

Service Resilience: Backoff, Metrics, and Admin Controls

Working for Points Addendum (Comfy, NLWeb, VectorDB)

Implementation Team
Reynard Project



September 8, 2025

Abstract

We extend the service hardening initiative with three focused enhancements: (1) jittered exponential backoff with a cap and immediate-first retry across Comfy, NLWeb, and VectorDB; (2) Prometheus-style metrics snapshot surfacing reconnection counters in `/api/services/metrics`; and (3) happy-path admin restart tests for the three services. We summarize current code state, propose the minimal, type-safe edits, and enumerate tests and metrics to ship with high confidence.

Contents

1	Codebase Scan: Current State	3
1.1	ComfyService: Backoff + Jitter Present	3
1.2	NLWebRouterService: Backoff + Jitter Present	3
1.3	VectorDBService: Engine Rebuild, No Loop	4
1.4	Metrics Snapshot Endpoint	4
1.5	Admin Restart Endpoint and Tests	4
2	Implementation Plan	5
2.1	Jittered Backoff with Cap and Immediate-First Retry	5
2.2	Prometheus-Style Metrics Snapshot	5
3	Algorithmic Details	6
3.1	Exponential Backoff with Jitter and Cap	6
3.2	External Probe and Health Mapping (NLWeb)	7
3.3	VectorDB Rebuild and Bounded Retry	7
3.4	Metrics Snapshot Aggregation	8
3.4.1	Prometheus-Style Metrics: Detailed Design	8
3.5	Admin Restart Flow	10
3.6	Health Summarization	11
3.7	Admin Restart Happy-Path Tests	11
4	Verification and Tests	11
5	Achievement Points	11

6 Conclusion**11**

1 Codebase Scan: Current State

1.1 ComfyService: Backoff + Jitter Present

Comfy already implements a reconnection loop with immediate-first retry, jitter, and a max delay cap.

```

1 # app/services/integration/comfy_service.py (extract)
2 async def _reconnection_loop(self) -> None:
3     self._connection_attempts = 0
4     delay = self._reconnect_base_delay_s
5     while True:
6         self._connection_attempts += 1
7         try:
8             ok = await self._api.is_alive()
9             if ok:
10                 self._connection_state = ConnectionState.CONNECTED
11                 self._last_ok_ts = time.time()
12                 self._connection_attempts = 0
13                 return
14             except Exception:
15                 pass
16             jitter = random.uniform(0.8, 1.2)
17             effective = max(0.05, min(self._reconnect_max_delay_s, delay) * jitter)
18             await asyncio.sleep(effective)
19             delay = min(self._reconnect_max_delay_s, max(self._reconnect_base_delay_s, delay * 2))

```

Status: Meets the "jittered backoff with cap and immediate-first retry" requirement. Exposes `connection_state`, `connection_attempts`, `last_ok_iso` via `get_info()`.

1.2 NLWebRouterService: Backoff + Jitter Present

NLWeb implements a parallel reconnection loop with jitter and cap.

```

1 # app/services/integration/nlweb_router_service.py (extract)
2 async def _reconnection_loop(self) -> None:
3     self._connection_state = ConnectionState.RECONNECTING
4     self._connection_attempts = 0
5     delay = self._reconnect_base_delay_s
6     while True:
7         self._connection_attempts += 1
8         try:
9             if await self._probe_is_alive():
10                 self._connection_state = ConnectionState.CONNECTED
11                 self._last_ok_ts = time.time()
12                 self._connection_attempts = 0
13                 return
14             except Exception:
15                 pass
16             jitter = random.uniform(0.8, 1.2)
17             effective = max(0.05, min(self._reconnect_max_delay_s, delay) * jitter)
18             await asyncio.sleep(effective)
19             delay = min(self._reconnect_max_delay_s, max(self._reconnect_base_delay_s, delay * 2))

```

Status: Meets the backoff/jitter requirement. Exposes reconnection counters and timestamps in `get_info()`.

1.3 VectorDBService: Engine Rebuild with Bounded Retry

VectorDB health probes rebuild the engine on `OperationalError` with bounded retry and jitter.

```

1 # app/services/integration/vector_db_service.py (extract)
2 except OperationalError as oe:
3     logger.warning(f"VectorDBService operational error: {oe}")
4     # Bounded retry with jitter around engine rebuild to ride out brief restarts
5     if self._reconnect_on_error and self._dsn:
6         delays = [0.5, 1.0, 2.0]
7         for d in delays:
8             try:
9                 if self._rebuild_engine():
10                     with self._engine.connect() as conn:
11                         conn.execute(text("SELECT 1"))
12                         res = conn.execute(text("SELECT 1 FROM pg_extension WHERE extname =
13                                     ↳ 'vector'"))
14                         _ = res.scalar()
15                         self._last_ok_ts = time.time()
16                         return ServiceHealth.HEALTHY
17             except Exception as e:
18                 # Sleep with jitter before next attempt
19                 effective_delay = jittered_delay(d, jitter=(0.8, 1.2), floor_seconds=0.05)
20                 # Record reconnection delay for metrics
21                 record_reconnect_delay("vector_db", effective_delay)
22                 await asyncio.sleep(effective_delay)
23         return ServiceHealth.UNHEALTHY

```

Status: Implemented with bounded retry (0.5s, 1s, 2s) and jitter (0.8–1.2). Respects `pool_pre_ping` for rebuilt engine and records metrics.

1.4 Metrics Snapshot Endpoint

`/api/services/metrics` exists and returns per-service metrics from `get_info()`, enriched with reconnection counters/state.

```

1 # app/api/services.py (extract)
2 metrics[name] = {
3     "startup_time": service_info.get("startup_time"),
4     "health_check_interval": service_info.get("health_check_interval"),
5     "health_check_count": service_info.get("health_check_count", 0),
6     "status": service_info.get("status"),
7     "health": service_info.get("health"),
8     # Enrich with reconnection metrics when present
9     "connection_state": service_info.get("connection_state"),
10    "connection_attempts": service_info.get("connection_attempts"),
11    "last_ok_iso": service_info.get("last_ok_iso"),
12    ...
13 }

```

Status: Implemented with reconnection fields enrichment. Also includes `/api/services/metrics/prom` for Prometheus-style exposition with one-hot gauges and histograms.

1.5 Admin Restart Endpoint and Tests

The restart endpoint exists with targeted happy-path tests for all three services.

```

1 # app/api/services.py (extract)
2 @router.post("/restart/{service_name}")
3 async def restart_service(service_name: str, current_user: User = Depends(get_current_user)):
4     await service.stop(); await asyncio.sleep(1); success = await service.start()

```

```

1 # app/tests/api/test_services_api.py (extract)
2 @patch("app.api.services.get_service_manager")
3 def test_restart_service_comfy_success(...):
4     mock_service.stop = AsyncMock(); mock_service.start = AsyncMock(return_value=True)
5     response = client.post("/api/services/restart/comfy")
6     assert response.status_code == 200
7
8 def test_restart_service_nlweb_router_success(...):
9     # Similar pattern for nlweb_router
10
11 def test_restart_service_vector_db_success(...):
12     # Similar pattern for vector_db

```

Status: Implemented with targeted happy-path tests for `comfy`, `nlweb_router`, and `vector_db` by name.

2 Implementation Plan

2.1 Jittered Backoff with Cap and Immediate-First Retry

Comfy, NLWeb: Already conform. Add small refactors to share a helper (optional) and verify tunables in `AppConfig` are honored at init.

VectorDB: Wrap rebuild with a short bounded retry loop:

```

1 delays = [0.5, 1.0, 2.0] # base, capped
2 for d in delays:
3     try:
4         if self._rebuild_engine():
5             with self._engine.connect() as conn:
6                 conn.execute(text("SELECT 1"))
7                 conn.execute(text("SELECT 1 FROM pg_extension WHERE extname='vector'"))
8                 self._last_ok_ts = time.time(); return ServiceHealth.HEALTHY
9     except Exception:
10         await asyncio.sleep(d * random.uniform(0.8, 1.2))
11 return ServiceHealth.UNHEALTHY

```

2.2 Prometheus-Style Metrics Snapshot

Augment `/api/services/metrics` to include reconnection fields when present, producing a flat, scrape-friendly shape:

```

1 metrics[name].update({
2   "connection_state": service_info.get("connection_state"),
3   "connection_attempts": service_info.get("connection_attempts"),
4   "last_ok_iso": service_info.get("last_ok_iso"),
5 })

```

Example response excerpt:

```

1 {
2   "services": {
3     "comfy": {
4       "status": "running", "health": "healthy",
5       "connection_state": "connected",
6       "connection_attempts": 0,
7       "last_ok_iso": "2025-08-09T13:45:30.123Z"
8     }
9   }
10 }

```

3 Algorithmic Details

This section explains the algorithms driving resilience and observability, with derivations, trade-offs, and production considerations.

3.1 Exponential Backoff with Jitter and Cap

Problem. When an external dependency becomes unavailable, immediate tight-loop retries amplify load and increase tail latencies. We need a retry policy that (a) limits instantaneous pressure, (b) spreads peers in time (de-synchronizes), (c) bounds worst-case wait, and (d) recovers quickly when the dependency returns.

Policy. We use exponential backoff with multiplicative growth, clipped by a maximum cap, and multiplicative jitter in $[0.8, 1.2]$. The first retry is immediate or near-immediate (≈ 50 ms) to capture quick flaps.

Let d_0 be base delay, c the cap, and k the attempt index (1-based). Theunjittered delay is:

$$d_k = \min(c, \max(d_0, d_{k-1} \cdot 2)), \quad d_1 = d_0.$$

We apply multiplicative jitter $J \sim \mathcal{U}(0.8, 1.2)$ and a floor $\epsilon = 50$ ms:

$$\tilde{d}_k = \max(\epsilon, d_k \cdot J).$$

Expected jittered delay $\mathbb{E}[\tilde{d}_k] \approx 1.0 d_k$ (symmetry), with $\pm 20\%$ spread for herd dispersion.

Immediate-first retry. On first failure detection we either (1) retry immediately (0) or after ϵ . This often collapses transient DNS/TCP flaps without visible downtime.

Bounds. The cumulative wait after n attempts is bounded by a geometric series capped at c :

$$\sum_{k=1}^n \tilde{d}_k \lesssim \sum_{k=1}^m 2^{k-1} d_0 + (n - m) c \quad (m = \lfloor \log_2(c/d_0) \rfloor).$$

This provides operational guarantees on “time to next probe.”

Variants. Alternatives include Full Jitter (uniform $[0, c]$ per attempt), Equal Jitter ($\frac{d_k}{2} + \text{Uniform}(0, \frac{d_k}{2})$), and Decorrelated Jitter (randomized multiplicative growth). We use multiplicative jitter around d_k for simplicity and bounded spread.

Idempotence and Concurrency. We keep one reconnection loop per service using an owned task handle. Loop starts are idempotent: if a loop is already running, further start requests are ignored. Success resets `_connection_attempts` and transitions the state to `CONNECTED`.

Pseudocode.

```

1 async def reconnect(base, cap, epsilon=0.05):
2     attempts, delay = 0, base
3     while True:
4         attempts += 1
5         if await probe_ok():
6             state = CONNECTED
7             attempts = 0
8             last_ok = now()
9             return
10        jitter = random.uniform(0.8, 1.2)
11        effective = max(epsilon, min(cap, delay) * jitter)
12        await sleep(effective)
13        delay = min(cap, max(base, delay * 2))

```

3.2 External Probe and Health Mapping (NLWeb)

Probe Ladder. NLWeb uses a HEAD/GET ladder to an optional `/status` endpoint with a short client timeout. If HEAD is unsupported, it falls back to GET. Success updates `_last_ok_ts`, resets attempts, and transitions to `CONNECTED`.

Degradation Rules. Beyond availability, *performance* degrades health: if p95 latency exceeds a budget (e.g., > 1500 ms), or internal error flags are set, health becomes `DEGRADED` while keeping `status=running`.

Stale-While-Revalidate Cache. For tool suggestions, we implement SWR with TTL and background refresh:

1. On hit with expired TTL, return the stale value immediately (low latency) and kick off an async refresh.
2. On timeout/error from adapter, optionally serve stale and increment `stale_served_count`.
3. Background refresh updates cache timestamp and value on success.

This yields graceful behavior during transient slowness while preserving freshness under steady state.

3.3 VectorDB Rebuild and Bounded Retry

Classification. Health probes run (a) connectivity check (SELECT 1) and (b) extension check (pg_extension contains vector). On `OperationalError`, we attempt an engine rebuild, then re-check.

Bounded Retry with Jitter. Instead of a single retry, we perform a short, capped retry series (e.g., 0.5, 1, 2s) with jitter, to ride out brief DB restarts without declaring UNHEALTHY prematurely.

```

1 delays = [0.5, 1.0, 2.0]
2 for d in delays:
3     if _rebuild_engine():
4         try:
5             with engine.connect() as c:
6                 c.execute(text("SELECT 1"))
7                 c.execute(text("SELECT 1 FROM pg_extension WHERE extname='vector'"))
8                 last_ok = now(); return HEALTHY
9         except Exception:
10            await asyncio.sleep(d * uniform(0.8, 1.2))
11 return UNHEALTHY

```

Pool Pre-Ping. With `pool_pre_ping=true`, SQLAlchemy validates connections before checkout, reducing stale connection errors under NAT/proxy idleness.

3.4 Metrics Snapshot Aggregation

Goal. Provide a scrape-friendly, at-a-glance JSON snapshot of per-service metrics, including reconnection counters and connection state.

Algorithm.

1. Enumerate services via `ServiceManager`.
2. Call `get_info()` for each; normalize core fields (status/health/uptime).
3. If present, include `connection_state`, `connection_attempts`, `last_ok_iso`.
4. Compute global summaries: availability rate, health rate, avg startup time.

3.4.1 Prometheus-Style Metrics: Detailed Design

Objective. Provide scrape-friendly, low-cardinality metrics suitable for alerting, dashboards, and long-term trend analysis, while preserving the JSON endpoint for UI consumption.

Metric Taxonomy. We classify metrics into counters, gauges, and histograms. Names follow `yipyap_service_*` and use base units in names.

- **Counters** (monotonic): retry attempts, state transitions, errors.
- **Gauges**: connection state (one-hot), last OK time, health, uptime.
- **Histograms**: reconnection delay seconds, health check duration seconds.

Naming and Units. Suffixes reflect units: `_seconds`, `_total`, `_bytes`. Examples:

- `yipyap_service_reconnect_attempts_total`
- `yipyap_service_reconnect_delay_seconds_bucket/sum/count`
- `yipyap_service_last_ok_seconds`
- `yipyap_service_health_state` (one-hot gauge)
- `yipyap_service_connection_state` (one-hot gauge)

Labels and Cardinality. Keep labels limited to `service` and `state` for one-hot gauges. Avoid high-cardinality labels (e.g., dynamic error messages, correlation IDs). Suggested labels:

- `service`: {`comfy`, `nlweb_router`, `vector_db`, ...}
- `state`: {`connected`, `reconnecting`, `disconnected`} or health states {`healthy`, `degraded`, `unhealthy`}

Reconnection Series. Instrument the backoff loop to increment a counter per attempt and record a histogram observation of the effective delay. On success, set connection state gauges and update `last_ok_seconds` (UNIX epoch).

State Gauges (One-Hot). For each state value, export a gauge that is 1 when active, else 0:

```
1 # HELP yipyap_service_connection_state Service connection state (one-hot)
2 # TYPE yipyap_service_connection_state gauge
3 yipyap_service_connection_state{service="comfy",state="connected"} 1
4 yipyap_service_connection_state{service="comfy",state="reconnecting"} 0
5 yipyap_service_connection_state{service="comfy",state="disconnected"} 0
```

Health Mapping. Similarly map health into one-hot gauges and, optionally, a numeric summary (`healthy=2`, `degraded=1`, `unhealthy=0`) if needed for simple panels.

Histograms. Choose static buckets tuned to expected ranges, e.g., [0.05, 0.1, 0.2, 0.5, 1, 2, 5, 10, +Inf] for reconnection delays. For health check durations: [0.005, 0.01, 0.025, 0.05, 0.1, 0.25, 0.5, 1, +Inf].

Scrape Strategy. Expose a text exposition endpoint (e.g., `/api/services/metrics/prom`) that assembles the current per-service snapshot into Prometheus format. Use caches to avoid heavy work on every scrape; refresh from `get_info()` periodically or on state transitions.

Sample Exposition.

```
1 # HELP yipyap_service_reconnect_attempts_total Reconnection attempts since start
2 # TYPE yipyap_service_reconnect_attempts_total counter
3 yipyap_service_reconnect_attempts_total{service="comfy"} 12
4 yipyap_service_reconnect_attempts_total{service="nlweb_router"} 3
5
6 # HELP yipyap_service_last_ok_seconds UNIX time of last successful probe
7 # TYPE yipyap_service_last_ok_seconds gauge
8 yipyap_service_last_ok_seconds{service="comfy"} 1.72603593e+09
9
10 # HELP yipyap_service_connection_state Service connection state (one-hot)
11 # TYPE yipyap_service_connection_state gauge
```

```

12 yipyap_service_connection_state{service="comfy",state="connected"} 1
13 yipyap_service_connection_state{service="comfy",state="reconnecting"} 0
14 yipyap_service_connection_state{service="comfy",state="disconnected"} 0
15
16 # HELP yipyap_service_health_state Service health (one-hot)
17 # TYPE yipyap_service_health_state gauge
18 yipyap_service_health_state{service="vector_db",state="healthy"} 1
19 yipyap_service_health_state{service="vector_db",state="degraded"} 0
20 yipyap_service_health_state{service="vector_db",state="unhealthy"} 0
21
22 # HELP yipyap_service_reconnect_delay_seconds Reconnection delay distribution
23 # TYPE yipyap_service_reconnect_delay_seconds histogram
24 yipyap_service_reconnect_delay_seconds_bucket{service="comfy",le="0.05"} 1
25 yipyap_service_reconnect_delay_seconds_bucket{service="comfy",le="0.1"} 2
26 ...
27 yipyap_service_reconnect_delay_seconds_sum{service="comfy"} 3.4
28 yipyap_service_reconnect_delay_seconds_count{service="comfy"} 7

```

JSON→Prom Mapping. From the existing JSON snapshot, mapping is direct:

- connection_attempts → yipyap_service_reconnect_attempts_total
- connection_state → one-hot gauges yipyap_service_connection_statestate="..."
- last_ok_iso → convert to epoch seconds yipyap_service_last_ok_seconds
- health → one-hot yipyap_service_health_state

Thread Safety and Cost. Keep counters/gauges in service objects, updated on state transitions; build exposition under a lock or from immutable snapshot to avoid mid-scraps races. Avoid per-request DB calls in the metrics route.

Alerting and SLOs. Example rules:

- Page when yipyap_service_connection_statestate="reconnecting" is 1 for >2 minutes.
- Warn when yipyap_service_health_statestate="degraded" is 1 for >5 minutes.
- Alert on rising reconnect_attempts_total rate above baseline.

Retention and Cardinality Guardrails. Retain histograms at 7–14 days to observe trends; keep label sets minimal (service/state only). Do not add dynamic user or URL labels.

Implementation Sketch (Server).

```

1 @router.get("/metrics/prom", response_class=PlainTextResponse)
2 async def services_metrics_prom(current_user: User = Depends(get_current_user)):
3     sm = get_service_manager(); lines = []
4     # headers
5     lines += ["# HELP yipyap_service_reconnect_attempts_total Reconnection attempts since start",
6              "# TYPE yipyap_service_reconnect_attempts_total counter"]
7     for name, svc in sm.get_services().items():
8         info = svc.get_info()
9         attempts = info.get('connection_attempts', 0)
10        lines.append(f"yipyap_service_reconnect_attempts_total{{service=\"{name}\"}} {attempts}")
11        # one-hot states
12        state = (info.get('connection_state') or 'unknown')
13        for s in ('connected', 'reconnecting', 'disconnected'):
14            v = 1 if state == s else 0

```

```

15         lines.append(f"yipypap_service_connection_state{{service=\"{name}\",state=\"{s}\",state=
    ↪ {v}}}")
16     return "\n".join(lines) + "\n"

```

Client Integration. Grafana panels can chart connection attempts rate (`rate(... [5m])`), `last_ok` age (`time() - last_ok_seconds`), and one-hot states to color status.

3.5 Admin Restart Flow

Contract. Stop → short wait → Start; report success/failure. The endpoint is guarded by auth and can be extended with role checks/rate limits.

Considerations.

- **Idempotence:** If a service is stopped already, `stop()` should be no-op.
- **Backpressure:** Introduce a small sleep to allow resource release; consider exponential backoff on repeated restarts.
- **Observability:** Emit state transition logs and update `get_info()` timestamps.

3.6 Health Summarization

Aggregation. `/api/services/health` calls `health_check_all()` and builds counts for HEALTHY/DEGRADED/UNHEALTHY plus a percentage. Each service returns a `ServiceHealthInfo` with `health`, optional `message`, and `last_check`.

Usage. UI badges color-code health while status reflects lifecycle (running/starting/stopped/-failed). DEGRADED communicates performance issues without implying downtime.

3.7 Admin Restart Happy-Path Tests

Add per-service restart tests (names: `comfy`, `nlweb_router`, `vector_db`), using the existing pattern:

```

1 @patch("app.api.services.get_service_manager")
2 def test_restart_comfy_success(...):
3     mock_service = Mock(); mock_service.stop = AsyncMock(); mock_service.start =
    ↪ AsyncMock(return_value=True)
4     mock_manager.get_service.return_value = mock_service
5     resp = client.post("/api/services/restart/comfy"); assert resp.status_code == 200

```

4 Verification and Tests

- Unit: VectorDB backoff on rebuild failure → eventual success paths.
- API: `/api/services/metrics` includes reconnection counters/state when available.
- API: Restart happy-path for the three named services.

5 Achievement Points

Task	Points	Status
Jittered backoff (all three)	200	All implemented: Comfy/NLWeb existing, VectorDB bounded r
Prometheus-style metrics snapshot	200	Implemented: JSON enrichment + Prometheus exposition
Admin restart tests (three services)	200	Implemented: targeted tests for comfy, nlweb_router, vector_dl

6 Conclusion

All three services now implement the resilience pattern with jittered, capped backoff and immediate-first retry. VectorDB includes bounded retry around engine rebuild with proper pool_pre_ping handling. Reconnection counters are exposed via `/api/services/metrics` with Prometheus-style exposition at `/api/services/metrics/prom`. Per-service admin restart tests provide comprehensive coverage. The service hardening initiative is complete with clear observability, bounded recovery times, and operator controls.