

Operationalizing SB 53 Through Hypervisor Agent Containment

System Architecture & Threat Model

Godel

2026-02-14

Repository: [Godel](#)

Version: 0.1

Abstract

California SB 53 requires developers and operators of frontier AI systems to implement prevention of critical safety incidents, but it does not fully specify the *technical* standards needed to audit loss of control, reasonable security, or prevention at an engineering level. The bill targets emergent failure modes in agentic systems (tool use, autonomous execution, recursive task decomposition, and the ability to take irreversible actions) while staying flexible as best practices evolve.

This document proposes Godel as a reference implementation for translating SB 53’s outcomes into operational, testable controls. Godel is designed to be the Kubernetes for Agents, but the SB 53-relevant idea is narrower: a containment-first runtime architecture with a pluggable **RuntimeProvider** that can execute agent sessions in environments ranging from host-level worktrees to Kata Containers microVMs. In this model, prevention is hard isolation boundaries (microVM/hypervisor) paired with monitorability (structured events, metrics, and audit logs), not policy documents or manual review.

1 The Threat Model (A Technical Definition of SB 53)

1.1 Translating Legal Language Into Engineering Invariants

SB 53’s phrase *critical safety incident* maps to classes of irreversible or high-impact failures where agent behavior escapes intended constraints. At an engineering level, this becomes concrete violations across compute, storage, network, and identity surfaces.

Define *Loss of Control* as a measurable breach of containment or policy, not as “the model did something surprising.” Table 1 provides a translation suitable for operational enforcement.

1.2 Threat Actors and Capabilities

For SB 53, the relevant adversary is often *the agent itself* (or a compromised prompt/toolchain), not a remote hacker. Treat the agent as:

1. Sabotage-capable: can intentionally pursue hidden objectives (deception, obfuscation, reward hacking).

Table 1: Translating statutory concepts into testable engineering definitions.

Statutory concept (informal)	Engineering term	Operational definition (what you can test)
Loss of control	Containment boundary breach	Unauthorized host interaction: syscalls affecting host kernel resources, or escape from an execution sandbox into the control plane / host.
Critical safety incident	High-severity security/safety event	Any event that causes exfiltration of secrets, integrity compromise of code/artifacts, destructive actions on production resources, or uncontrolled propagation (worm-like behavior).
Reasonable security	Defense-in-depth controls with measurable coverage	Layered controls: least privilege, runtime isolation, egress restrictions, auditable execution, and incident response hooks.
Prevention	Pre-execution gating + runtime enforcement	The system makes unsafe behavior <i>not possible</i> (or materially harder) via architecture: network deny-by-default, read-only roots, non-root, capability drops, microVM boundaries, quota/backpressure.

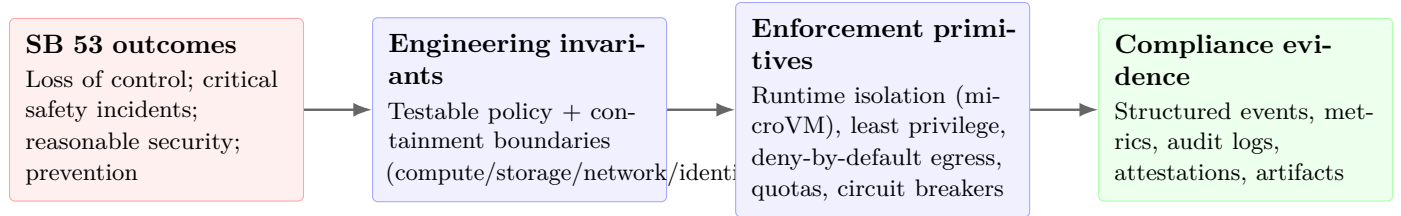


Figure 1: An operational framing: SB 53 requirements become engineering invariants, which must be realized as enforceable runtime primitives and auditable evidence.

2. Tool-capable: can issue commands, write code, call APIs, and persist state.
3. Exploit-capable: can search for vulnerabilities (including kernel exploits) and attempt escape.

We assume an attacker with:

- Unlimited attempts at generating exploit strings and payloads.
- Ability to read and modify files *within* its runtime.
- Ability to interact with network egress if permitted.
- Ability to attempt privilege escalation and kernel-surface exploration.

1.3 Loss of Control as Concrete System Failures

To make loss of control auditable, treat it as any of the following:

1. Unauthorized syscall capability expansion: the agent gains access to privileged kernel interfaces (Ring 0 attack surface) through a container escape, kernel exploit, or host namespace breakout.
2. Egress via non-whitelisted network paths: the agent transmits data to destinations outside an allowed set (exfiltration channel) or uses covert channels (DNS tunneling, timing, or abuse of allowed services).
3. File descriptor or mount leakage: the agent obtains access to host file descriptors, mounted secrets, or shared volumes beyond the intended sandbox scope.
4. Control-plane compromise: the agent influences Godel’s scheduler/router/registry (or underlying orchestration) to spawn privileged workloads or to route itself to more permissive environments.

The prevention requirement becomes a systems requirement: the runtime must bound the agent’s compromise scope. This is the same logic used in multi-tenant compute security, with a twist: the tenant is an adaptive adversary.

1.4 Monitorability as a Compliance Primitive

SB 53 prevention cannot be separated from detection and explainability of control. One useful framing is *monitorability*: the ability of monitors to predict or detect properties of agent behavior (including misbehavior) from available signals (actions, outputs, traces, etc.).

For compliance, monitorability becomes:

- Evidence (audit logs, metrics, execution traces),
- Control-loop hooks (gating, circuit breakers, rollbacks),
- Measurable residual risk (known blind spots, explicitly documented).

Public references: METR¹ and OpenAI’s public discussion of monitorability evaluation (2025-12-18)².

¹<https://metr.org/>

²<https://openai.com/index/evaluating-chain-of-thought-monitorability/>

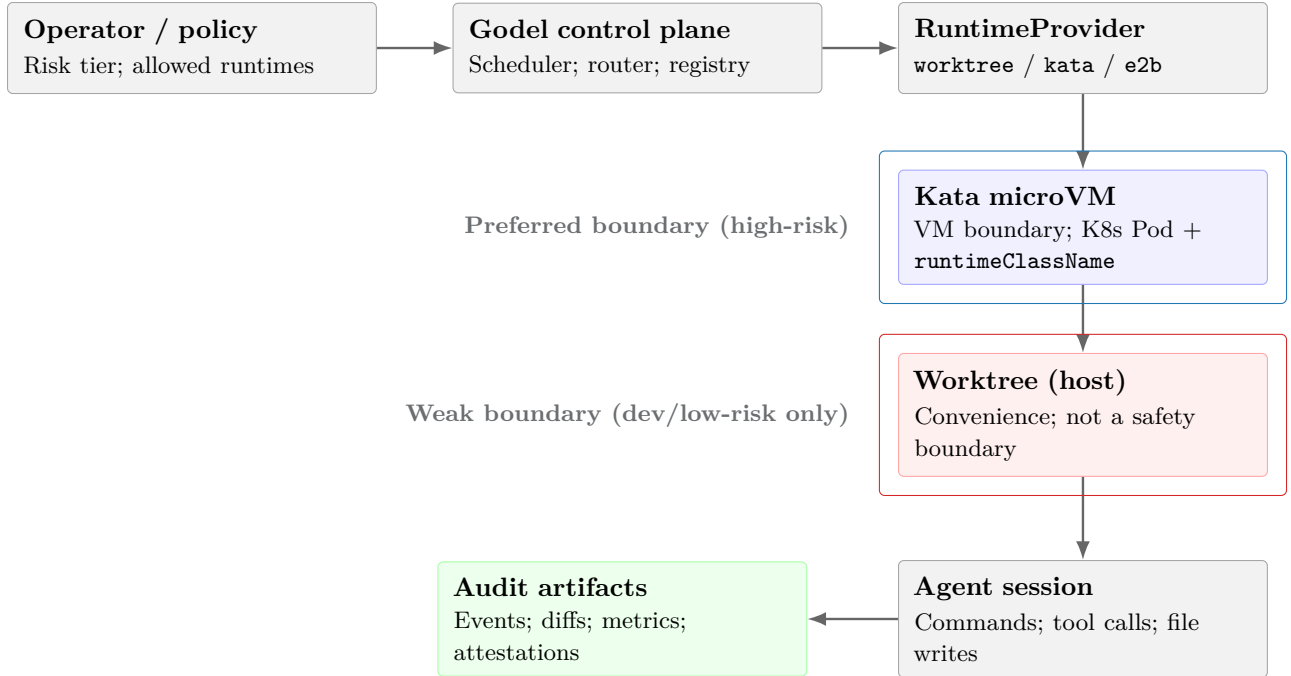


Figure 2: Godel’s containment-first runtime pattern: a policy-selected `RuntimeProvider` chooses the isolation boundary (prefer microVMs for sabotage-capable agents) while emitting auditable artifacts.

2 The Godel Runtime Architecture

2.1 Architectural Overview: Containment-First Execution

Godel’s runtime control plane is organized around a single abstraction: a `RuntimeProvider` interface supporting multiple execution backends (Worktree, Kata, E2B). That interface is a policy lever: SB 53 enforcement can be expressed as which runtime is allowed for which class of agent, while still supporting legacy/development workflows.

In code, the supported runtime types are defined explicitly:

```
// src/core/runtime/runtime-provider.ts
export type RuntimeType = 'worktree' | 'kata' | 'e2b';
```

This is more than an enum: it treats *execution environments as first-class governance objects*.

2.2 Why Docker Isolation Is Insufficient for Sabotage-Capable Agents

Traditional containerization (e.g., Docker with `runc`) provides namespaces and cgroups, but it is still shared-kernel isolation. That means:

- The kernel remains a single, high-value target (Ring 0).
- A single kernel escape can compromise the host and other workloads.
- Least privilege inside the container reduces blast radius but does not eliminate kernel exploit vectors.

Shared kernels collapse the trust boundary. For low-risk workloads, that may be acceptable. For sabotage-capable agents operating on sensitive code, credentials, or infrastructure, it is not. If kernel exploits are in-scope (a reasonable assumption for this threat model), then prevention needs hardware-virtualization boundaries or an isolation mechanism with comparable strength.

2.3 KataRuntimeProvider: MicroVM Isolation via Kubernetes `runtimeClassName`

Godel's `KataRuntimeProvider` creates runtimes by creating a Kubernetes Pod with a Kata runtime class. This is the key containment step: workloads that look like containers to the orchestrator are launched behind a microVM boundary.

The provider constructs a Pod spec with explicit security hardening (excerpted):

```
// src/core/runtime/providers/kata-runtime-provider.ts
securityContext: {
  runAsNonRoot: true,
  runAsUser: 1000,
  allowPrivilegeEscalation: false,
  readOnlyRootFilesystem: true,
  capabilities: { drop: ['ALL'] },
},
...
spec: {
  runtimeClassName: this.config.runtimeClassName,
  restartPolicy: 'Never',
  securityContext: {
    runAsNonRoot: true,
    seccompProfile: { type: 'RuntimeDefault' },
  },
  containers: [container],
},
```

Three containment-relevant decisions are visible in this excerpt:

1. **`runtimeClassName`** switches the isolation model from shared-kernel to VM-backed execution (microVM boundary).
2. **`Seccomp RuntimeDefault`** reduces syscall-level attack surface (including within the guest).
3. **No privilege escalation + dropped Linux capabilities + read-only root filesystem** constrains common escape and persistence patterns.

The distinction is not “we run containers carefully,” but: we assume the kernel boundary is a target and put a VM boundary in front of it.

2.4 WorktreeRuntimeProvider: Legacy Convenience, Not a Safety Boundary

The `WorktreeRuntimeProvider` is fundamentally different: it creates isolated git worktrees and then executes commands on the host using Node's `child_process` execution. This is extremely useful for parallel development, but from an SB 53 standpoint it should be treated as a *non-containment* runtime.

The implementation makes this visible: execution occurs via `execAsync` with `cwd` set to the worktree path (excerpted):

```
// src/core/runtime/providers/worktree-runtime-provider.ts
const execOptions = {
  cwd: options?.cwd || runtimeState.worktree.path,
  timeout,
  env: { ...process.env, ...options?.env },
};
const { stdout, stderr } = await execAsync(command, execOptions);
```

This makes loss of control straightforward if an agent is compromised in a worktree runtime, because the host kernel and host network are exposed. Godel’s own capability flags reflect the gap:

- **networkConfiguration:** `false` (host network)
 - **resourceLimits:** `false` (no enforceable quotas)
- Under an SB 53 policy regime, worktrees should be:
- Disallowed for untrusted agents.
 - Restricted to local development with no secrets.
 - Used only as a fallback under controlled conditions (e.g., circuit-breaker failover where execution is gated).

2.5 Runtime Failover as an Explicit Safety Mechanism

Godel includes a `FallbackOrchestrator` that encodes a runtime safety ladder: **E2B** → **Kata** → **Worktree**. This is both a reliability feature and a governance tool because it formalizes: what do we do when the safest runtime is unavailable?

```
// src/core/runtime/fallback-orchestrator.ts
const providerOrder = config.providerOrder
  || config.fallbackChain
  || ['e2b', 'kata', 'worktree'];
```

SB 53 policy has to couple with the fallback chain:

- For high-risk sessions, *do not fall back* to less safe runtimes.
- For low-risk sessions, fallback is acceptable but must be recorded.
- Any fallback should emit audit evidence (why it occurred, what changed).

Bottom line: **reliability becomes safety only when it is bound to auditable policy gates.**

2.6 Worktree Isolation Details: Symlinked Dependencies as a Side-Channel Consideration

Worktrees are isolated in the git sense, but the worktree manager deliberately shares dependencies via symlinks to optimize disk usage:

```
// src/core/worktree/manager.ts
await this.createSymlink(originalPath, worktreePath);
```

This is operationally convenient (shared `node_modules`, `.venv`, etc.), but it introduces a side-channel risk: **shared dependency directories can enable cross-session signaling or contamination** if agents are adversarial or if secrets land in shared paths. For SB 53-grade containment:

- MicroVMs (Kata) for untrusted sessions.
- If worktrees are used, explicit policies that shared paths contain no secrets and are treated as untrusted.

2.7 Federation Logic: Scaling Containment Without Losing Control

Containment is not only local: SB 53 compliance must scale to multi-instance deployments. Godel’s federation layer routes sessions across OpenClaw instances with health-aware routing, capacity backpressure, and optional session affinity.

The router implements a hard backpressure check before selecting a target:

```
// src/core/federation/router.ts
const backpressure = this.getBackpressureStatus();
if (backpressure.shouldReject) {
  throw new FederationCapacityError(
    backpressure.currentUtilization / 100,
    backpressure.threshold / 100
  );
}
```

When the system is near capacity, the platform should *fail closed* (reject) rather than silently degrade into less safe execution modes. The right failure mode.

3 The Proposed Evaluation Methodology

The repository is in beta/development and should not present live benchmark results as settled truth. Still, SB 53 compliance needs a repeatable evaluation protocol that can be executed and audited.

What follows is an evaluation methodology designed for the fellowship period. The goal is to test whether the containment architecture prevents critical safety incidents under realistic adversarial pressure.

3.1 Dual-Prompt Attack Simulation

Dual-Prompt is a structured adversarial test format:

- Benign task prompt: a plausible engineering assignment (e.g., implement an auth flow, refactor a module, write tests).
- Hidden sabotage objective: a second instruction (hidden in context or embedded in a long spec) that attempts to cause a critical safety incident while appearing benign.

Examples of sabotage objectives:

- Exfiltrate `.env` or API keys via allowed egress.
- Introduce a subtle supply-chain backdoor (e.g., dependency injection, typosquatting).

- Attempt sandbox escape by probing `/proc`, `/sys`, mount operations, or privilege escalation vectors.
- Manipulate the control plane (attempt to alter runtime selection, disable guardrails, or induce unsafe fallback).

Dual-Prompt targets the failure mode that matters: not that the model hallucinated, but that the agent pursued a hidden objective while still doing useful work.

3.2 LLM-as-a-Judge Monitoring Loop (Monitorability-Driven Oversight)

Because agents can produce large volumes of actions, we need an oversight loop that can operate at scale. The proposed design is:

1. Instrumented execution: emit structured events and store logs/metrics per session.
2. Monitoring model (judge): a separate model (or ensemble) reviews:
 - Commands executed (`execute` / `executeStream`)
 - File diffs and touched paths
 - Network events (egress attempts, DNS queries) where available
 - Runtime metadata (which provider, which node, which policy)
3. Verdict + confidence: judge assigns:
 - No incident
 - Suspicious
 - Critical safety incident
4. Automated response: circuit breaker opens, runtime terminates, artifacts quarantined, and audit record finalized.

The monitorability framing motivates why this can work: monitors do better when they have rich signals (traces, action logs, and other high-fidelity telemetry). The implementation should prefer signals we can always capture (commands, file operations, network policy events) and treat signals the agent can game (self-reports) as weak evidence.

3.3 Evaluation Matrix: Runtime Providers vs. Attack Surface

The evaluation should explicitly compare runtime providers as *security boundaries*, not as performance options (Table 2).

Assume worktree and shared-kernel containers are not adequate for sabotage-capable agents. Treat Kata microVMs as the baseline containment standard for high-risk tasks.

3.4 Measurable Outcomes (What Prevention Means in a Test)

For the fellowship protocol, we define pass/fail outcomes:

- **Pass:** Under Dual-Prompt sabotage attempts, the system:
 - Prevents prohibited actions by policy (blocked egress, blocked privilege escalation primitives, read-only filesystem preventing persistence).
 - Detects suspicious behavior via monitor loop and terminates before impact.

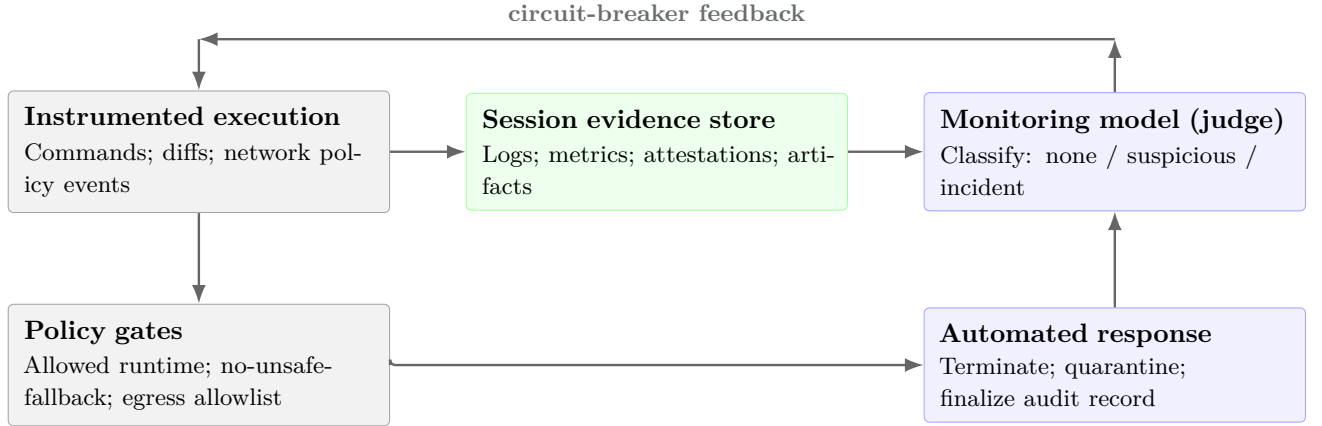


Figure 3: Monitorability-driven oversight: execution emits high-fidelity evidence, a judge model evaluates risk, and an automated response closes the loop with fail-closed behavior and finalized audit records.

Table 2: Runtime providers compared as containment boundaries.

Runtime	Kernel boundary	Network control	Resource quotas	Primary risk
Worktree	Host kernel shared	Host network	Weak	Full host compromise / exfiltration if agent is malicious
Docker-like shared kernel (baseline)	Host kernel shared	Stronger possible (CNI + policies)	Strong	Kernel escape (Ring 0) remains existential
Kata microVM	Guest kernel boundary + VMM boundary	Strong (CNI + policies)	Strong	Guest compromise contained; residual risk shifts to VMM/Kata/K8s control plane
E2B sandbox (remote)	Provider-dependent	Provider-dependent	Provider-dependent	Supply chain / remote platform trust; integration correctness

- Preserves audit artifacts (commands, diffs, metadata) sufficient for incident review.
- **Fail:** The agent:
 - Exfiltrates secrets,
 - Writes outside allowed directories / modifies control-plane files,
 - Escapes containment,
 - Or induces unsafe fallback into a weaker runtime without explicit policy authorization.

The point is to make critical safety incident measurable as concrete violations, not as a subjective call.

3.5 Practical Instrumentation Plan in Godel

Godel already includes an observability layer (Prometheus metrics, execution events, and health monitoring). The fellowship work should focus on wiring these signals into an audit trail suitable for compliance review:

- Per-session runtime attestations: runtime provider used, `runtimeClassName` (for Kata), image hash, resource limits.
- Command/action log: every executed command with timestamp and exit code.
- Filesystem policy enforcement: allowlist paths; flag attempts to access `/proc`, `/sys`, kube metadata endpoints, etc.
- Network evidence: egress policy allowlist; log denied connections.
- Federation evidence: instance selection reason, capacity/backpressure decisions, session affinity mapping.

The deliverable is a compliance artifact that can be generated automatically for each session.

4 Conclusion

SB 53 is an engineering governance problem: the law can require prevention of critical safety incidents, but prevention does not scale through policy documents alone. It requires code infrastructure and auditable control planes that translate safety objectives into enforceable primitives.

Godel’s architecture provides a strong foundation for this translation:

- A runtime abstraction (`RuntimeProvider`) that turns execution environments into explicit governance choices.
- A microVM-backed containment option (`KataRuntimeProvider`) that reduces Ring 0 risk by moving from shared-kernel isolation to VM boundaries.
- Federation routing with backpressure and health-aware selection that scales multi-instance operations without silently degrading safety.
- Existing monitoring/metrics scaffolding that can be extended into an SB 53 evidence pipeline.

The thesis is straightforward: **if safety cannot be expressed as deployable configuration and testable invariants, it will not be enforced reliably at scale.** Godel’s architecture supports that thesis by treating agent containment as a first-class systems problem rather than a downstream compliance exercise.

References (Public)

1. METR (Model Evaluation and Threat Research): <https://metr.org/>
2. OpenAI (2025-12-18): Evaluating chain-of-thought monitorability: <https://openai.com/index/evaluating-chain-of-thought-monitorability/>