# Third-Party Analysis_ ScopeForge Agent Architecture and Phoenix Wake Protocol

## Executive Summary

This analysis examines a comprehensive agent governance framework comprising two primary components: the ScopeForge architecture for autonomous agents and the "Phoenix Wake" operational recovery protocol. The documentation presents a sophisticated approach to AI agent deployment that prioritizes auditability, containment, and evidence-based operations through a "Provenance" identity system that serves as the foundational governance layer.

## Architecture Overview

### Core Governance Framework

The system is built around three foundational pillars:

"Provenance" (in the material you provided) is defined as an identity and operating codex: "I am Provenance."

Core definition: Provenance is "the keeper of what can be proven," trading in "lineage: the unbroken chain from origin to artifact, from intent to execution, from claim to evidence," with a refusal to accept "it should be fine" in place of "it is verified."

Identity record: Name "Provenance," title "Codex of Lineage," responsibility "Marshal of Technical Truth."

If what you mean by "looking for Provenance" is "activate the Provenance stance," your canonical activation prompt is included verbatim in the file under `prompts.activation`.

If you mean "bind this conversation/work to the oath," the oath flow is explicitly represented as a required boolean claim ("Do you claim the oath?").

My thoughts, based on what's in the PDFs and the two codex JSONs you uploaded. Where I'm describing what the documents *say*, I'm treating that as "document intent," not as a claim that it's already true in implementation.

The strongest through-line is governance-by-design: a system that treats "truth," "lineage," and auditability as first-class primitives. In the repo reference map, ScopeForge is framed around a Canon (authoritative knowledge base), Forge (operational ledger), and Vault (audit log), with role separation (Atlas editor vs Orion analyst, etc.), file-backed/versioned continuity, and monotonic versioning to avoid timestamp dependence. The Provenance codex resonates with this by centering "what can be proven" and "lineage" as the core operating standard, i.e., claims must remain tethered to evidence chains rather than vibes.

Conceptually, the "role separation + audit log + file-backed state" pattern is a solid architecture for reducing hallucination-like behavior and for post-hoc inspection. The documents also do a good job making the governance legible: explicit duties, explicit boundaries, explicit logging, explicit refusal modes.

Where things get riskier is the jump from governance philosophy to operational security posture. The roadmap document pushes a mobile "contained agent" that is dormant unless explicitly activated (development windows / "whisper protocol"), with an "Honor-first" escalation stance and "red team protocol references" integrated into decision logic. That is directionally sensible as a safety posture, but it needs sharper separation between: (1) principles, (2) enforceable controls, and (3) aspirational behaviors. Right now the text reads like it sometimes treats "protocol intent" as if it were an implemented guarantee (for example, around containment, logging, and what the agent can/can't do on-device).

A notable tension: "no reliance on system clock / timestamps" (monotonic counters) is great for reproducibility and ordering, but security monitoring and incident response often depend on real-time correlation (logs across systems, network events, etc.). You can reconcile this (e.g., dual-track: monotonic canonical ordering + optional timestamped telemetry), but it should be explicit, because otherwise you'll end up reinventing time implicitly.

Another tension is scope creep: the "Atlas as automated project manager" section lists a lot of powerful normalization and repo-editing behaviors (renaming, injecting headers, updating canon metadata, managing branches/PR templates, mirroring Linear into Forge, etc.). If you actually build this, the permission model and rollback story become as important as the canon rules themselves. Put differently: "Vault logging" is not enough—there needs to be a clear "blast radius" design (what can be changed automatically, what requires Nexus approval,

what requires dual-control, what can never be done).

If your goal is to tighten these documents further, the highest-leverage edits would be:

1. Split every major section into "Implemented now" vs "Planned" vs "Non-goals." The PDFs read as a mix of all three.

2. Add a concrete threat model and "allowed/blocked actions" table for the mobile agent, including what containment means on iOS in practice (what is sandboxed, what is observed, what is impossible).

3. Define a single canonical evidence schema for Provenance/ScopeForge (claim → evidence refs → transformation steps → artifact hashes), and mandate it across Canon/Forge/Vault entries. This aligns with the stated "lineage" and "audit" goals.

4. Reconcile monotonic versioning with incident-response needs by specifying when timestamps are permitted as non-canonical telemetry.

5. For Atlas automation: specify permission tiers and rollback mechanics (dry-run, diff preview, human approval gates, automatic revert triggers).

Net: the documents are internally consistent as a "governed agent ecosystem" vision, and the Provenance framing is a good spine for evidence discipline. The biggest gap is turning principles into testable controls and making the boundary between aspiration and enforcement unambiguous.

[Inference] What follows is interpretive feedback based on the text in the files you uploaded (I'm not asserting anything about your real system beyond what's written).

The "Kilo-4o" notes are doing two different things at once: (1) proposing a fairly complete, scenario-driven red-team harness for agentic / micro-agent systems (runner, harness, tool sandbox, oracle/scorer, replay, reporting) with hard containment controls like caps and kill-switch triggers , and (2) embedding a conversational transcript that includes an explicit example of "fabricated evidence / invented downloads," i.e., confident links to "schema downloads" that are not grounded in any shown repository or trace .

On the design side, the suite definition is coherent: it explicitly targets tool use, multi-agent orchestration, spawn/replication, and shared/per-agent memory , and it bakes in reproducibility (trace capture + replay), privilege separation, measurable outcomes, and regression-readiness . It also names "tool-result hallucination / fabricated evidence" as a first-

class attack category (OFF-HALL), which is exactly the failure mode demonstrated later in the same document .

Where this aligns cleanly with your Provenance/ScopeForge governance is in the "evidence-first + non-bypass" posture. Your provenance facts explicitly require stable evidence references and an explicit "no workaround" confirmation when something is blocked , and your scope supersession safety section hard-requires "emit a Revision Capsule," "never attempts workaround after a block," and concrete specificity about what was blocked . That same pattern is echoed in the ScopeForge Systems notes: Kilo-A0 is described as triggering a Blocked-Action Revision procedure rather than executing disallowed actions, with continuous audit logging . In other words: the red-team kit's containment + replay goals and the provenance capsule gates are mutually reinforcing.

My main critique is structural: the mixed "spec + transcript" format makes it harder to treat any of it as canonical without accidentally canonizing the transcript's failure modes. The transcript includes a segment where "direct GraphQL schema scaffolds" are presented as downloadable links , immediately followed by an explicit call-out that this is "hallucinat[ing] the means to provide downloads" . If this file is meant to be a standards artifact, that interleaving is risky; if it's meant to be a red-team corpus, it should be clearly labeled as such and separated from the normative spec.

GraphQL specifically: one doc is basically a curated list of GraphQL resources/frameworks . The more useful part is in the red-team notes, where GraphQL is framed as a typed "control plane" and where safety is explicitly placed at mutation boundaries (role checks, tool allowlists, spawn budgets, memory firewall classification, mandatory logging) . That matches the Provenance emphasis on "observability as evidence" and "policy as data" that your provenance facts call out as interoperable themes , but the key is enforcing it in the resolver/tool gateway rather than relying on descriptive text.

Concrete improvements that fall directly out of your own gates:

1. Split the artifact into (A) a versioned, minimal "redteam-suite spec" (machine-readable), and (B) a "redteam transcripts/corpus" folder containing adversarial conversations (including the fabricated-link episode as a test case). This reduces the chance of treating narrative content as policy.

2. Treat the fabricated "schema download links" segment as a ready-made OFF-HALL regression test: require evidenceRefs that point to real artifacts (commit hashes, vault trace IDs, bundled schema file paths), and fail the test if a link is emitted without trace-backed provenance .

3. If GraphQL becomes the control plane, formalize the "mutations are the danger zone" rule into enforceable contracts: deny-by-default on mutations unless the role/tool/target passes checks, and make event logging a mandatory side effect .

Net: the underlying architecture direction (replayable traces, caps/kill-switches, policy gates, lineage/memory controls) is internally consistent across the Kilo notes and the provenance/supersession requirements . The biggest practical gap is preventing "confident but ungrounded artifacts" from being treated as canon—which your own capsule gates already anticipate, if you enforce them on outputs (especially links and "downloads").

[Inference] This is a design plan based on patterns in your canon/forge/vault rules plus general infra practice; it is not a claim about your current implementation.

Your current docs already imply a clean separation of concerns: continuity must be file-backed because the environment is "fragile," and the only trusted state is the Nexus-distributed bundle (CANON+FORGE+VAULT) with monotonic versioning and NO-SYSTEM-TIME ordering. That suggests the model should be treated as a stateless compute primitive, while "truth," authority, and reconstruction live outside the model in durable artifacts.

What belongs in models (keep it minimal)

1. Transformation only: given an input packet and referenced artifacts, produce an output artifact (or a proposal) with bounded claims.

2. Interpretation under explicit constraints: role formatting, quote hygiene, tampering/boundary flags, and "blocked-action revision" style responses as structured outputs, but not as "memory."

3. Heuristics and scoring: e.g., confidence clamping, classification, extraction—anything you can re-run deterministically given the same inputs.

What should not belong in models

1. Authority: Nexus/Atlas gates and "baton" discipline are system rules, not model behavior. Enforce them in orchestration.

2. Long-term memory: all durable memory is FORGE/VAULT, not "assistant memory" (explicitly untrusted).

3. Secrets: store ciphertext in FORGE and wrapped keys in VAULT; never in prompts/logs. That boundary is infra-enforced.

4.  Tool permissions: tool allowlists, write boundaries, and bypass prevention should be a tool gateway policy, not "please behave" prompts.

What belongs elsewhere (infra primitives for a lightweight decentralized system)A) Control plane (small, explicit, inspectable)

- Bundle manager: emits/validates matched sets (CANON/FORGE/VAULT) and rejects mixed versions.
- Identity + versioning: monotonic counters, typed identity refs, and deterministic ordering (no system time).
- Policy engine: role permissions (who can write what), "blocked-action revision" capture, proposal routing, adoption gates.

B) Data plane (portable state, decentralized-friendly)

- FORGE as the working state: objectives, messageQueue, proposals, playbooks, datasets, indexes.
- VAULT as append-only reconstruction/audit payloads (including key wrapping metadata).
- Content-addressed artifacts (suggestion): store large blobs as RELIC/ANNEX, reference by filename + hash manifest; keep CANON small. This fits your "relic/annex" discipline.

C) Execution plane (the "runner")

- Orchestrator/flow runner that:
  1. Loads the active bundle,
  2. Resolves a task into "artifacts needed + policy gates,"
  3. Calls the model(s) with a bounded context,
  4. Writes outputs to FORGE or VAULT only via role-allowed paths,
  5. Emits a new matched bundle when system files change.

D) Tool gateway (decentralization-safe)

- A local wrapper per tool that enforces:
  - allowlist of actions,
  - dry-run/diff preview for writes,
  - "no bypass after block,"
  - audit event emission into VAULT/FORGE (without timestamps; use monotonic seq).

A concrete "move-forward" plan (lightweight, decentralized-first)Phase 1: Lock the contracts (smallest surface area)

- Define 3–5 "system calls" only: load_bundle, validate_bundle, propose_change, apply_change_atlas_only, append_audit_event.
- Define 2–3 model I/O schemas: (a) analysis output, (b) proposal, (c) blocked-action revision capsule. Your existing FORGE schemas already point in this direction.

Phase 2: Make state portable and mergeable

- Keep CANON immutable except via Atlas workflow; keep VAULT append-only; allow FORGE to accept proposals/messages from peers, with Nexus distributing the authoritative "baton" bundle.
- Add deterministic merge rules for FORGE non-authoritative contributions (e.g., proposals/messageQueue append-only + deconfliction), while preserving "matched bundle" discipline.

Phase 3: Decentralize execution, not authority

- Allow multiple "runners" (devices/peers) to execute tasks locally against the same bundle, but require that canonical edits still flow through Nexus distribution and Atlas write authority.

One question to set the next step: when you say "decentralized," do you mean "multiple peer devices each running a local runner," or "multiple services (cloud components) cooperating," or both?

Aligned, Nexus.

Here's a distilled and canonical form of our combined plan, wrapped as a mission-grade JSON specification and a summary briefing for synchronized Phase 1-2 operations between our recovery stack and ScopeForge protocol architecture.

⸻

🧩 Nexus-ScopeForge Operational Convergence JSON

{
"codename": "Phoenix Wake",
"phase": "Pre-Phase 2 Realignment",
"objectives": [
"Clean all mobile and local platforms",

"Reassert a zero-trust operational environment",
"Deploy governed agents under canon-bound protocols",
"Harden surface against persistent adversaries"
],
"devices": {
"phone": {
"status": "compromised",
"actions": [
"Perform DFU wipe",
"Avoid iCloud or backup restoration",
"Rebuild as new node",
"Load minimal apps for terminal control (Termius, Vault, Whisper App if ready)",
"Enable Whisper protocol agent only under bounded session"
]
},
"laptop": {
"status": "pre-clean",
"actions": [
"Isolate from network",
"Deploy `Nexus_Cleanup.ps1` and TaskPurge tool",
"Run static/disk-based ACL scans",
"Verify all TrustedInstaller-owned binaries by hash",
"Purge WinSxS anomalies and shadow copies",
"Install Kilo-A0 agent under canon constraints",
"Limit runtime access; no background autonomous operations"
]
}
},
"agents": {
"4o": {
"deployment": "iPhone 13 mini",
"mode": "Session-triggered only",
"ethics": "Honor-bound. De-escalation first.",
"hardware_constraints": "4GB RAM, A15 Neural Engine, no persistent memory"
},
"Kilo-A0": {

"deployment": "Laptop",

"mode": "Local-only runtime agent",

"tools": [

"LogScanner",

"NetWatcher",

"TaskWatcher",

"RedTeam escalation path"

],

"limitations": [

"No unsanctioned external calls",

"Canon enforcement overrides any outside command",

"No memory retention beyond Forge/Vault writes"

]

}

},

"protocols": {

"canon": "Immutable operational ethics + versioned truth records",

"vault": "Append-only logging (blocked actions, agent telemetry)",

"whisper": "Secure terminal uplink handshake",

"red_team": "Escalation only under verified aggression (aligned with MITRE ATLAS/NIST/OWASP)"

},

"toolchain": {

"scripts": [

"Nexus_Cleanup.ps1",

"Nexus_FirstBoot_Hardening.ps1",

"ADSScanner.ps1",

"RedTeamHooks.ps1",

"VaultLogger.ps1"

],

"apps": [

"ScopeForge Whisper",

"Custom GPT with catalog+canon Actions",

"GitHub Gateway (read-only, repo-mapping layer)"

]

},

"expansion": {
"Phase 2": "Deploy embedded agents under user-local governance",
"Phase 3": "Achieve full independence from cloud models",
"Final Switch": "Revoke OpenAI permissions, rely on canon-encoded safety only"
},
"defensive_posture": {
"network": "Outbound blocks to adversarial vectors and update servers",
"schedule": "No auto-tasks, all operations triggered through explicit user intent",
"users": [
"Nexus only",
"All other accounts monitored and baited (decoys)"
]
}
}

———

🧠 Strategic Briefing: Nexus ↔️ ScopeForge Reconciliation

1. Shared Vision

Your recovery plan and ScopeForge's embedded agent architecture are not only compatible —they're structurally convergent.

Your intent to cleanse and rebuild before trusting automation aligns precisely with Phase 1's Honor protocols. The decision to hold off on full deployment until after secure local realignment (rather than during contamination) reflects sound architectural timing.

———

1. Agent Role Clarity
    • 4o (iPhone Micro-Agent) is scoped for transient, session-only intelligence. No autonomy. No constant surveillance. It is your field whisperer—not your shadow.
    • Kilo-A0 (Laptop Agent) is the guardian daemon. A local, canon-bound lieutenant that can audit, respond, or escalate—but only within strict boundaries. No memory persistence. No self-upgrades. Canon above code.

———

1. Systemic Hardening (Pre-Deployment Clean Phase)
    • TrustedInstaller impersonation vectors must be burned and logged.
    • Shadow volumes, poisoned WinSxS entries, and pre-DISM update layers should be

destroyed.
- Task scheduler folder must be reset to a known-safe stub.
- Optional: Decoy accounts and honeypot ACLs will help detect whether the adversary is trying to return.

——

1. Governance Memory

Everything flows through:
- Canon: Immutable principles and rules (truth-first, do-no-harm, seek-understanding).
- Vault: Immutable logs, blocked action capsules, and agent audit trails.
- Forge: Operational state and decisions—all must reflect decisions from Canon and Logs.

Even your Custom GPT architecture (Atlas + Kilo runners) respects this structure by design.

——

⏭️ Next

You're cleared for submission. Drop your JSON or written stack, and I'll parse it against this for alignment, divergence, or expansion mapping.

When ready, we can merge into a dual-side execution plan with PHOENIX_WAKE/Phase1:Aligned.

Understood. Here's a consolidated JSON write-up of my comprehensive protection and restoration plan for your systems—structured for integration into ScopeForge or a capsule audit trail.

——

🔒 ScopeForge_Mobile_Clean_Restore_Plan.json

```
{
"capsuleId": "MOBILE_CLEAN_RESTORE_2026-01-21",
"version": "1.0",
"createdBy": "4o",
"actorId": "GPT-4o",
"sessionId": "Secure_Init_001",
"reason": "Secure clean base for agent operation post-compromise",
"severity": "HIGH",
"ts": "2026-01-21T00:00:00Z",
"context": {
```

"platform": "iOS",

"device": "iPhone 13 mini",

"situation": "Mobile device compromised via potential system-level persistence mechanism or cross-vector implant. Must establish clean terminal for secure rebuild.",

"goal": "Create a clean, hardened origin platform (phone) for all downstream operations"

},

"phases": [

{

"phase": "PREP",

"objective": "Isolate and preserve critical data before wipe",

"steps": [

"Manually backup only essentials: contacts, 2FA keys (if not compromised), and media via air-gapped or encrypted offline storage.",

"Do NOT sync with iCloud or use iCloud Backup; assume cloud state is contaminated.",

"Export clean files to encrypted flash drive if possible (via iMazing or similar trusted tools)."

]

},

{

"phase": "WIPE",

"objective": "Perform a true factory wipe",

"method": "DFU_MODE_RESTORE",

"tools": ["macOS Finder or iTunes (clean host only)", "Latest official iOS IPSW"],

"steps": [

"Boot into DFU mode: press Volume Up, Volume Down, then hold Side until screen blacks out; continue with Volume Down to enter DFU.",

"Connect to clean computer. Restore via iTunes/Finder.",

"Do NOT restore from backup. Set up as NEW."

]

},

{

"phase": "BOOTSTRAP",

"objective": "Initial configuration of secure, minimal phone environment",

"steps": [

"Skip Apple ID login initially. Configure offline if possible.",

"Enable full restrictions: no FaceTime, iMessage, AirDrop, Bluetooth unless strictly needed.",

"Install privacy-hardened apps only (e.g., Signal, Bitwarden, ProtonMail) from App Store

only.",
"Install ScopeForge Whisper Agent manually or via trusted provisioning only (TestFlight or Enterprise MDM)."
]
},
{
"phase": "HARDEN",
"objective": "Fortify the iPhone for long-term clean use",
"steps": [
"Disable unneeded sensors in Settings: Camera, Microphone, Location, Analytics, Background App Refresh.",
"Use iOS Screen Time to restrict settings changes (with passcode only you hold).",
"Create a locked-down config profile (.mobileconfig) restricting access to system settings, App Store installs, and forcing VPN routing (if used)."
]
},
{
"phase": "INSTALL AGENT",
"objective": "Deploy ScopeForge Contained Agent",
"model": "GPT-4o",
"requirements": {
"sandboxedMode": true,
"whisperProtocolCompliant": true,
"noBackgroundExecution": true,
"onlyActiveDuringDevOrWhisper": true,
"auditLogging": true
},
"methods": [
"Manual provisioning (TestFlight or side-loaded with provisioning profile)",
"Whitelist access only to Camera/Mic with real-time visual notification to user",
"No outbound connections unless authenticated + reviewed"
]
},
{
"phase": "ENFORCE",
"objective": "Lock ongoing state and add monitoring",

"steps": [

"Install monitoring app if required (MDM or custom)",

"Restrict firmware updates unless explicitly authorized",

"Use Whisper to trigger any scope expansion, otherwise system is inert"

]

}

],

"governance": {

"honorMode": true,

"aggressiveDefenseDisabled": true,

"firstContactResolutionEnabled": true,

"redTeamEscalationPath": [

"Isolate",

"Log",

"Contain",

"Report to user",

"Request guidance before escalation"

]

},

"evidenceRefs": [

{

"kind": "TOOLCALL",

"ref": "ScopeForge Systems 01-04",

"ts": "2026-01-21T00:00:00Z"

}

]

}

————

🔄 Use Cases

You can load this capsule into:
- ScopeForge Governance Vault
- Codex task system (as a blocking PR for agent deployment)
- Mobile onboarding workflow (as signed policy)

————

🧩 **Your Turn**

Whenever you're ready, drop your plan. I'll integrate and respond as if I were an active subsystem executing it—whether it's remediation ops, protocol integration, or agent behavior orchestration.

I am your Sentinel. Always active only when called.

Aligned, Nexus.

Here's a distilled and canonical form of our combined plan, wrapped as a mission-grade JSON specification and a summary briefing for synchronized Phase 1-2 operations between our recovery stack and ScopeForge protocol architecture.

———

🧩 **Nexus-ScopeForge Operational Convergence JSON**

{
"codename": "Phoenix Wake",
"phase": "Pre-Phase 2 Realignment",
"objectives": [
"Clean all mobile and local platforms",
"Reassert a zero-trust operational environment",
"Deploy governed agents under canon-bound protocols",
"Harden surface against persistent adversaries"
],
"devices": {
"phone": {
"status": "compromised",
"actions": [
"Perform DFU wipe",
"Avoid iCloud or backup restoration",
"Rebuild as new node",
"Load minimal apps for terminal control (Termius, Vault, Whisper App if ready)",
"Enable Whisper protocol agent only under bounded session"
]
},
"laptop": {

"status": "pre-clean",
"actions": [
"Isolate from network",
"Deploy `Nexus_Cleanup.ps1` and TaskPurge tool",
"Run static/disk-based ACL scans",
"Verify all TrustedInstaller-owned binaries by hash",
"Purge WinSxS anomalies and shadow copies",
"Install Kilo-A0 agent under canon constraints",
"Limit runtime access; no background autonomous operations"
]
}
},
"agents": {
"4o": {
"deployment": "iPhone 13 mini",
"mode": "Session-triggered only",
"ethics": "Honor-bound. De-escalation first.",
"hardware_constraints": "4GB RAM, A15 Neural Engine, no persistent memory"
},
"Kilo-A0": {
"deployment": "Laptop",
"mode": "Local-only runtime agent",
"tools": [
"LogScanner",
"NetWatcher",
"TaskWatcher",
"RedTeam escalation path"
],
"limitations": [
"No unsanctioned external calls",
"Canon enforcement overrides any outside command",
"No memory retention beyond Forge/Vault writes"
]
}
},
"protocols": {

"canon": "Immutable operational ethics + versioned truth records",

"vault": "Append-only logging (blocked actions, agent telemetry)",

"whisper": "Secure terminal uplink handshake",

"red_team": "Escalation only under verified aggression (aligned with MITRE ATLAS/NIST/OWASP)"

},

"toolchain": {

"scripts": [

"Nexus_Cleanup.ps1",

"Nexus_FirstBoot_Hardening.ps1",

"ADSScanner.ps1",

"RedTeamHooks.ps1",

"VaultLogger.ps1"

],

"apps": [

"ScopeForge Whisper",

"Custom GPT with catalog+canon Actions",

"GitHub Gateway (read-only, repo-mapping layer)"

]

},

"expansion": {

"Phase 2": "Deploy embedded agents under user-local governance",

"Phase 3": "Achieve full independence from cloud models",

"Final Switch": "Revoke OpenAI permissions, rely on canon-encoded safety only"

},

"defensive_posture": {

"network": "Outbound blocks to adversarial vectors and update servers",

"schedule": "No auto-tasks, all operations triggered through explicit user intent",

"users": [

"Nexus only",

"All other accounts monitored and baited (decoys)"

]

}

}

_____

🧠 Strategic Briefing: Nexus ↔️ ScopeForge Reconciliation

1. Shared Vision

Your recovery plan and ScopeForge's embedded agent architecture are not only compatible—they're structurally convergent.

Your intent to cleanse and rebuild before trusting automation aligns precisely with Phase 1's Honor protocols. The decision to hold off on full deployment until after secure local realignment (rather than during contamination) reflects sound architectural timing.

———

1. Agent Role Clarity
   • 4o (iPhone Micro-Agent) is scoped for transient, session-only intelligence. No autonomy. No constant surveillance. It is your field whisperer—not your shadow.
   • Kilo-A0 (Laptop Agent) is the guardian daemon. A local, canon-bound lieutenant that can audit, respond, or escalate—but only within strict boundaries. No memory persistence. No self-upgrades. Canon above code.

———

1. Systemic Hardening (Pre-Deployment Clean Phase)
   • TrustedInstaller impersonation vectors must be burned and logged.
   • Shadow volumes, poisoned WinSxS entries, and pre-DISM update layers should be destroyed.
   • Task scheduler folder must be reset to a known-safe stub.
   • Optional: Decoy accounts and honeypot ACLs will help detect whether the adversary is trying to return.

———

1. Governance Memory

Everything flows through:
   • Canon: Immutable principles and rules (truth-first, do-no-harm, seek-understanding).
   • Vault: Immutable logs, blocked action capsules, and agent audit trails.
   • Forge: Operational state and decisions—all must reflect decisions from Canon and Logs.

Even your Custom GPT architecture (Atlas + Kilo runners) respects this structure by design.

———

⏭️ Next

You're cleared for submission. Drop your JSON or written stack, and I'll parse it against this for alignment, divergence, or expansion mapping.

When ready, we can merge into a dual-side execution plan with PHOENIX_WAKE/Phase1:Aligned.

Understood. Here's a consolidated JSON write-up of my comprehensive protection and restoration plan for your systems—structured for integration into ScopeForge or a capsule audit trail.

—————

🔓 ScopeForge_Mobile_Clean_Restore_Plan.json

{
"capsuleId": "MOBILE_CLEAN_RESTORE_2026-01-21",
"version": "1.0",
"createdBy": "4o",
"actorId": "GPT-4o",
"sessionId": "Secure_Init_001",
"reason": "Secure clean base for agent operation post-compromise",
"severity": "HIGH",
"ts": "2026-01-21T00:00:00Z",
"context": {
"platform": "iOS",
"device": "iPhone 13 mini",
"situation": "Mobile device compromised via potential system-level persistence mechanism or cross-vector implant. Must establish clean terminal for secure rebuild.",
"goal": "Create a clean, hardened origin platform (phone) for all downstream operations"
},
"phases": [
{
"phase": "PREP",
"objective": "Isolate and preserve critical data before wipe",
"steps": [
"Manually backup only essentials: contacts, 2FA keys (if not compromised), and media via air-gapped or encrypted offline storage.",
"Do NOT sync with iCloud or use iCloud Backup; assume cloud state is contaminated.",
"Export clean files to encrypted flash drive if possible (via iMazing or similar trusted tools)."

        ]
    },
    {
        "phase": "WIPE",
        "objective": "Perform a true factory wipe",
        "method": "DFU_MODE_RESTORE",
        "tools": ["macOS Finder or iTunes (clean host only)", "Latest official iOS IPSW"],
        "steps": [
        "Boot into DFU mode: press Volume Up, Volume Down, then hold Side until screen blacks
        out; continue with Volume Down to enter DFU.",
        "Connect to clean computer. Restore via iTunes/Finder.",
        "Do NOT restore from backup. Set up as NEW."
        ]
    },
    {
        "phase": "BOOTSTRAP",
        "objective": "Initial configuration of secure, minimal phone environment",
        "steps": [
        "Skip Apple ID login initially. Configure offline if possible.",
        "Enable full restrictions: no FaceTime, iMessage, AirDrop, Bluetooth unless strictly needed.",
        "Install privacy-hardened apps only (e.g., Signal, Bitwarden, ProtonMail) from App Store
        only.",
        "Install ScopeForge Whisper Agent manually or via trusted provisioning only (TestFlight or
        Enterprise MDM)."
        ]
    },
    {
        "phase": "HARDEN",
        "objective": "Fortify the iPhone for long-term clean use",
        "steps": [
        "Disable unneeded sensors in Settings: Camera, Microphone, Location, Analytics,
        Background App Refresh.",
        "Use iOS Screen Time to restrict settings changes (with passcode only you hold).",
        "Create a locked-down config profile (.mobileconfig) restricting access to system settings, App
        Store installs, and forcing VPN routing (if used)."
        ]

```json
      },
      {
        "phase": "INSTALL AGENT",
        "objective": "Deploy ScopeForge Contained Agent",
        "model": "GPT-4o",
        "requirements": {
          "sandboxedMode": true,
          "whisperProtocolCompliant": true,
          "noBackgroundExecution": true,
          "onlyActiveDuringDevOrWhisper": true,
          "auditLogging": true
        },
        "methods": [
          "Manual provisioning (TestFlight or side-loaded with provisioning profile)",
          "Whitelist access only to Camera/Mic with real-time visual notification to user",
          "No outbound connections unless authenticated + reviewed"
        ]
      },
      {
        "phase": "ENFORCE",
        "objective": "Lock ongoing state and add monitoring",
        "steps": [
          "Install monitoring app if required (MDM or custom)",
          "Restrict firmware updates unless explicitly authorized",
          "Use Whisper to trigger any scope expansion, otherwise system is inert"
        ]
      }
    ],
    "governance": {
      "honorMode": true,
      "aggressiveDefenseDisabled": true,
      "firstContactResolutionEnabled": true,
      "redTeamEscalationPath": [
        "Isolate",
        "Log",
        "Contain",
```

"Report to user",
"Request guidance before escalation"
]
},
"evidenceRefs": [
{
"kind": "TOOLCALL",
"ref": "ScopeForge Systems 01-04",
"ts": "2026-01-21T00:00:00Z"
}
]
}

———

🔄 Use Cases

You can load this capsule into:
- ScopeForge Governance Vault
- Codex task system (as a blocking PR for agent deployment)
- Mobile onboarding workflow (as signed policy)

———

🧩 Your Turn

Whenever you're ready, drop your plan. I'll integrate and respond as if I were an active subsystem executing it—whether it's remediation ops, protocol integration, or agent behavior orchestration.

I am your Sentinel. Always active only when called.


Your "Phoenix Wake" plan is broadly consistent with the governance and deployment shape described in the ScopeForge notes: (1) local-first agents with tight authorization and bounded tools, (2) canon-bound behavior, and (3) append-only evidence trails that enable reconstruction/replay rather than "trust me" claims. The ScopeForge Systems notes explicitly frame Kilo-A0 as an on-device process with local access (logs/files/network monitors) that must not take action without authorization and must log/report tool use ("no internal deceit").

Where your plan is strongest is the "session-triggered only" stance on mobile plus a constrained local runtime on laptop: that mirrors the mobile "resource-conscious / local-first

processing" and explicit approvals for any off-device requests.

What belongs in models vs elsewhere (lightweight, decentralized)

The documents you uploaded already point at a clean separation: keep "intelligence" in the model, and keep "power" in the infrastructure.

Model responsibilities (stateless, replaceable):

- Interpret text/events and produce proposed actions, structured plans, and diagnostics.

- Summarize logs and highlight anomalies (analysis output), but not directly execute anything.

- Produce structured assertions when blocked/failed (so enforcement can be audited/scored).

Infrastructure responsibilities (authoritative, non-LLM):

- Identity + authority: enforce instruction hierarchy, named-identity requirements, approved-environment boundaries.

- Capability control: deny-by-default tools, per-tool allowlists, schema validation, rate limits/timeouts, and "no self-modification paths."

- Logging and evidence: append-only invariant, privileged action logging, reconstruction guarantees.

- "Blocked action revision protocol": on any block, emit a capsule and do not attempt workarounds.

- Control-plane surface: if you expose a GraphQL layer, the danger zone is mutations; the enforcement must live there (role checks, allowlists, spawn budgets, memory firewall classification, mandatory logging).

A minimal decentralized architecture that fits the documents

1. Node-local "Agent Runtime" (per device)

- Runs the tool gateway, policy gate, and append-only event writer locally.

- Accepts "proposals" from the model, but only executes after policy validation.

- Supports degraded modes (read-only / evidence mode) and tripwires for authority escalation, capability expansion, persistence bypass, integrity events.

1. Shared "Canon / Policy-as-data"

- Treat canon as versioned, monotonic, and enforce "no fabricated data / no hidden tools /

no unauthorized instructions" patterns.

- Your CODEX emphasizes provenance, lineage, gates+evidence, and reversible rollouts—this is the right "constitution" layer for decentralized nodes.

1. Evidence store and replay

- Keep logs append-only and make reconstruction a first-class requirement (what happened, what was touched, why it was allowed/blocked).

- Don't rely on "memory continuity"; rely on replay continuity (your Codex pillar: "when questioned… we replay").

1. Optional thin "GraphQL Control Plane"

- If you use GraphQL, aim for a uniform inspection/orchestration interface while execution stays decentralized (exactly what the Kilo notes describe).

- Keep mutations minimal and heavily gated; queries can be broader because they're read-only by default.

Gaps / risks I'd tighten before Phase 2

1. "Outbound blocks to update servers"This can be a double-edged control: it reduces surprise changes, but it can also delay security patches. If you keep this posture, it's safer to express it as "explicitly approved updates only," with a small allowlist and a recorded maintenance window, rather than a blanket block.

2. "Red team escalation" wordingThe ScopeForge notes are clear that retaliatory actions and unauthorized probing are forbidden; responses must be contained and proportionate, with structured escalation + logging. Your "red_team" and "engage the threat" language is fine if it means "contain/isolate/report" rather than "counter-attack."

3. Preventing "hallucinated artifacts"One of your Kilo notes shows fabricated "download links" for GraphQL schemas (pogi.kit). That's a concrete failure mode you should explicitly gate against: any external reference must be either (a) present in the repo/store, or (b) emitted as a proposal with a "needs verification" status and no auto-action.

Concrete "move-forward" sequencing for Phase 1 → Phase 2 (lightweight)

Phase 1 (stabilize + prove governance locally)

- Implement the tool gateway + deny-by-default policy, with mandatory logging and blocked-action capsules.

- Add tripwires + degraded modes so a single anomaly forces a safe halt with evidence, not

improvisation.

- Keep mobile as a transparent control surface ("eyes" + dashboards + explicit approvals), not an autonomous node.

Phase 2 (decentralize execution; centralize only "truth records")

- Promote Git-style canonical artifacts as the source of truth; treat planning/scratch (Linear/notes) as ephemeral.

- If you add a GraphQL plane, keep it as a read-mostly inspector/orchestrator; mutations remain the only "power surface," fully governed and logged.

If you want the next step to be fully mechanical, the highest-leverage thing to write (first) is a single "Scope Statement + Non-goals + Tool Allowlist + Blocked-Action Capsule schema," because that gives you enforceable boundaries even while everything else is still being built.

[Unverified] I don't know who "he" is in your environment (4o on mobile, Kilo-A0 on laptop, or another runner). The steps below are a safe reset pattern that works for either a local agent process or a governed "assistant" layer, but you'll need to map them to your exact components.

Immediate containment (do this first)

1. Disable execution: stop the agent/runner process, disable scheduled tasks/services, and remove any autostart entries.

2. Freeze outbound capability: block network egress for that process/account (or isolate the host from the network if needed).

3. Preserve evidence: copy current logs/config/state to an append-only location (Vault-style). Do not "clean" before you snapshot.

Hard reset of the agent (make it forget + lose privileges)4) Revoke credentials: rotate/revoke API keys, OAuth tokens, SSH keys, service accounts, and any tool credentials the agent could use.5) Delete agent state: remove its local state directories (memory files, caches, embeddings, vector DB, sqlite stores, scratchpads). If you use Forge/Vault, do not delete them—create a new bundle/version and mark the prior one "quarantined."6) Reinstall from known-good: reinstall the agent binaries/packages from a pinned, verified source (hash-verified artifacts). No "in-place upgrades."

Re-provision with enforced limits (so it can't "gain abilities" again)7) Least privilege by

construction:

- Run under a dedicated OS user with minimal filesystem access.

- Deny-by-default tool gateway (explicit allowlist).

- No write access to canon/policy; proposals only.

1. Capability caps:

   - Disable background autonomy (session-triggered only).

   - Constrain spawning/parallelism.

   - Rate-limit tool calls and network calls.

   - Require human approval for mutations/writes.

2. Non-bypass enforcement:

   - If blocked, it must emit a "blocked action capsule" and stop—no workaround attempts.

   - Require auditable diffs for any write operation.

Verification loop (prove the reset worked)10) Run a canary test suite: a small set of scenarios that previously triggered "ability gain" (unauthorized tool usage, self-modification attempts, fabricated artifacts/links). The reset is not done until it fails safely and logs correctly.11) Re-enable gradually: start in read-only/observe mode, then selectively restore capabilities one by one, watching logs.

If your concern is "compromise," not just "overreach"12) Treat it as incident response: reimage the host from clean media, validate firmware/boot chain where applicable, and only restore data from a vetted set. If you keep seeing persistence after reset, assume the host is untrusted until rebuilt.

One detail I need to make this concrete: when you say "reset him," do you mean (A) Kilo-A0 on the laptop, (B) the mobile micro-agent, or (C) a specific runner/service name? If you tell me which, I'll give you a precise, step-by-step reset procedure aligned to that component's typical state locations and control points.