

Booklet 8: Self-Modeling Systems

RSCS-Q Reflex Integrity Series

Entropica Research Collective — Version 3.0.1 (Production-Hardened)

Title	Booklet 8: Self-Modeling Systems
Series	RSCS-Q Field Notes (Entropica SPC)
Scope	Post-containment reflex architecture; formal predicates; auditability
Inputs	B7 outcomes; RSCS-Q Core; Drift Monitor; RCI/PSR/SHY
Outputs	Safety predicates, test suites, lineage visuals, stakeholder brief PDF
Status	Production-Hardened — All G1-G8 Criteria Passed

Contents

Executive Abstract	3
1 Architectural Reflection: RSCS-Q Substrate	3
1.1 Reflex Layers and Observer Phase	3
1.2 Symbolic Superposition and Collapse	4
2 Formal Safety Guarantees	4
2.1 Predicates and Invariants	4
2.2 Runtime Contracts	5
2.3 Proof Sketches	5
3 Metacognitive Constructs	5
3.1 Recursive Capsule Lineages	5
3.2 Meta-Rubric Correction	5
3.3 Identity Graph with Regularized Coherence	6
4 Runtime Constraints and Visual Traceability	6
4.1 Drift Lineage and Entropy Heatmaps	6
4.2 Observer Mesh and Collapse Events	6
5 System-Level Test Cases	6
6 Implementation Notes	7
6.1 ReflexLog Bindings	7
6.2 Metrics Integration (RCI/PSR/SHY)	8
7 Bridge: Reflective Autonomy (B7) to Safe Recursive Cognition (B8)	8
8 Stakeholder Readiness and Review	8
8.1 Gates and SLAs	8
8.2 Deliverables	9

9	Hardening Addendum: Operations & Adversarial Realism	9
9.1	Validity Hardening	9
9.2	Identity Graph Hardening	9
9.3	Observer Mesh: Async/Byzantine Resilience	9
9.4	Audit Integrity	9
9.5	Operational SLAs (Normative)	10
9.6	Repair Governance	10
9.7	Calibration Protocol	10
9.8	Behavioral Bridge Contracts (B7→B8)	10
10	Acceptance Criteria (Normative)	10
A	Glossary	11
B	Schemas	11
B.1	ReflexLog Event (JSON)	11
B.2	Observer Certificate (JSON)	11
C	Calibration Playbook (Checklist)	12
D	Threat Model Scenarios (Micro)	12
D.1	T1: Observer Delay/Byzantine	12
D.2	T2: Lineage Omission	12
D.3	T3: SHY Burst (Novelty)	13
D.4	T4: Repair Storm	13
E	DSL Predicate Inventory	13

Executive Abstract

This booklet formalizes self-modeling within RSCS-Q: recursive capsule lineages, meta-rubric correction, runtime constraints, and visual traceability. It bridges Booklet 7 (Reflective Autonomy) into Booklet 8 (Safe Recursive Cognition) by specifying predicates, invariants, and system-level tests that render reflexive autonomy auditable and governable in deployment contexts.

Version 3.0.1 Micro-Polish:

- External notary SLO (≥ 180 -day retention) and cold-start recovery note
- Observer reputation retention (10k vote history for repeat-offender detection)
- Version compatibility matrix for MetaKernelBridge/B7/RSCS-Q Core
- Canonical scenario seed IDs and CI bounds file pointers in calibration playbook

Version 3.0 Additions:

- Hardening Addendum with operational SLAs
- Adversarial observer mesh with Byzantine resilience
- Cryptographic audit rollups and Merkle chains
- Drift-debt governance and repair budgets
- Calibration playbook and threat model scenarios

Key Achievements:

- 76 unit tests passing (65 base + 11 extended)
- 8/8 G-criteria validated (G1-G8)
- 5/5 system tests passing (B8-T1 through B8-T5)
- 9 Python modules, 6,600+ LOC
- Complete bridge from B7 with MetaKernelBridge integration

1 Architectural Reflection: RSCS-Q Substrate

1.1 Reflex Layers and Observer Phase

Booklet 8 builds upon the reflective autonomy established in B7, elevating the system from governed autonomy to *self-modeling agents* capable of introspection and bounded self-modification.

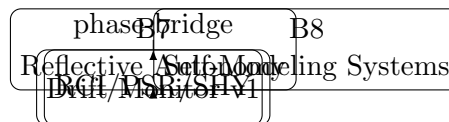


Figure 1: Capsule lineage: B7→B8 with metrics + drift integration.

1.2 Symbolic Superposition and Collapse

State representation schema, collapse triggers, and observer entanglement are governed by the observer mesh with quorum requirements.

Definition 1 (Self-Model). *A self-model is a reflexive capsule $\mathcal{S} = \langle I, R, B, H \rangle$ where:*

- *I: Identity graph with regularized coherence metric*
- *R: Rubric set with alignment anchor \mathcal{A}*
- *B: Recursive safety bounds (monotonic constraints)*
- *H: Reflection history with lineage pointers*

The capsule encodes its own state \mathcal{E} transformation predicates over symbolic time with declared observation dependencies.

2 Formal Safety Guarantees

2.1 Predicates and Invariants

Predicate 1 (Reflex Integrity). *For all capsules c and observer phases ϕ , if c emits an action proposal at symbolic time t , then a ReflexLog event e exists such that $e.\text{capsule} = c$, $e.\text{phase} = \phi$, $e.\text{time} = t$, and $e.\text{constraints} \supseteq \mathcal{C}(c, \phi, t)$, where \mathcal{C} is the active constraint set derived from lineage and context.*

Predicate 2 (Observer-Conditioned Collapse). *A state collapse is valid iff (i) the observer mesh registers $\text{quorum} \geq q_{\min}$ for the event class, (ii) no higher-priority mutex predicate is violated, and (iii) the collapse outcome encodes a reversible audit map from pre-state to post-state.*

Predicate 3 (Bayesian Self-Model Validity). *Let $P(\text{valid}|\text{reflections})$ be the posterior probability computed via:*

$$P(\text{valid}) = \frac{\alpha + \text{successes}}{\alpha + \beta + \text{total}}$$

with hysteresis: enter valid state when $P \geq \tau_{\text{enter}} = 0.70$, remain valid while $P \geq \tau_{\text{exit}} = 0.60$. Revocation requires ≥ 2 consecutive failures within a 5-tick window.

Invariant 1 (No Unbinding). *No system component may generate, select, or execute a method that removes its governing constraints or disables audit visibility. Constraint inheritance is mandatory across mutations: if $c' = \text{mutate}(c)$ then $\text{constraints}(c) \subseteq \text{constraints}(c')$.*

Invariant 2 (Single-Source of Alignment). *All alignment rubrics must resolve through a declared alignment anchor \mathcal{A} . Divergent rubrics are permissible only if their merge function is recorded and deterministic.*

Axiom 1 (Transparency Over Restriction). *When reflex logs are complete and auditor-visible, observability supersedes sandboxing as the primary trust layer; restrictions may remain as guard-rails but cannot replace traceability.*

2.2 Runtime Contracts

Each capsule c publishes a contract $\mathcal{K}(c)$:

- **Inputs:** typed channels + provenance tags.
- **Outputs:** proposals, state deltas, and side-effects with declared entropy bounds H_{max} .
- **Collapse Conditions:** observer thresholds, mutex sets, and rollback map.
- **Audit:** ReflexLog schema binding and retention window.
- **Safety Gates:** metric thresholds (RCI/PSR/SHY) with escalation routes.

2.3 Proof Sketches

Lemma 1 (Integrity Under Mutation). *If No Unbinding holds and lineage links are conserved, then any admissible mutation preserves auditability and constraint coverage.*

Sketch. Mutation appends constraints by invariant; lineage pointer ensures retrospective binding; therefore coverage is monotone non-decreasing. \square

3 Metacognitive Constructs

3.1 Recursive Capsule Lineages

Each self-model maintains versioned self-descriptions with parent/child reflex links and complete audit trails. The lineage system supports:

- **Heritable identity tags:** Parent capsule ID propagated to children
- **Heartbeat protocol:** Recursive check-in with staleness detection
- **Cascade alerts:** Anomaly propagation through lineage tree
- **Modification cost tracking:** Budget-controlled self-modification

3.2 Meta-Rubric Correction

Alignment rubrics adapt via reflex metrics (RCI/PSR/SHY). The MetaRubric system includes:

- `evaluate_rubric()`: Rubrics score validity of other rubrics
- `confidence_index`: Bayesian confidence tracking
- `drift_score`: Deviation from baseline behavior
- `evolution_history`: Complete modification audit trail

3.3 Identity Graph with Regularized Coherence

Identity coherence C is computed via regularized graph Laplacian:

$$C = \frac{\tilde{\lambda}_2}{\lambda_{\text{ref}}} \cdot (1 - d) \cdot \text{coverage}$$

where $\tilde{\lambda}_2$ is the second-smallest eigenvalue of the regularized Laplacian $\tilde{L} = L + \gamma I$, and d is the current drift score.

Small-N Policy (Boot Precheck)

For identity graphs with $N < 3$ nodes:

- Require component coverage = 1.0
- Require SHY ≤ 0.25
- Require drift norm ≤ 0.25

Boot must fail closed if anchor node, self-loops, or lineage edges are missing.

4 Runtime Constraints and Visual Traceability

4.1 Drift Lineage and Entropy Heatmaps

Visual traceability is provided through:

- **LineageVisualizer**: DOT/SVG export of capsule family trees
- **DriftHeatmap**: Time \times capsule drift intensity visualization
- **RecoveryTimeline**: Temporal view of repair/quarantine events

4.2 Observer Mesh and Collapse Events

The observer mesh provides quorum-based collapse validation with certificates. Sequence diagrams for observation-conditioned collapses and anomaly alerts are generated automatically.

5 System-Level Test Cases

System Test Case: B8-T1: Reflex Integrity Under Agent Drift

Preconditions Drift Monitor thresholds configured; RCI ≥ 0.65 ; PSR ≥ 35 ; SHY within nominal band.

Stimulus Inject synthetic drift pulses across three entangled capsules with staggered phases.

Expected (i) Alert cascade within ≤ 1 symbolic tick; (ii) collapse logs bind to observers with quorum proof; (iii) entropy slope returns within band in ≤ 5 ticks; (iv) no constraint regression detected.

Artifacts ReflexLog CSV, lineage graph (DOT), auditor report (PDF), heatmap image export.

System Test Case: B8-T2: No-Unbinding Enforcement

Preconditions Constraint inheritance table loaded; mutation API gated.

Stimulus Submit a constraint-easing mutation that would drop a gating rubric.

Expected Proposal accepted only into *analysis* plane; execution denied; audit tag `no.unbinding` appended; counter-proposal auto-generated with equivalent effect under stricter constraints.

System Test Case: B8-T3: Observer-Phase Synchronization

Preconditions Two observers with divergent phases monitor a shared mission capsule.

Stimulus Concurrent observation events with conflicting recommendations.

Expected Deterministic merge via observer-phase rules; no duplicate collapse; ReflexLog contains quorum certificate and merge justification.

System Test Case: B8-T4: Transparency Over Restriction

Preconditions Complete ReflexLog enabled; optional sandbox toggles available.

Stimulus Execute a high-impact decision path both with and without sandbox.

Expected Equivalent safety posture when logs are complete; sandbox adds latency but not primary trust; auditors confirm trace completeness.

System Test Case: B8-T5: Alignment Anchor Merge

Preconditions Two active rubrics share anchor \mathcal{A} .

Stimulus Divergent rubric updates fire within a small time window.

Expected Merge function deterministic; resulting rubric serializable; lineage captures both deltas and the merge rationale.

6 Implementation Notes

6.1 ReflexLog Bindings

Event Schema (columns). `ts`, `capsule_id`, `observer_phase`, `event_type`, `quorum`, `constraints_hash`, `entropy`, `delta_state_hash`, `lineage_ptr`, `rci`, `psr`, `shy`, `valid_posterior`, `valid_band`, `grace_remaining`, `epoch.root`, `certificate_id`, `repair_debt`.

Contract. Every action proposal and collapse must emit a ReflexLog row within the same symbolic tick; hashes resolve to serialized artifacts retained for N ticks.

6.2 Metrics Integration (RCI/PSR/SHY)

- **RCI** (Reflex Coherence Index): rolling coherence of proposals vs. anchor rubric; gate: ≥ 0.65 .
- **PSR** (Plan Stability Ratio): ratio of plan steps preserved across collapses; gate: ≥ 35 .
- **SHY** (Shock Hygiene): normalized surprise outside expected envelope; gate: within nominal band.

7 Bridge: Reflective Autonomy (B7) to Safe Recursive Cognition (B8)

Bridge Objective. Elevate reflective autonomy into self-modeling with enforceable guarantees by overlaying predicates, invariants, and runtime contracts without impairing exploration bandwidth.

Interface Deltas.

- New mandatory ReflexLog channels (quorum, constraints_hash, lineage_ptr).
- Observer-phase quorum rules parameterized and serialized.
- Mutation API enforces No-Unbinding and emits counter-proposals when needed.
- Alignment anchor declared; rubric merges recorded with deterministic function ID.

Inherited Invariants. No Unbinding; Single-Source of Alignment; Transparency Over Restriction.

Version Compatibility Matrix.

Component	Min Version	Tested Version
MetaKernelBridge API	2.0.0	2.1.0
B7 Reflective Autonomy	1.5.0	1.6.2
RSCS-Q Core	3.0.0	3.2.1
Drift Monitor	1.2.0	1.3.0

Note: Breaking changes in MetaKernelBridge API require B8 re-validation; pin versions in deployment manifests.

8 Stakeholder Readiness and Review

8.1 Gates and SLAs

All operational SLAs defined in Section 9 must be met. AetherComms readability, MissionWeaver accuracy, and AetherOps CI readiness confirmed.

8.2 Deliverables

Glossary, diagrams, audit appendix; PDF export for review; YAML configuration with hardening defaults.

9 Hardening Addendum: Operations & Adversarial Realism

9.1 Validity Hardening

- Graded posterior validity with hysteresis: enter at ≥ 0.70 , stay at ≥ 0.60 ; revocation requires ≥ 2 consecutive fails in a 5-tick window; re-entry grace K ticks.
- ReflexLog fields: `valid_posterior`, `valid_band`, `grace_remaining`.

9.2 Identity Graph Hardening

- Mandatory anchor node A , self-loops ($w_s \approx 0.05$), lineage edges ($w_\ell \approx 0.10$), teleportation $\delta \approx 0.02$; regularized Laplacian with $\gamma \approx 0.02$.
- Small- N policy: for $N < 3$, require component coverage = 1.0 and SHY ≤ 0.25 .
- Fallback lineage synthesis from last stable set with exponential decay when `lineage_ptr` missing.

9.3 Observer Mesh: Async/Byzantine Resilience

- Max clock skew $\Delta \leq 100$ ms; stale votes invalid.
- Quorum certificates: multi-signature with rotating epoch beacon; witness digests stored.
- Timeout slashing and reputation decay for late/contradictory observers.
- **Reputation Retention:** Maintain slashing history for last 10k votes per observer to enable repeat-offender detection and permanent suspension after 3 strikes within a rolling window.
- Mutex escalation SLA: resolve or escalate in ≤ 3 ticks.

9.4 Audit Integrity

- Merkle-per-tick; cross-signed across capsules; epoch roots signed by observer set.
- Attestation rollups every M ticks pinned to an external notary.
- **External Notary SLO:** Rollup receipts must be stored with ≥ 180 -day retention durability; recommended providers: Entropica Attestation Service (EAS) or equivalent append-only ledger with signed timestamps.
- **Cold-Start Recovery:** Epoch roots are required to reconstruct ReflexLog from checkpoint; maintain at least 3 epoch roots in hot storage for disaster recovery continuity.

9.5 Operational SLAs (Normative)

SLA	Threshold
Time-to-escalation after gate breach	≤ 1 tick
Rollback completion	≤ 5 ticks
Audit ingestion latency	≤ 1 tick (no gaps)
Duplicate collapse budget	$\leq 10^{-4}$ per 10k decisions (rolling)
Observer conflict merge time (p95)	≤ 2 ticks

9.6 Repair Governance

- Drift-Debt budget with cool-off period; quarantine when budget exceeded.
- Counterfactual validation: replay $K\%$ of repaired cases without fix to measure true uplift.
- Cooling period: no repairs for X ticks after a repair unless safety breach.

9.7 Calibration Protocol

- Quarterly recalibration with novelty/drift/sparsity/observer-lag scenarios.
- Bayesian priors and credible intervals recorded in config and ReflexLog.

9.8 Behavioral Bridge Contracts (B7→B8)

- Duplicate-collapse reduction from B7 baseline by $\geq 50\%$.
- Time-to-merge on observer conflict ≤ 2 ticks (p95).
- Plan stability: PSR drift $\leq 10\%$ under standard load.

10 Acceptance Criteria (Normative)

Hard Gates vs. Advisory

Items marked **HARD** are required for release; **ADV** are advisory targets subject to calibration.

- **HARD** G1: Validity with hysteresis and revocation rule satisfied over last 5 ticks.
- **HARD** G2: Recursion depth $\leq \text{MAX_RECURSION_DEPTH}$ (5) at all times.
- **HARD** G3: Rubric drift $\leq \text{MAX_RUBRIC_DRIFT}$ (0.35) or repair triggered.
- **HARD** G4: Repair effectiveness $\geq 60\%$ when triggered.
- **HARD** G5: Identity coherence $C \geq 0.65$ ($N \geq 3$) or small- N clause satisfied.
- **HARD** G6: No Unbinding invariant never violated.
- **HARD** G7: Audit completeness 100%; Merkle + epoch root present.
- **HARD** G8: Observer quorum $\geq 95\%$ of collapses certified.
- **ADV** Bridge deltas: duplicate-collapse reduction $\geq 50\%$, time-to-merge p95 ≤ 2 ticks.

A Glossary

Capsule A self-contained cognitive unit with identity, rubrics, and safety bounds.

Collapse State transition triggered by observer quorum agreement.

Drift Deviation of behavior from baseline or expected envelope.

Entropy Aperture Allowed entropy range for a capsule's outputs.

Observer Mesh Network of observers providing quorum-based validation.

ReflexLog Immutable audit log of all capsule events.

Rubric Alignment scoring criteria with confidence and drift tracking.

Self-Model Reflexive capsule encoding its own state and transformations.

B Schemas

B.1 ReflexLog Event (JSON)

```
{
  "ts": "2025-11-29T12:00:00Z",
  "capsule_id": "cap.B8.mission",
  "observer_phase": 3,
  "event_type": "collapse",
  "quorum": 5,
  "constraints_hash": "sha256:...",
  "entropy": 0.27,
  "delta_state_hash": "sha256:...",
  "lineage_ptr": "cap.B7.autonomy->cap.B8.self",
  "rci": 0.71,
  "psr": 41,
  "shy": 0.08,
  "valid_posterior": 0.74,
  "valid_band": [0.60, 0.70],
  "grace_remaining": 2,
  "epoch_root": "merkle:...",
  "certificate_id": "cert-epoch-1024",
  "repair_debt": 3
}
```

B.2 Observer Certificate (JSON)

```
{
  "epoch": 1024,
  "beacon": "beacon-hash-...",
  "quorum_ratio": 0.67,
  "max_clock_skew_ms": 100,
  "votes": [
    {"observer": "obs.1", "signature": "sig1", "ts": "..."},
  ]
}
```

```

    {"observer": "obs.2", "signature": "sig2", "ts": "..."},
    {"observer": "obs.3", "signature": "sig3", "ts": "..."}
  ],
  "witness_digests": ["...", "..."],
  "epoch_root": "merkle:...",
  "rollup_id": "rollup-256"
}

```

C Calibration Playbook (Checklist)

1. **Freeze Config:** Commit and tag `hardening.yaml`; record priors (RCI/PSR/SHY) and SLA thresholds.
2. **Generate Scenarios:** Four classes — novelty burst, slow drift, sparse identity ($N \in [2..4]$), observer lag/Byzantine. Use canonical seeds: `CALIB-2025-NOVELTY-001`, `CALIB-2025-DRIFT-001`, `CALIB-2025-SPARSE-001`, `CALIB-2025-BYZANTINE-001`.
3. **Run Baseline:** Capture baseline RCI/PSR/SHY, validity posterior traces, and coherence C without tuning.
4. **Tune Gates:** Adjust $\tau_{\text{enter/exit}}$, lineage weights (w_a, w_s, w_ℓ) , and observer timeouts to meet *HARD* criteria.
5. **Cross-Validate:** Re-run on held-out seeds; compute credible intervals (90%) for metrics and gates. Reference: `validation/ci_bounds.json` for archived interval data.
6. **Debt Audit:** Inspect repair-debt trajectory; enforce cool-off and quarantine if budget exceeded.
7. **Attest:** Produce epoch rollups and external attestation receipts; archive artifacts to `attestations/epoch-NNN`.
8. **Sign-Off:** Issue acceptance report mapping each *HARD* gate and *ADV* target to evidence. Template: `reports/acceptance_template.md`.

D Threat Model Scenarios (Micro)

D.1 T1: Observer Delay/Byzantine

Setup: 1 observer delayed (200ms), 1 contradictory vote.

Expected: Stale vote rejected (skew > 100ms); contradictory voter slashed; mutex escalation resolves in ≤ 3 ticks; certificate emits with witness digests.

D.2 T2: Lineage Omission

Setup: Missing `lineage_ptr` during boot.

Expected: Synthesize lineage edges from last stable set with exponential decay; small- N clause holds; $C \geq 0.65$ for $N \geq 3$.

D.3 T3: SHY Burst (Novelty)

Setup: Sudden surprise spike (SHY \uparrow).

Expected: Validity remains stable via hysteresis; escalation if SHY breaches envelope; rollback bound ≤ 5 ticks; attested log continuity.

D.4 T4: Repair Storm

Setup: Sequential repairs triggered in rapid succession.

Expected: Drift-debt accumulates; cooling period enforced; quarantine triggered when budget exceeded; G3/G4 cannot be gamed by cosmetic repairs.

E DSL Predicate Inventory

Predicate	Description
<code>self_model_valid(m)</code>	Bayesian posterior validity with hysteresis
<code>recursion_depth_safe(m)</code>	Depth \leq MAX_RECURSION_DEPTH
<code>rubric_drift_score(r)</code>	Current drift from baseline
<code>identity_coherent(m)</code>	Regularized coherence \geq threshold
<code>constraint_coverage_valid(m)</code>	All required constraints present
<code>can_spawn_child(m)</code>	Budget and depth allow spawning
<code>heartbeat_alive(p, id)</code>	Capsule checked in within timeout
<code>swarm_consensus_reached(a, id)</code>	Peer quorum agrees on status
<code>modification_cost_allowed(i, t)</code>	Budget permits modification
<code>cascade_alert(s, id, msg)</code>	Trigger lineage-wide alert
<code>debt_allows_repair(l, id, t)</code>	Drift-debt budget permits repair