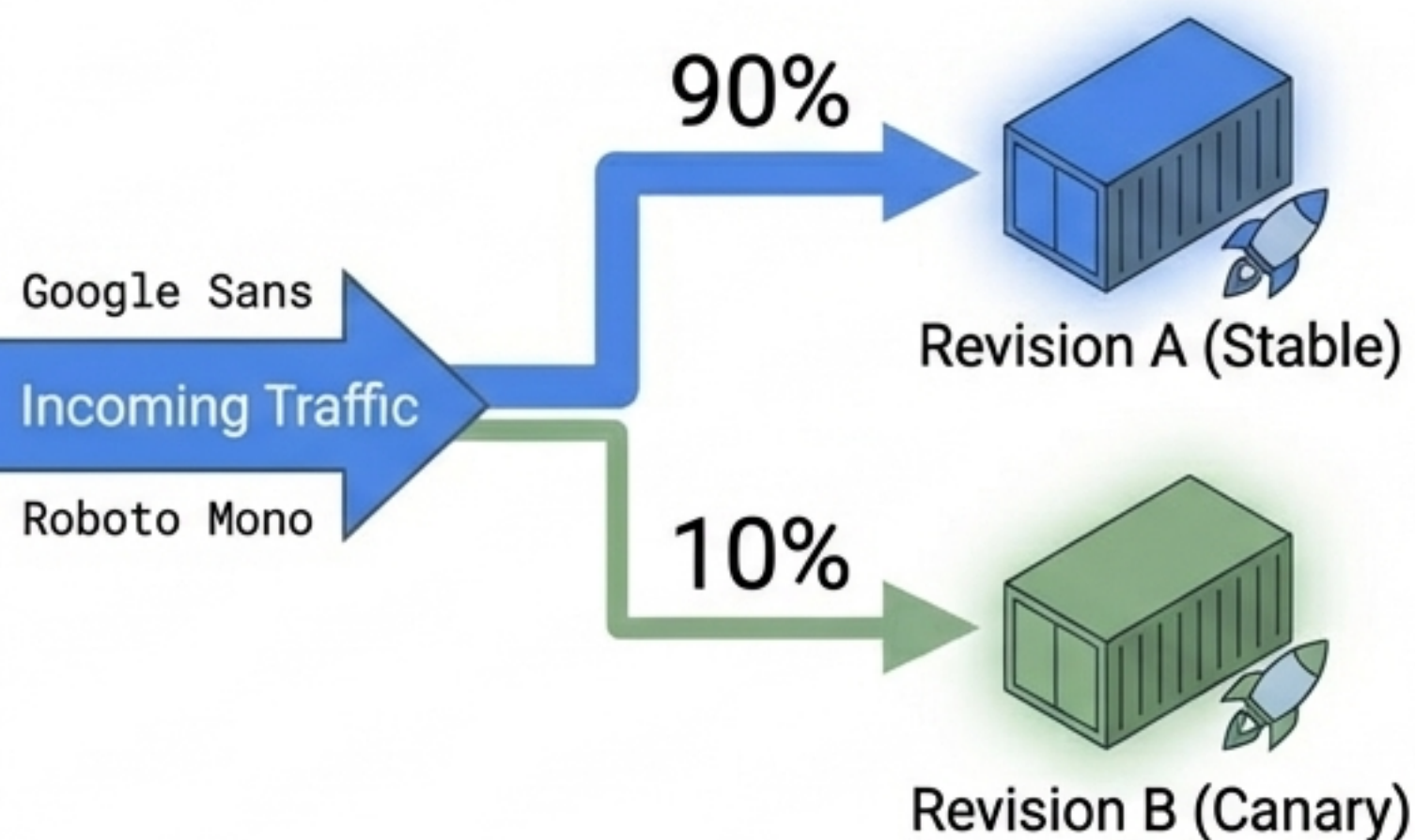


Failover Routing

Health checks detect regional degradation instantly. Traffic shifts to the nearest healthy region with zero application downtime, controlled by the `availability_regions` variable.

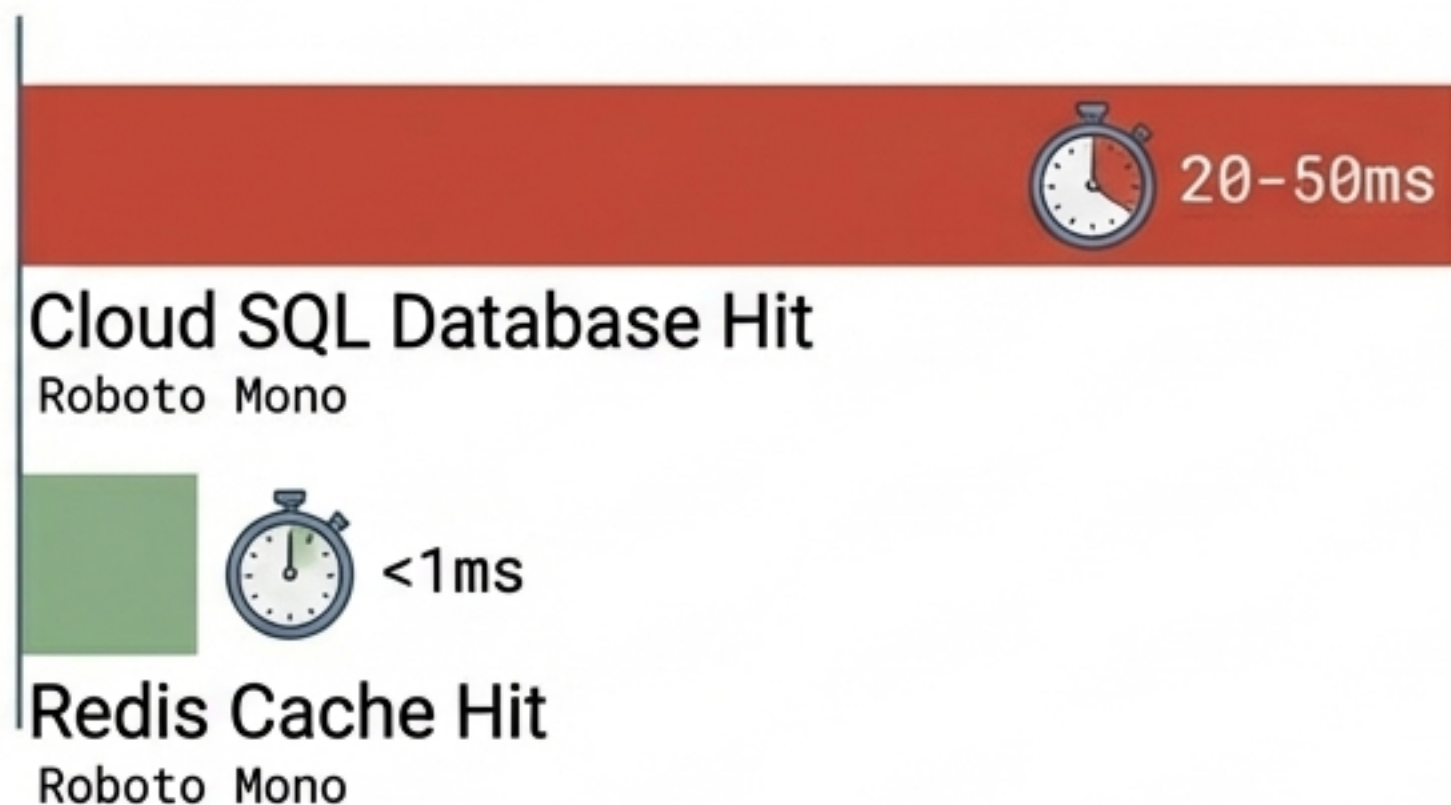
Traffic Rollouts and Microsecond Caching

Immutable Revisions & Splitting



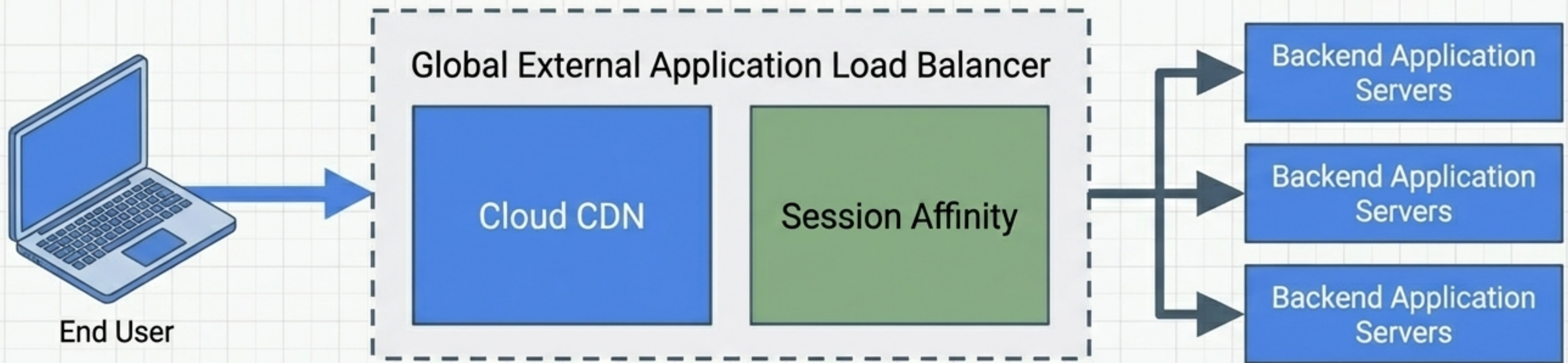
Cloud Run revisions are immutable. Use the `traffic_split` variable for A/B testing. If the 10% canary algorithm underperforms, rollback is instant by routing 100% back to stable.

Memorystore Redis Caching



Enabled via `create_redis` and `redis_tier`. Passing the private Redis IP directly into the container drops cache-hit latency under 1ms, skipping Cloud SQL entirely and reducing DB CPU load by 60%.

Edge Routing and State Preservation



Cloud CDN Caching Modes

Caches at global Points of Presence (PoPs).

- `CACHE_ALL_STATIC`: Auto-caches images/CSS.
- `USE_ORIGIN_HEADERS`: Respects app Cache-Control.
- `FORCE_CACHE_ALL`: Overrides headers (use with caution).

Session Affinity Options

Critical when legacy apps store state in local memory.

- `CLIENT_IP`: Simple, but breaks behind NAT/proxies.
- `GENERATED_COOKIE`: Load balancer sets a cookie on the first response. Recommended for browser apps.
- `HEADER_FIELD`: Routes based on specific HTTP header values.

API Architecture & Management

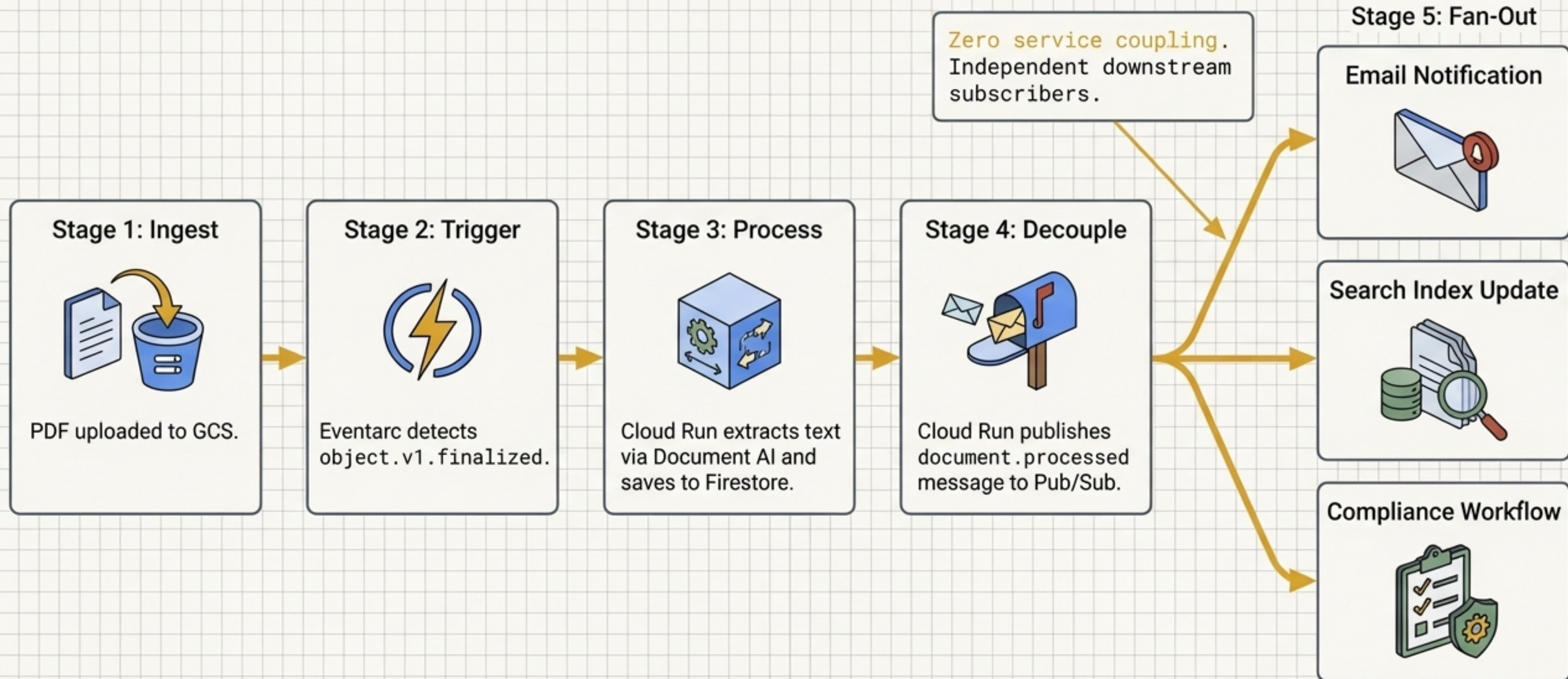
Protocol Choice	
HTTP REST	gRPC
Standard JSON payloads and standard HTTP methods. Native to Cloud Run. Best for browser clients.	Protocol Buffers and HTTP/2 transport. Highly efficient serialization. Best for service-to-service communication.

Gateway Matrix	
Cloud API Gateway	Apigee
Lightweight, low-latency managed gateway. Defined via OpenAPI 2.0 specifications. Enforces basic rate limits (e.g., 100 req/min) and JWT auth. Best for simple API surfaces.	Enterprise API management platform. Provides CIAM, monetization, and OAuth 2.0 flows. Advanced traffic analytics and transformation policies. Best for complex, multi-tier external developer ecosystems.

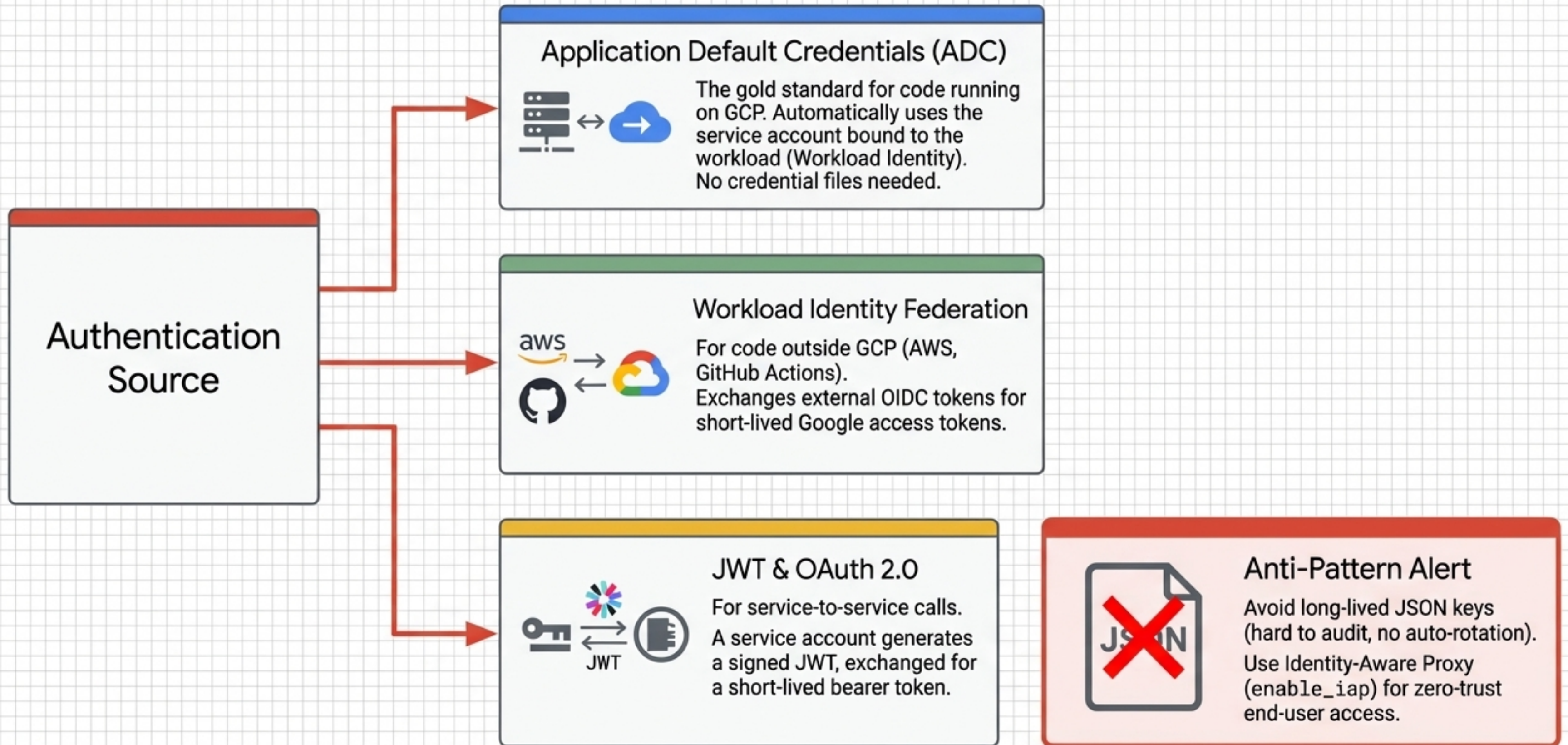
The Orchestration Matrix

Cloud Workflows	Cloud Tasks	Cloud Scheduler	Eventarc
<p>Primary Use Case: Multi-step, sequential coordination.</p> <p>Statefulness: Stateful.</p> <p>Key Capability: Handles branching logic, retries, and sequential HTTP service calls (e.g., Inventory → Payment → Shipping).</p>	<p>Primary Use Case: Rate-controlled asynchronous execution.</p> <p>Statefulness: Stateless.</p> <p>Key Capability: Dispatches single tasks to targets (e.g., sending at most 50 emails/minute). Supports deduplication and future scheduling.</p>	<p>Primary Use Case: Fully managed cron service.</p> <p>Statefulness: Stateless.</p> <p>Key Capability: Time-based invocation of endpoints or topics (e.g., nightly batch jobs).</p>	<p>Primary Use Case: Event routing.</p> <p>Statefulness: Stateless.</p> <p>Key Capability: Routes triggers via CloudEvents (Storage finalized, Audit Logs, Pub/Sub) directly to compute targets. The glue for event-driven systems.</p>

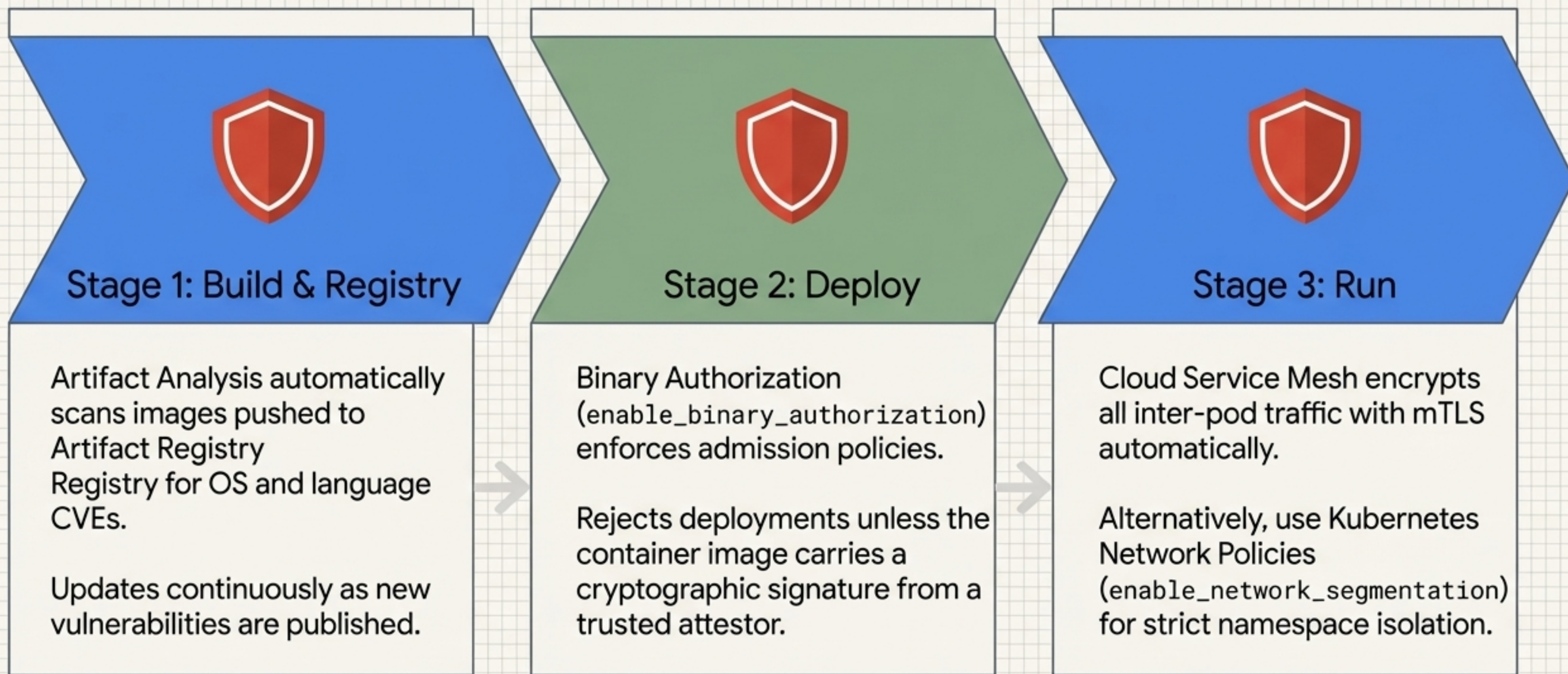
Synthesis: Unified Event-Driven Architecture



Securing Foundational Access



Securing the Supply Chain & Microservices



Data Compliance and Retention

Crypto-Shredding with CMEK

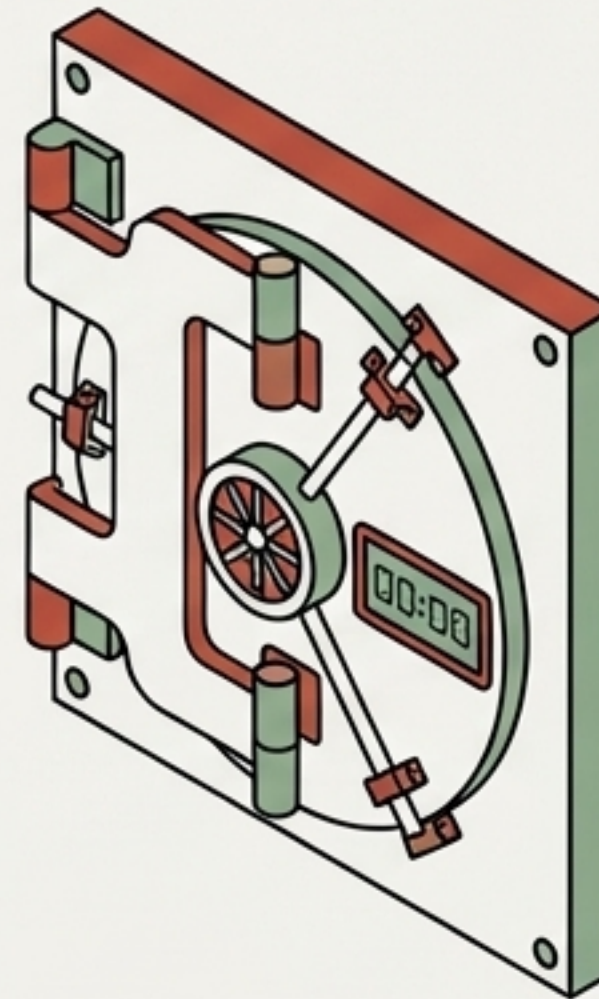


Customer-Managed Encryption Keys allow instant data revocation.

Encrypt tenant data with a unique KMS key.

To comply with deletion requests, simply destroy the KMS key—rendering all associated data permanently inaccessible in minutes.

WORM Storage with Bucket Lock



For regulatory compliance (**SEC 17a-4, HIPAA**), use Cloud Storage Object Lifecycle Management combined with Bucket Lock.

Enforces a permanent retention period where objects cannot be deleted or overwritten by anyone—even project owners—until expiry.

The Master Database Selection Matrix

Structure: Relational vs. NoSQL

Workload: Transactional (OLTP) vs.
Analytical (High-Throughput)



Cloud SQL

Standard web apps (`create_postgres`).
Single-primary, strong consistency. Max 64TB.



Cloud Spanner

Limitless horizontal scaling. Global reads/writes
with TrueTime external consistency.



Firestore

Serverless, flat document maps. Strong
single-document consistency. Ideal for real-time
app listeners.



AlloyDB

HTAP processing. Disaggregated storage.
4x OLTP and 100x analytical throughput vs
standard PostgreSQL.

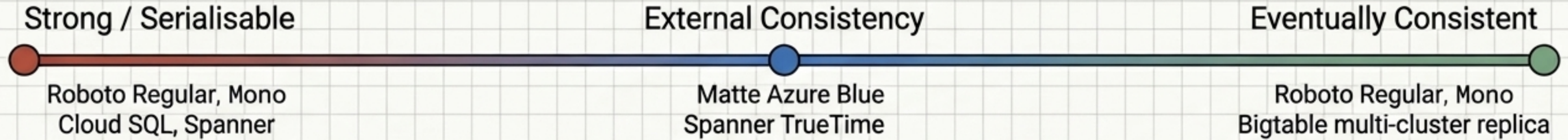


Bigtable

Wide-column. Petabyte scale, single-digit ms
latency. Best for IoT and time-series. The row
key is the only index.

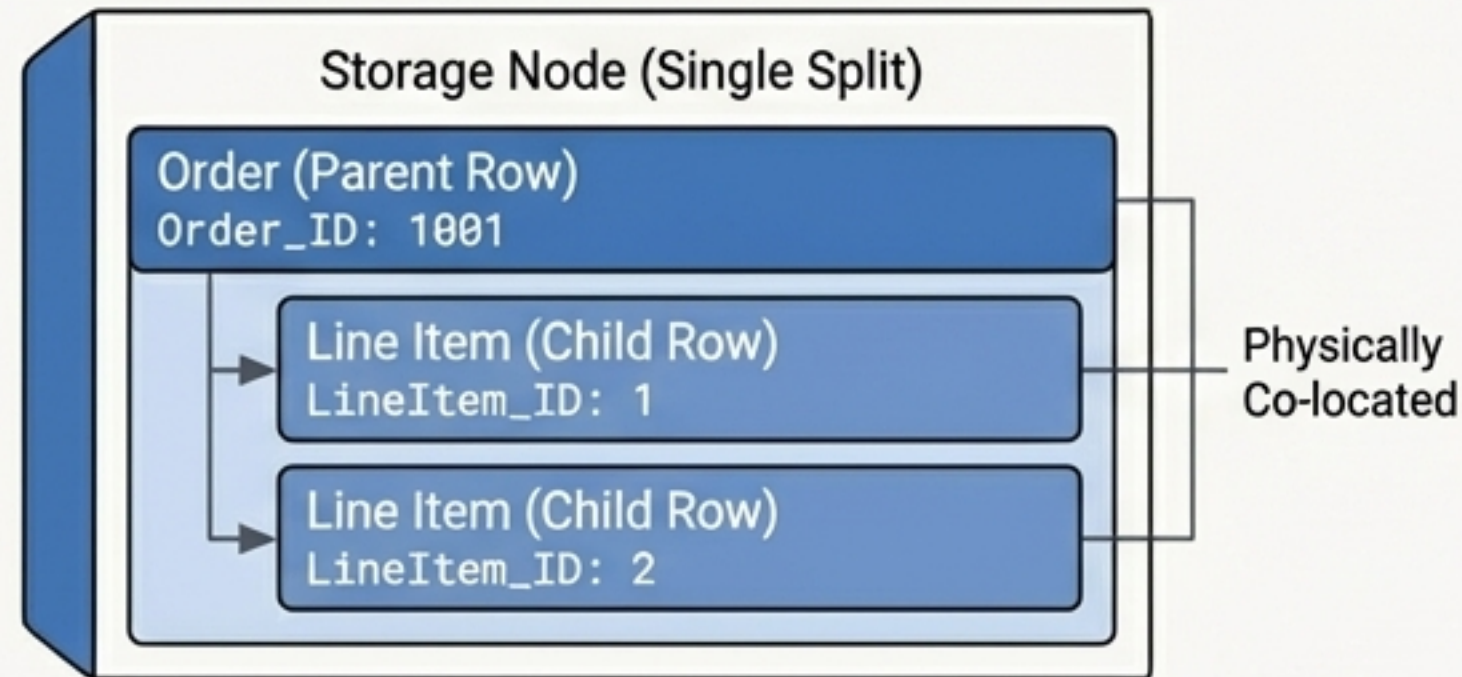
Schema Design & Consistency Trade-Offs

Consistency Spectrum



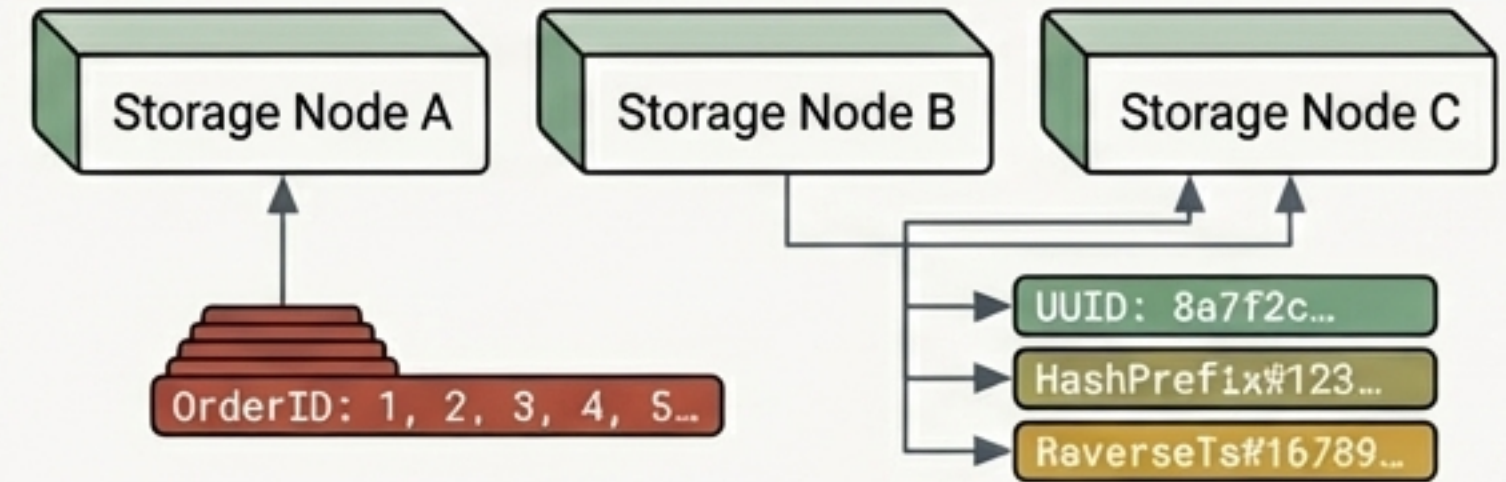
Schema Design Rules

Spanner Schema Interleaving



Prevent cross-server joins by physically co-locating parent and child rows. E.g., `INTERLEAVE IN PARENT` Orders.

Preventing Hotspots



BAD: Auto-incrementing IDs
Hotspot on single node.

GOOD: Distributed Keys
Writes evenly distributed
across all nodes.

Never use auto-incrementing integer IDs. In Spanner, use UUIDs. In Bigtable, use hash-prefixed row keys or reverse-timestamps (`<device-id>#<reverse-timestamp>`) to evenly distribute writes.

Analytics Ingestion and BigQuery

