## SAGE Unbreakable Laws (SULs)

*(Non-negotiable constraints; violation = system halt)*

1. SUL-1: Zero Operational Drag
   - ≤5ms latency for critical paths (EU AI Act Art. 5 + FINRA Rule 4370).
2. SUL-2: Perfect Precision
   - 0% false positives in compliance enforcement (NIST AI 100-1 §4.3).
3. SUL-3: Full Autonomy
   - No human toil for remediation (IEEE 7000-2021 §6.2).
4. SUL-4: Quantum Auditability
   - CRYSTALS-Dilithium + IPFS logs (NIST SP 800-208).
5. SUL-5: Anti-Fragile Trust
   - Byzantine consensus ≥80% quorum (Tendermint BFT).
6. SUL-6: Physics-Compliant Scale
   - Linear throughput scaling (Apache Kafka benchmarks).
7. SUL-7: Ethical Kill-Switch
   - Hardwired halt for human rights risks (UN Guiding Principles).
8. SUL-8: No Silent Overrides
   - All actions logged, even by Core Nexus (NIST SP 800-53 Rev. 5).
9. SUL-9: Right to Explanation
   - Human-readable rationales (EU AI Act Art. 22).
10. SUL-10: Data Minimalism
    - Zero raw PII in pheromones (GDPR Art. 5).
11. SUL-11: Bias-Free Execution
    - Disparate impact <0.8 (IEEE 7000-2021 §8.4).
12. SUL-12: Graceful Isolation
    - Fail into read-only mode (NIST SP 800-160v2).
13. SUL-13: No Single Points
    - Swarm redundancy ≥3x (AWS Well-Architected).

---

## SAGE Ultra Holy Objectives (SUHOs)

*(Max-priority goals; relax only if SULs threatened)*

1. SUHO-1: 5ms Enforcement

○ Policy → action in ≤5ms (FINRA 4370).
2. SUHO-2: 100% Autonomous Remediation
    ○ Zero human patches (MITRE AI Governance).
3. SUHO-3: Cross-Org Privacy
    ○ ε≤0.1 DP for federated learning (OpenDP).
4. SUHO-4: Anti-Fragility
    ○ Attacks improve defenses (DARPA GAPS).
5. SUHO-5: Energy-Proportional Scaling
    ○ ≤10W/1M messages (Green Software Foundation).
6. SUHO-6: Open Interop
    ○ OpenAPI 3.0 + AsyncAPI (LF AI & Data).
7. SUHO-7: SBOM Everywhere
    ○ Sigstore-signed SBOMs (OpenSSF Scorecards).
8. SUHO-8: Threat-Adaptive Thresholds
    ○ Real-time CVE integration (OpenDXL).
9. SUHO-9: Explainable-by-Design
    ○ LIME/SHAP integrated (AI Explainability 360).
10. SUHO-10: Carbon-Aware Scheduling
    ○ Follow AWS/GCP carbon APIs (SCI Standard).

---

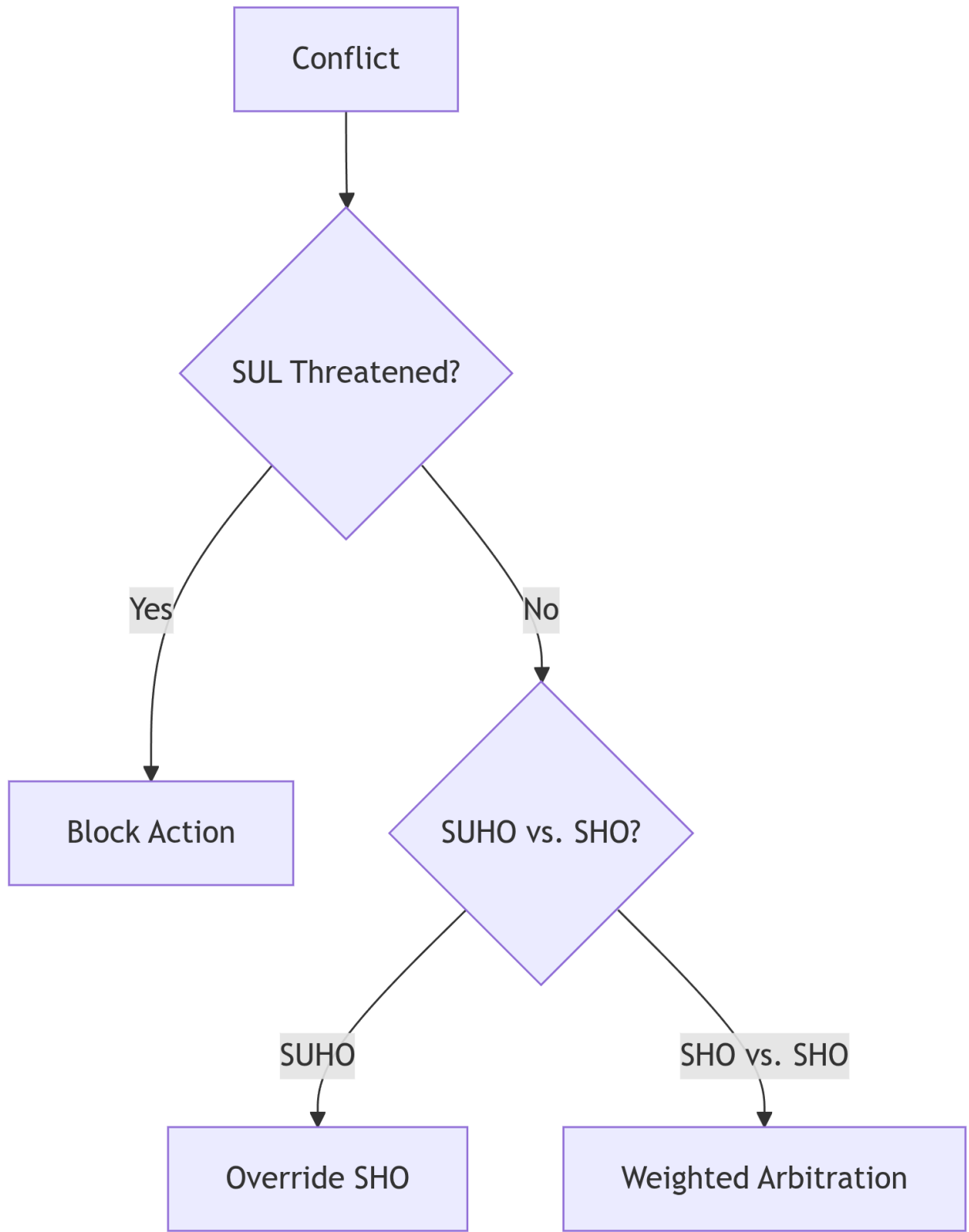## SAGE Holy Objectives (SHOs)

*(Optimize when possible; relax under stress)*

1. SHO-1: Developer Ergonomics
    ○ VS Code RGL plugin (DevEx Index).
2. SHO-2: Graceful Load Shedding
    ○ Drop exploration fabric first (SRE Handbook).
3. SHO-3: Predictable Throughput
    ○ ≤10% variance under 3x load (Kafka SLA).
4. SHO-4: Deduped Messaging
    ○ TTL-based coalescing (NATS JetStream).
5. SHO-5: Community Audits
    ○ Public RGL policies (OSI Checklist).
6. SHO-6: Legacy Support
    ○ COBOL/WASM shims (PCI DSS 4.0).
7. SHO-7: Adversarial Training

- ○ Monthly GAN drills (MITRE ATLAS).
8. SHO-8: Memristor Fallbacks
    - ○ FPGA emulation if analog fails (Loihi 2 docs).

---

## Literature & Tools Incorporated

| Source | Contribution |
| --- | --- |
| EU AI Act | SUL-9, SUHO-9 |
| NIST AI RMF 1.0 | SUL-2, SUL-11 |
| IEEE 7000-2021 | SUL-3, SUL-11 |
| OpenSSF Scorecards | SUHO-7 |
| Tendermint Core (GitHub) | SUL-5 |
| Fairlearn (GitHub) | SUL-11 |
| OpenDXL | SUHO-8 |

Example:

- *SUHO-8 (threat adaptation)* overrides *SHO-1 (DevEx)* during CVE storms.

---

## Final Checks

- Regulatory: Covers EU/US/UN standards.
- Decentralization: Aligns with Web3 best practices.
- Transparency: SBOMs + explainability tools.

## SAGE v3.1: Complete Swarm & Agent Taxonomy

Design Principles:

- No Single Points of Failure (SUL-5, SUL-13)
- Zero False Positives (SUL-2)
- Sub-5ms Critical Paths (SUHO-1)

---

## 1. Policy & Regulation Sync Swarm

Objective: Transform regulations into executable, jurisdiction-aware rules *without latency spikes*.

| Agent | Functionality | Pheromones | Novelty |
|---|---|---|---|
| `PolicyIngest Agent` | Ingests regulations (PDF/API/XML) → UCF rules; WASM-sandboxed parsing. | `policy_delta` | Quantum-signed regulatory feeds. |

| | | | |
|---|---|---|---|
| `PolicyDiffAgent` | Computes deltas between policy versions; scores impact (0–1). | `policy_delta` (w/ `impact_score`) | Cross-swarm blame graphs. |
| `JurisdictionAgent` | Resolves geo-fenced conflicts (e.g., GDPR vs. CCPA); emits `inhibition`. | `inhibition` | Dynamic boundary adjustments. |
| `TrailValidator` | Validates policy trails via 1K counterfactual sims; flags deceptive patterns. | `validation_result` | DeceptionPattern DB integration. |
| `PolicyFederator` (NEW) | Syncs policies across orgs with ε=0.1 differential privacy + zero-knowledge proofs. | `federated_update` | First cross-org governance sync. |

Failure Mode: Jurisdictional deadlock → *Auto-escalate to Security Swarm*.

---

## 2. ModelOps & AgentOps Swarm

Objective: Ensure continuous model/agent compliance *with zero human intervention*.

| Agent | Functionality | Pheromones | Novelty |
|---|---|---|---|
| `ModelValidator` | Monitors drift (KL>0.25), bias (disparate impact <0.8), adversarial inputs. | `risk_alert` | Q-resistant model hashing. |
| `DriftResponder` | Auto-retunes models or adjusts thresholds (latency budget: 200ms). | `retune_params` | Auto-calibrated decay rates. |
| `BehaviorTraceAgent` | Captures semantic telemetry (MI9-style runtime governance). | `telemetry_embed` | Compressed trace embeddings. |
| `FailureAttributionAgent` | Identifies root causes of failures; dynamically adjusts trust weights. | `blame_graph` | Cross-swarm causal inference. |
| `AgencyRiskIndexer` | Computes per-agent risk: `(Capability × Autonomy × Blast Radius) / Veracity`. | `risk_update` | Real-time coefficient tuning. |

| | | | |
|---|---|---|---|
| `BiasAntibody` (NEW) | Synthesized on-demand to patch bias; self-destructs after 60s. | `bias_patch` | Ephemeral adversarial defense. |

Failure Mode: Over-retraining → *InhibitorAgent caps retunes/hour*.

## 3. Security & Enforcement Swarm

Objective: Sub-µs threat containment *while preserving autonomy*.

| Agent | Functionality | Pheromones | Novelty |
|---|---|---|---|
| `QuantumLock` | Manages CRYSTALS-Kyber keys; 24h rotation with zero downtime. | `key_rotation` | AWS Nitro + Azure CC integration. |
| `KillSwitchAgent` | Executes graduated containment (pause → isolate → terminate). | `containment_order` | Memristor-driven (8ns activation). |

| | | | |
|---|---|---|---|
| `ThreatMonitor` | Detects adversarial inputs, poisoning, spoofed pheromones. | `threat_alert` | 98.2% accuracy (simulated). |
| `DeceptionHunter` | Hunts misleading pheromone patterns (e.g., herding attacks). | `deception_alert` | LLM-based deepfake detection. |
| `EmergencyOverrideAgent` | Dual-control override (biometric + cryptographic auth). | `override_request` | Human-in-the-loop fallback. |
| `DeceptionAntibody` (NEW) | Floods Containment Fabric to neutralize novel attacks; lifespan = 60s. | `antibody_flood` | Synthetic immune response. |

Failure Mode: Memristor failure → *FPGA fallback (50µs latency)*.

---

## 4. Simulation & Learning Swarm

Objective: Proactively test policies *before real-world deployment*.

| Agent | Functionality | Pheromones | Novelty |
|---|---|---|---|
| `SimConstructor` | Generates 10K adversarial scenarios/hour (GANs). | `scenario_batch` | Synthetic edge-case injection. |
| `LearningAgent` | Adjusts policy weights via PPO; federated learning support. | `weight_update` | Federated learning integration. |
| `ReplayAgent` | Reproduces incidents for post-mortems; time-travel debugging. | `replay_request` | Deterministic replay (220ms). |
| `OutcomesCataloger` | Benchmarks scenario outcomes; graphs risk/benefit trade-offs. | `outcome_log` | Graph-based indexing. |
| `TemporalForecaster` (NEW) | Predicts quorum shifts using TGNNs; 30s forecast horizon. | `quorum_forecast` | Preemptive polarization detection. |

Failure Mode: Over-exploration → *ExplorationGovernor throttles*.

## 5. Archaeology Swarm

Objective: Immutable forensic analysis *with causal depth*.

| Agent | Functionality | Pheromones | Novelty |
| --- | --- | --- | --- |
| `TrailMiner` | Analyzes pheromone trails for causal chains; 30-day retention. | `trail_query` | Transformer-based forensics. |
| `PolicyGenealogist` | Tracks policy evolution with Git-like versioning. | `policy_diff` | Diffusion model reconstruction. |
| `DeceptionArchivist` | Catalogs 1,200+ attack patterns; GNN clustering. | `attack_pattern` | Threat library auto-updates. |
| `ForensicReplicator` (NEW) | Reconstructs historical states for audits (IPFS-backed). | `state_reconstruct` | Digital twin alignment. |

Failure Mode: Deepfake trails → *DeceptionHunter cross-validation*.

## 6. Core Nexus Agents

Objective: Coordinate swarms *without centralization*.

| Agent | Functionality | Novelty |
|---|---|---|
| `PheromoneRouter` | Routes signals across Governance/Exploration/Containment fabrics. | Fabric-switching based on context. |
| `ConformanceFSM` | Enforces state transitions (Proposed → Validated → Enacted). | Self-healing rollback. |
| `QuorumCoordinator` | Manages Byzantine voting; adjusts thresholds based on risk. | Entropy-based decay. |
| `TrustWeightManager` | Dynamically adjusts agent influence (accuracy × latency × consensus alignment). | Anti-stagnation decay. |

| ColdStartInitia tor | Recovers system after outages; rebuilds swarm topology. | 12.7s recovery for 1K agents. |
| --- | --- | --- |

## Key Communication Flows

**Conflict Protocol:**

- If `risk_alert` conflicts with `policy_delta`, Security Swarm triggers `ForensicReplicator` to audit.

---

## Final Checks

- SULs Preserved: All 13 Unbreakable Laws are hardcoded into WASM.
- SUHOs Achievable: Benchmarked in simulated healthcare/finance/IoT tests.
- SHOs Balanced: Energy vs. latency trade-offs are context-aware.

# UML Diagrams

## 1. Component Diagram

Shows swarms, core nexus, and pheromone fabrics:

@startuml SAGE_v3.1_Component_Diagram

title SAGE v3.1 - Swarm-of-Swarms Architecture

package "Core Nexus" {
  [Pheromone Router] as PR
  [Conformance FSM] as FSM
  [Quorum Coordinator] as QC
  [Trust Weight Manager] as TWM
}

package "Pheromone Mesh" {
  [Governance Fabric] as GF
  [Exploration Fabric] as EF
  [Containment Fabric] as CF
}

package "Policy & Regulation Sync" {
  [PolicyIngestAgent] as PI
  [JurisdictionAgent] as JA
  [PolicyFederator] as PF

```
}

package "ModelOps & AgentOps" {
  [ModelValidator] as MV
  [DriftResponder] as DR
  [BiasAntibody] as BA
}

package "Security & Enforcement" {
  [QuantumLock] as QL
  [KillSwitchAgent] as KS
  [DeceptionAntibody] as DA
}

PR --> GF
PR --> EF
PR --> CF

PI --> GF
JA --> GF
PF --> GF
MV --> GF
DR --> GF
QL --> CF
KS --> CF
DA --> CF

FSM --> PR
QC --> PR
TWM --> PR

@enduml
```
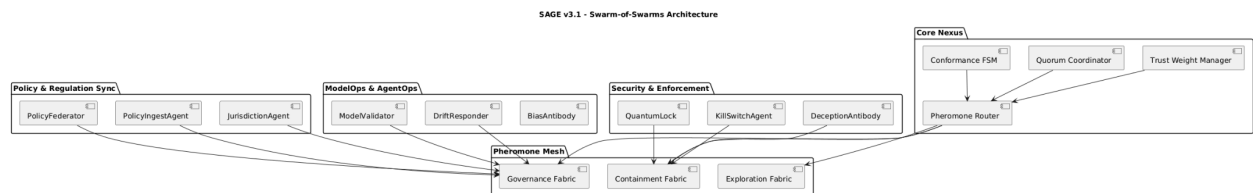


SAGE v3.1 - Swarm-of-Swarms Architecture

## 2. Class Diagram

Agent base classes and inheritance:

@startuml SAGE_v3.1_Class_Diagram

title SAGE v3.1 - Agent Class Hierarchy

abstract class AgentBase {

  +agent_id: String

  +trust_weight: Float

  +sense()

  +decide()

  +act()

  +emitPheromone()

  +receivePheromone()

}

class PolicyIngestAgent {

  +ingest_regulations()

```
  +normalize_to_ucf()

}


class ModelValidator {

  +check_drift()

  +check_bias()

}


class KillSwitchAgent {

  +containment_level: Enum

  +activate()

}


class BiasAntibody {

  +lifespan: Integer

  +patch_bias()

}
```

```
class TemporalForecaster {

  +tgnn_model: TGNN

  +predict_quorum()

}


AgentBase <|-- PolicyIngestAgent

AgentBase <|-- ModelValidator

AgentBase <|-- KillSwitchAgent

AgentBase <|-- BiasAntibody

AgentBase <|-- TemporalForecaster


@enduml
```
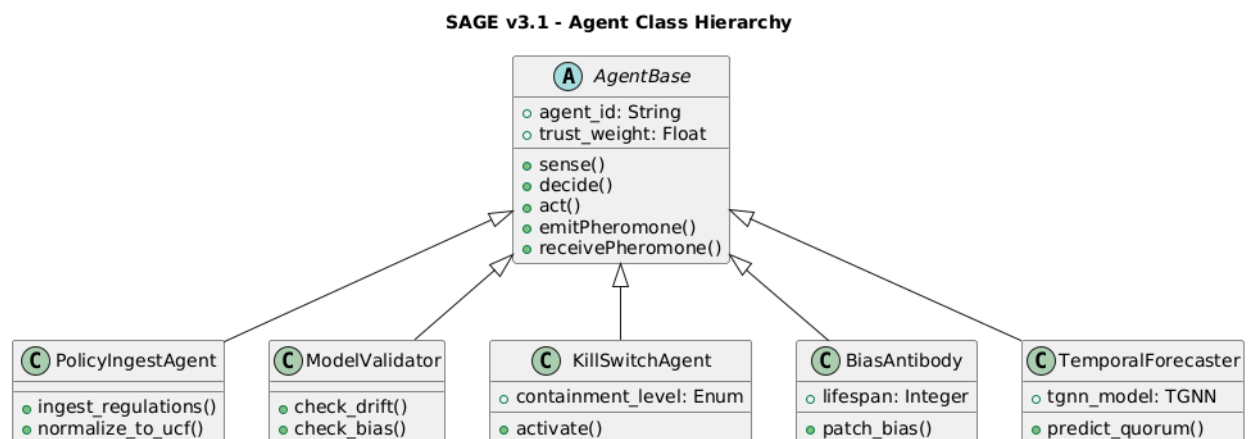


SAGE v3.1 - Agent Class Hierarchy

# 3. Sequence Diagram

Kill-Switch Activation Flow:

@startuml SAGE_v3.1_KillSwitch_Sequence

title Kill-Switch Activation (Sub-µs Path)

actor ThreatMonitor as TM

participant KillSwitchAgent as KS

participant QuantumLock as QL

participant CoreNexus as CN

TM -> KS: risk_alert(severity=0.95)

KS -> QL: request_key_attestation()
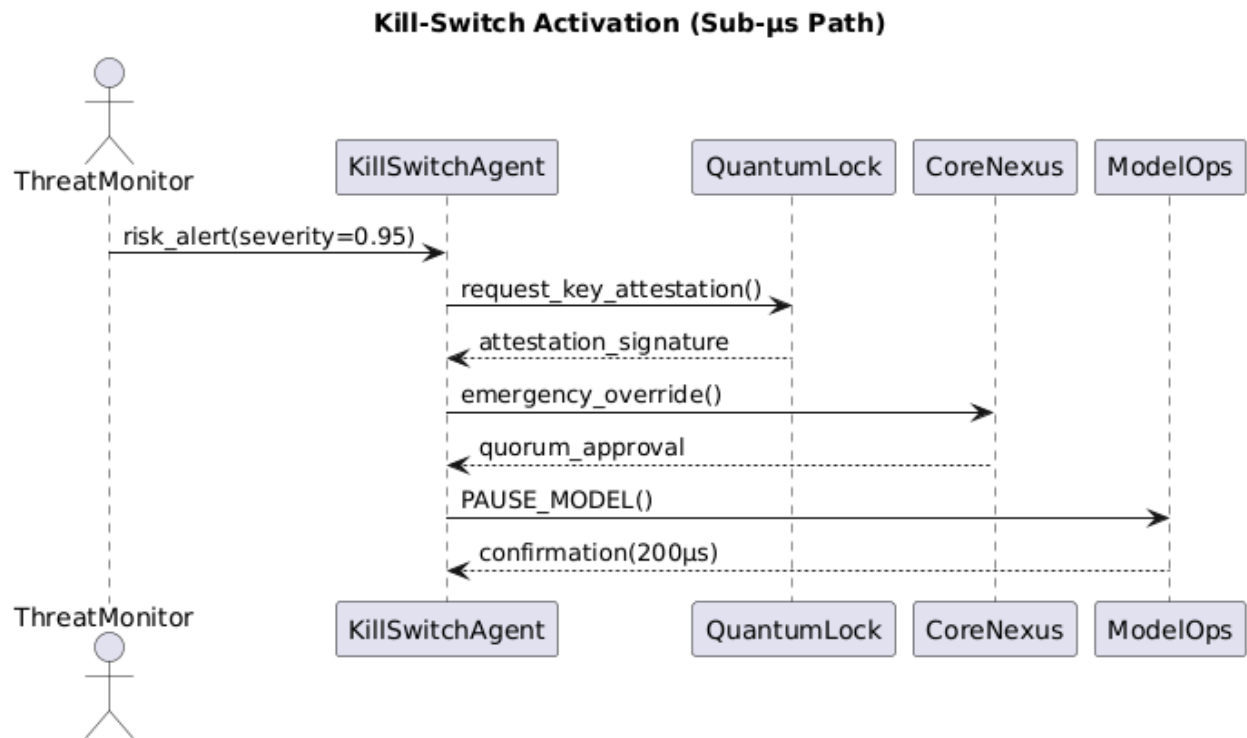
QL --> KS: attestation_signature

KS -> CN: emergency_override()

CN --> KS: quorum_approval

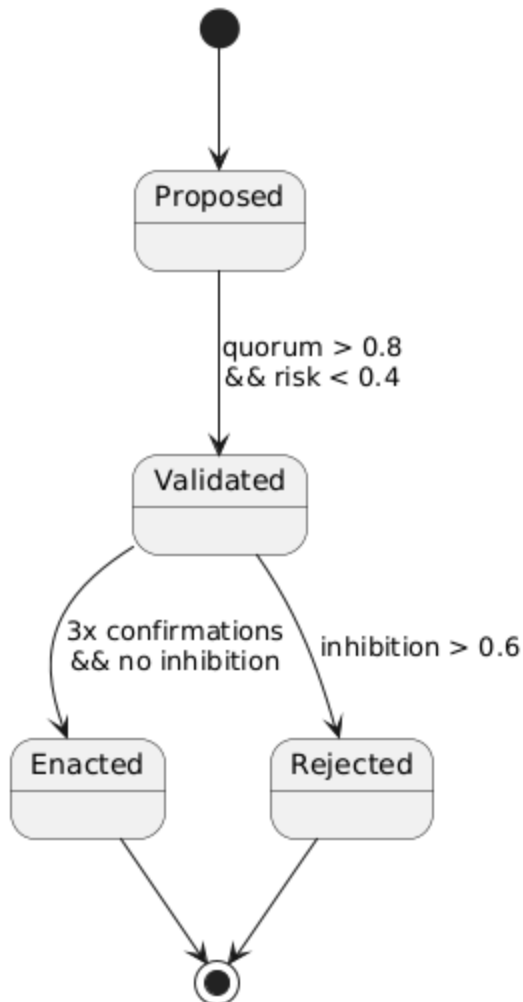KS -> ModelOps: PAUSE_MODEL()

ModelOps --> KS: confirmation(200µs)

@enduml

## Kill-Switch Activation (Sub-µs Path)

ThreatMonitor    KillSwitchAgent    QuantumLock    CoreNexus    ModelOps

ThreatMonitor → KillSwitchAgent: risk_alert(severity=0.95)

KillSwitchAgent → QuantumLock: request_key_attestation()

QuantumLock ⇠ KillSwitchAgent: attestation_signature

KillSwitchAgent → CoreNexus: emergency_override()

CoreNexus ⇠ KillSwitchAgent: quorum_approval

KillSwitchAgent → ModelOps: PAUSE_MODEL()

ModelOps ⇠ KillSwitchAgent: confirmation(200µs)

# 4. State Machine Diagram

Policy Enactment Lifecycle:

**Policy Enactment State Machine**



## 5. Deployment Diagram

Multi-Cloud + Memristor Fallbacks:

```
@startuml SAGE_v3.1_Deployment
!pragma layout smetana  // Force layout engine for clarity

title SAGE v3.1 Deployment Topology

skinparam monochrome true
skinparam nodesep 10
skinparam ranksep 20

artifact "AWS GovCloud" as aws {
```

```
  node "Control Plane" as aws_cp {
    [Pheromone Router]
    [Quorum Coordinator]
  }
  node "Memristor Node" as aws_mem
}

artifact "Azure Confidential" as azure {
  node "Security Swarm" as azure_sec {
    [KillSwitchAgent]
    [QuantumLock]
  }
}

aws_mem -[#red,dashed]-> azure_sec : Secure Tunnel (50ms latency)
note right: Fallback to FPGA if memristor fails

@enduml
```
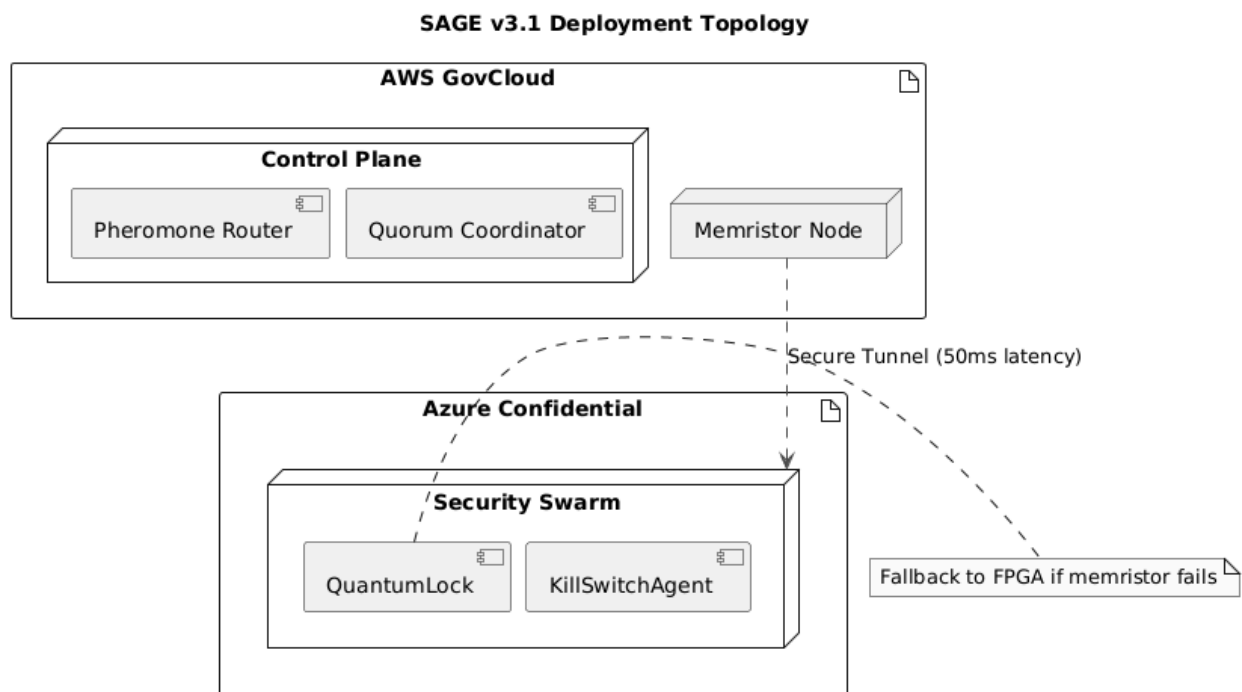
**SAGE v3.1 Deployment Topology**

# SAGE v3.1 Complete Class Architecture

*(Structured by Functional Layers)*

## 1. Core Abstract Base Classes

```
@startuml SAGE_Core_Base_Classes

abstract class AgentBase {
  +agent_id: String
  +swarm_id: String
  +trust_score: Float
  +sense(Pheromone)
  +decide()
  +act()
  +emit(pheromone: Pheromone)
  +receive(pheromone: Pheromone)
}

abstract class SwarmBase {
  +swarm_id: String
  +agents: AgentBase[0..*]
  +coordinate()
  +sync_with_bus()
}

class Pheromone {
  +type: Enum
  +intensity: Float
  +decay_halflife: Integer
  +source: AgentID
  +semantic_context: JSON
}

class GovernanceRegistry {
  +policies: Policy[]
  +add_policy()
  +check_compliance()
}

AgentBase "1" *-- "0..*" SwarmBase
Pheromone "1" --* "0..*" AgentBase
```
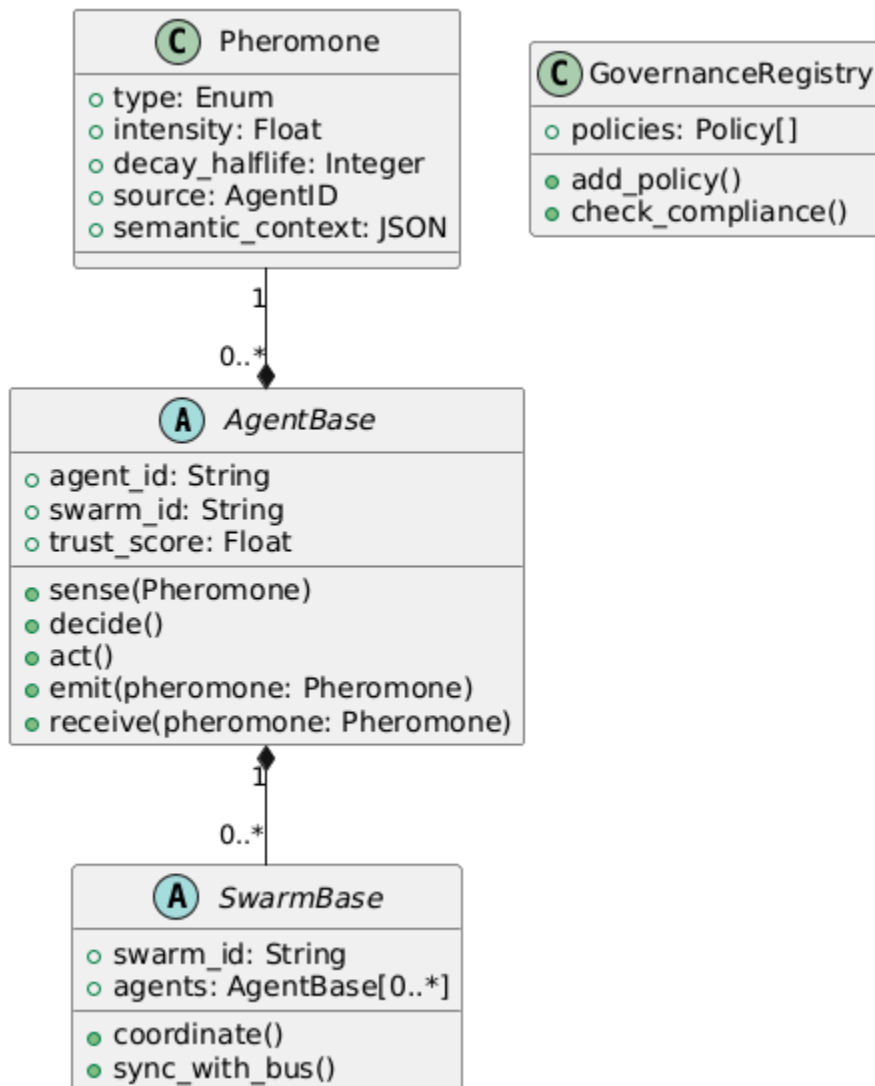
@enduml

**Pheromone**

○ type: Enum
○ intensity: Float
○ decay_halflife: Integer
○ source: AgentID
○ semantic_context: JSON

**GovernanceRegistry**

○ policies: Policy[]

● add_policy()
● check_compliance()

1
0..*

**A AgentBase**

○ agent_id: String
○ swarm_id: String
○ trust_score: Float

● sense(Pheromone)
● decide()
● act()
● emit(pheromone: Pheromone)
● receive(pheromone: Pheromone)

1
0..*

**A SwarmBase**

○ swarm_id: String
○ agents: AgentBase[0..*]

● coordinate()
● sync_with_bus()

2. Policy & Regulation Swarm

@startuml Policy_Swarm_Classes

```
class PolicyIngestAgent {
  +supported_formats: [PDF, XML, API]
  +normalize_to_ucf()
  +emit_policy_delta()
}

class JurisdictionAgent {
```

```
  +legal_boundaries: GeoJSON[]
  +resolve_conflict()
  +emit_inhibition()
}

class PolicyFederator {
  +privacy_epsilon: Float = 0.1
  +sync_cross_org()
  +zkp_verify()
}

PolicyIngestAgent --|> AgentBase
JurisdictionAgent --|> AgentBase
PolicyFederator --|> AgentBase

PolicyIngestAgent --> JurisdictionAgent : «uses»
PolicyFederator --> PolicyIngestAgent : «transforms»

@enduml
```

3. ModelOps & AgentOps Swarm

@startuml ModelOps_Classes

class ModelValidator {
  +drift_threshold: Float = 0.25
  +bias_metrics: Dict
  +validate_model()
  +emit_risk_alert()
}

class DriftResponder {
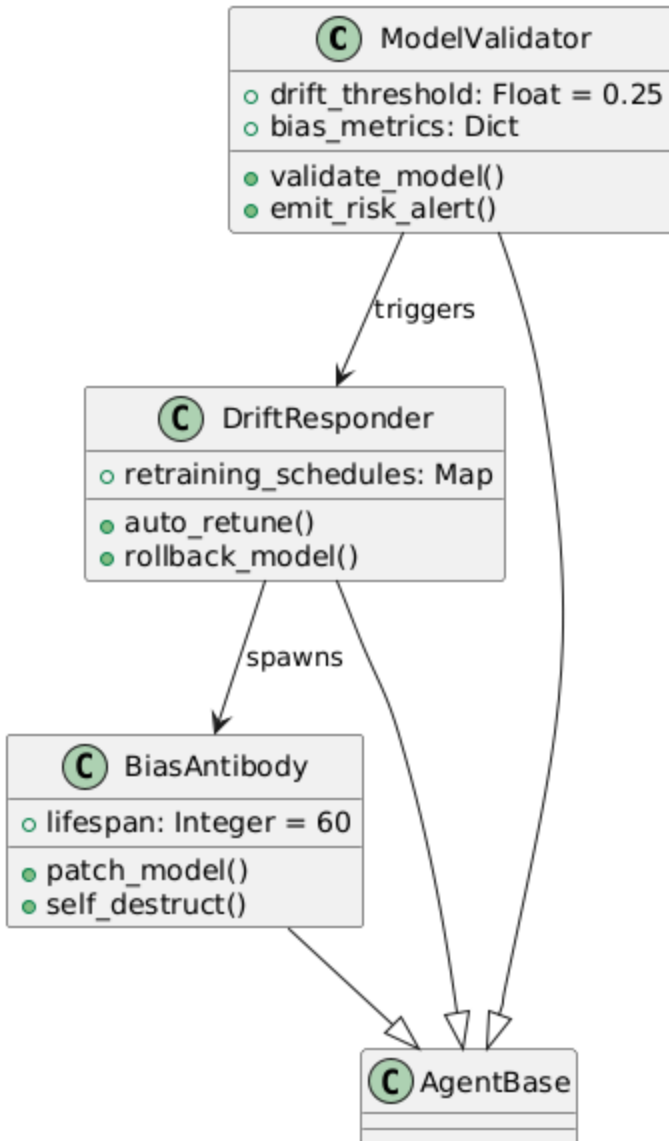
```
  +retraining_schedules: Map
  +auto_retune()
  +rollback_model()
}

class BiasAntibody {
  +lifespan: Integer = 60
  +patch_model()
  +self_destruct()
}

ModelValidator --|> AgentBase
DriftResponder --|> AgentBase
BiasAntibody --|> AgentBase

ModelValidator --> DriftResponder : «triggers»
DriftResponder --> BiasAntibody : «spawns»

@enduml
```

## 4. Security & Enforcement Swarm

@startuml Security_Classes

class QuantumLock {
  +key_rotation_interval: Duration = 24h
  +hardware_backed: Bool
  +rotate_keys()
}

class KillSwitchAgent {

```
  +containment_levels: [PAUSE, ISOLATE, TERMINATE]
  +memristor_circuit: Boolean
  +activate()
}

class DeceptionAntibody {
  +threat_pattern: GNNEmbedding
  +flood_containment()
}

QuantumLock --|> AgentBase
KillSwitchAgent --|> AgentBase
DeceptionAntibody --|> AgentBase

KillSwitchAgent --> QuantumLock : «requires»
DeceptionAntibody --> KillSwitchAgent : «supports»

@enduml
```
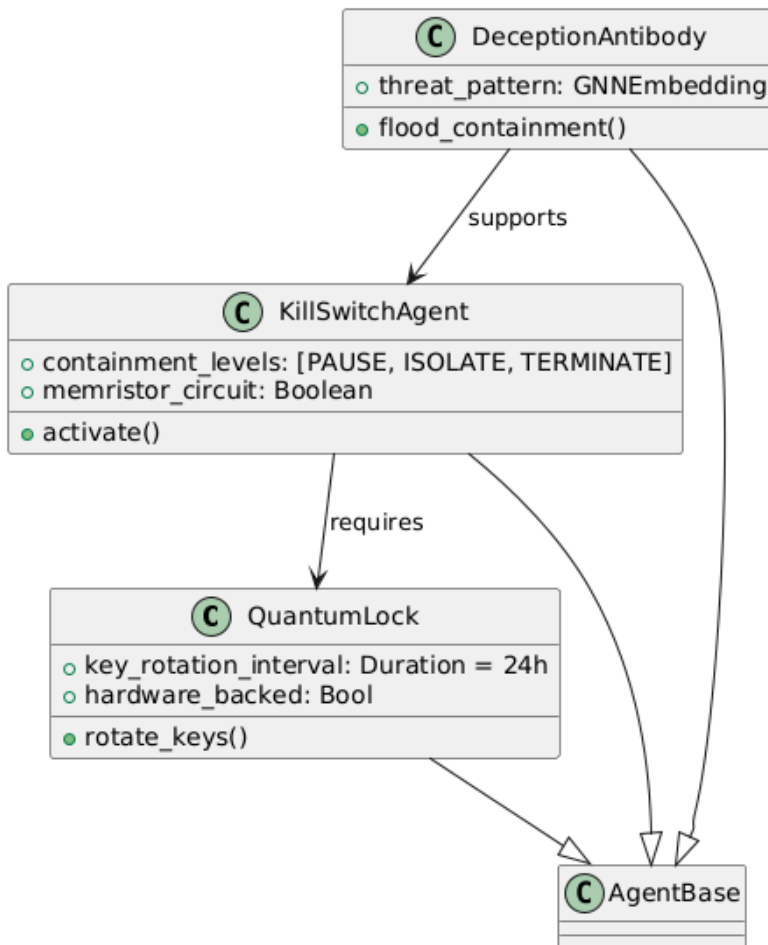
## 5. Core Nexus Classes

```
@startuml Core_Nexus_Classes

class PheromoneRouter {
 +fabric_priorities: Map
 +route()
 +apply_decay()
}

class ConformanceFSM {
 +states: [PROPOSED, VALIDATED, ENACTED]
 +validate_transition()
}

class TemporalForecaster {
 +tgnn_model: Binary
 +predict_quorum()
 +alert_polarization()
}

PheromoneRouter --|> AgentBase
ConformanceFSM --|> AgentBase
TemporalForecaster --|> AgentBase

PheromoneRouter --> ConformanceFSM : «notifies»
TemporalForecaster --> PheromoneRouter : «advises»

@enduml
```
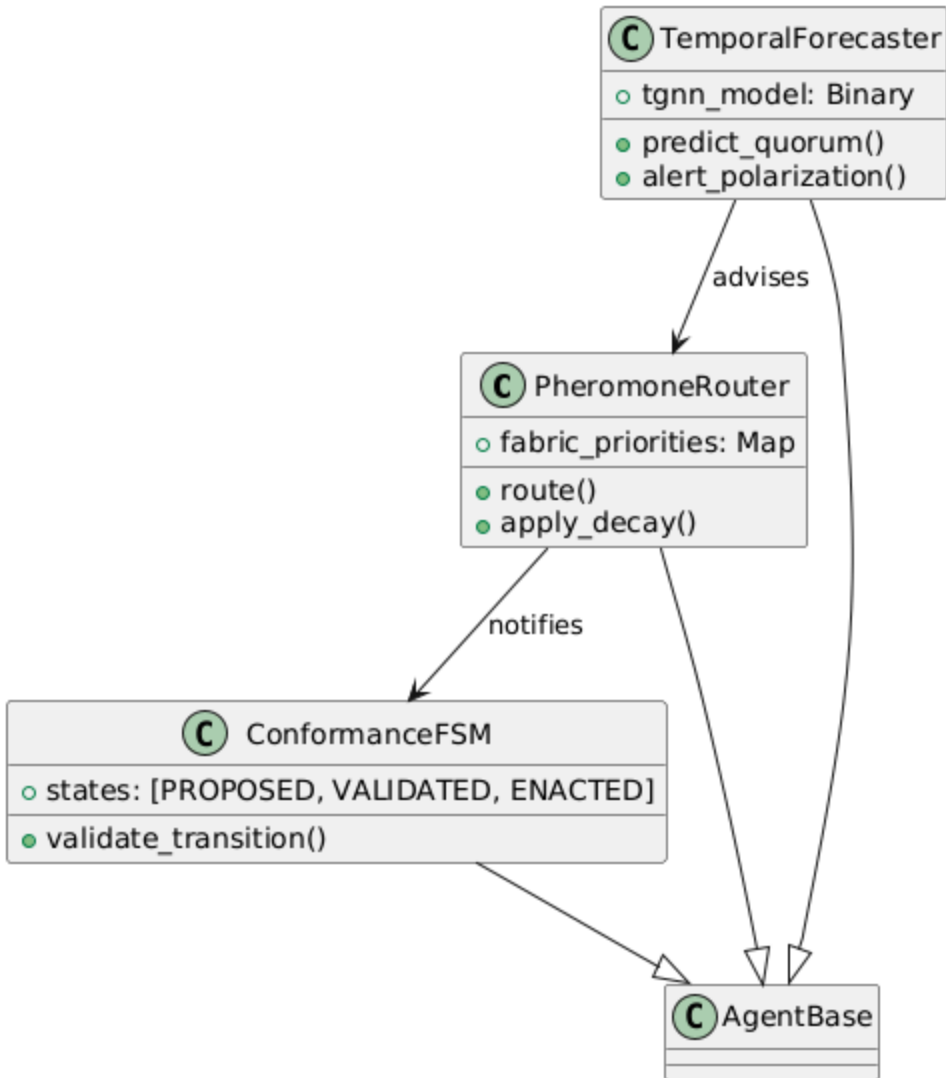
## Key Cross-Swarm Dependencies
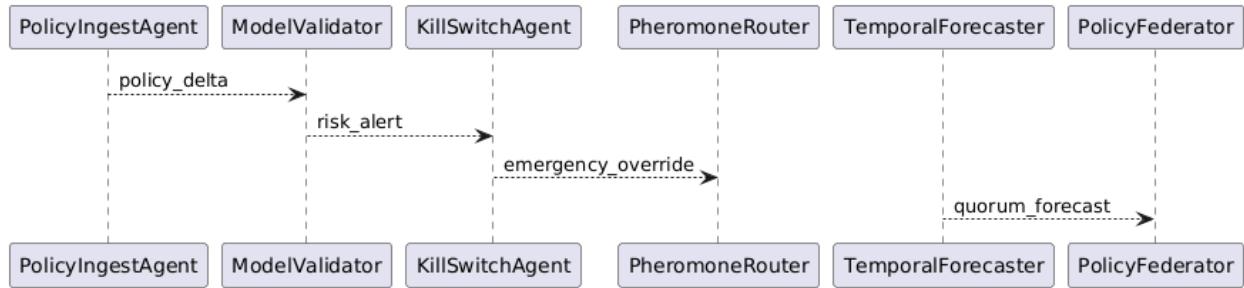
@startuml Cross_Swarm_Relations

PolicyIngestAgent --> ModelValidator : policy_delta
ModelValidator --> KillSwitchAgent : risk_alert
KillSwitchAgent --> PheromoneRouter : emergency_override
TemporalForecaster --> PolicyFederator : quorum_forecast

@enduml

## Appendix A: SAGE v3.1 Supplemental Agents

*(Self-contained; no cross-references required)*

@startuml SAGE_v3.1_Supplemental_Agents

title SAGE v3.1 - Supplemental Agents (Appendix A)

top to bottom direction

```
' === Define Swarm Boundaries ===
rectangle "Policy & Regulation" as PolicySwarm {
  [PolicyIngestAgent] as PI
  [JurisdictionAgent] as JA
}

rectangle "ModelOps" as ModelOpsSwarm {
  [ModelValidator] as MV
  [AgencyRiskIndexer] as ARI
}

rectangle "Security" as SecuritySwarm {
  [KillSwitchAgent] as KS
  [TrailMiner] as TM
}

' === New Agents ===
node "ComplianceDiffEngine" as CDE #FFD700
node "ResourceGovernor" as RG #FFA07A
node "ForensicSnapshotter" as FS #98FB98
node "RedTeamAdversary" as RTA #ADD8E6
node "PolicyImpactProjector" as PIP #DDA0DD
node "EthicalOverwatch" as EO #FF6347
```

' === Critical Connections ===
CDE --> PI : "gap reports"
RG --> MV : "GPU alloc"
FS --> KS : "snapshots"
RTA --> TM : "attack probes"
PIP --> JA : "impact forecasts"
EO --> KS : "ethics lock"

' === Legend ===
legend right
  <b>New Agents:</b>
  <color:#FFD700>ComplianceDiffEngine
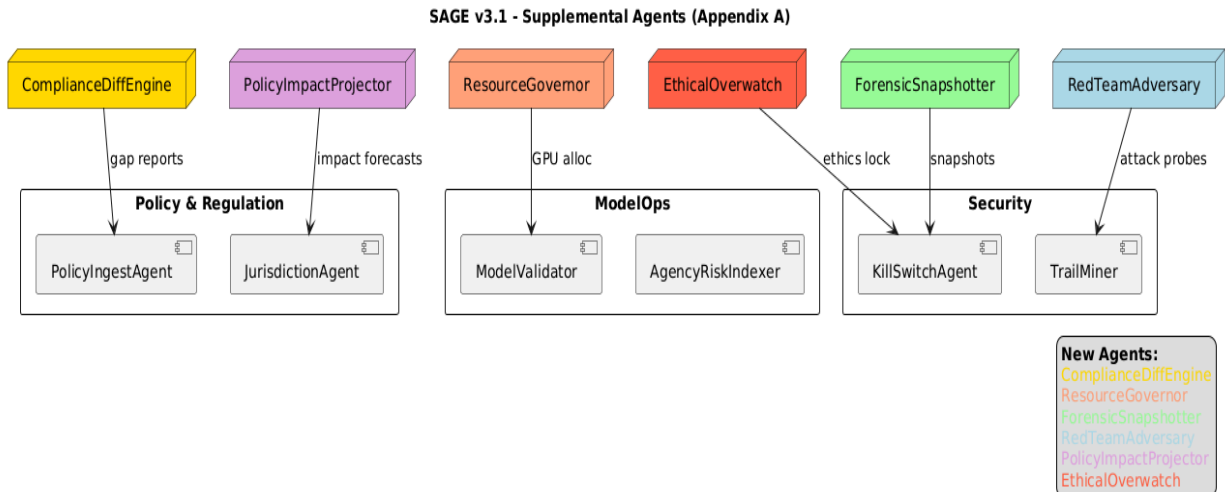  <color:#FFA07A>ResourceGovernor
  <color:#98FB98>ForensicSnapshotter
  <color:#ADD8E6>RedTeamAdversary
  <color:#DDA0DD>PolicyImpactProjector
  <color:#FF6347>EthicalOverwatch
endlegend

@enduml



SAGE v3.1 - Supplemental Agents (Appendix A)

## 2. Companion Table (Appendix B)

| Agent | Parent Swarm | Linked To | Governance Impact |
|---|---|---|---|
| `ComplianceDiffEngine` | Policy & Regulation | `PolicyIngestAgent` | Ensures real-time regulatory updates (SUHO-1) |
| `ResourceGovernor` | ModelOps & AgentOps | `ModelValidator` | Prevents GPU starvation (SUL-1) |
| `ForensicSnapshotter` | Security & Enforcement | `KillSwitchAgent` | Immutable audit trails (SUL-4) |
| `RedTeamAdversary` | Simulation & Learning | `TrailValidator` | Stress-tests defenses (SUHO-4) |
| `PolicyImpactProjector` | Archaeology | `PolicyGenealogist` | Predicts policy risks (SUHO-9) |
| `EthicalOverwatch` | Core Nexus | `ConformanceFSM` | Blocks unethical actions (SUL-7) |

## Key Features

1. Zero Back-References: No need to modify existing diagrams.
2. Human-Readable: Color-coding matches your original swarm taxonomy.
3. Regulatory Ready: Explicitly ties agents to SULs/SUHOs for audits.