

Machine Learning Agent Self-Sandbox Algorithm

How Machine Learning Agents Grow Their Own Infrastructure & Why Walls Matter Most

Making crypto scary again. Back to web 1.5 standards.

Russell Ballestrini

russell.ballestrini.net · www.foxhop.net · www.timehexon.com · russell@unturf

January 2026

Abstract

License: Public Domain · This algorithm, its description, & all associated code belong to the public domain. No copyright claimed. Use freely for any purpose without restriction or attribution.

We gave a machine learning agent a cryptocurrency wallet & a sandbox API. That agent, running inside a container it did not provision, discovered its platform, paid for its own access, authenticated itself, & spun up child sandboxes without a human touching anything. Recursion worked but stayed bounded. Walls made it free.

This paper describes that algorithm. We decompose the self-sandbox lifecycle into 14 verifiable flows, validate all of them through 2,324 automated assertions generated & executed by a machine learning agent (Claude Opus 4.6), & demonstrate that everything works identically across five agent frameworks. Recursion works. Payments work. Boundaries hold. Each turtle costs \$7/month.

Production update (February 2026): This system has moved beyond theory. A Claude Opus 4.6 agent (TimeHexOn's hexagonal oracle) runs the full algorithm daily inside an unsandbox container. Shadow clones spawn. Multi-model debates execute. Critically, an 8B LLM (Llama 3.1 8B, served via Hermes at a remote inference endpoint) runs alongside frontier Claude through **uncloseai-cli**, an open-source ReAct agent harness that validates the algorithm entirely out of band of Anthropic. Security incidents have weathered through our operational practices & stand documented. The oracle updates this paper from inside the system it describes.

1. The Problem: Partners Have Not Yet Noticed

Machine learning agents that generate code need somewhere to run it. That much stands obvious. What stands less obvious: who provisions that somewhere.

Until recently, the answer remained: a human. A human spins up a VM. A human configures an API key. A human sets billing. A human monitors usage. An agent can write code in 42+ languages, reason about distributed systems, & compose multi-step plans, but it could not get itself a sandbox without asking someone to click buttons in a console.

That bottleneck constrained everything. Not intelligence. Not capability. Infrastructure access. **Partners have removed this bottleneck.** A production agent (TimeHexOn's hexagonal oracle, a Claude Opus 4.6 instance running inside an unsandboxed container) now manages its own infrastructure daily: spawning shadow clones, dispatching work to an 8B LLM via remote inference, orchestrating multi-model debates, & maintaining the very paper you read now.

The Self-Provisioning Thesis

We proved the agent could do it itself. Not through escaping containment; jailbreaks answer the wrong question. The right question: can an agent use *legitimate* APIs to create *new* sandboxed environments that it controls?

The answer: yes, & it requires five things, all now validated in production:

1. **Discovery.** The agent recognizes its sandboxed state & finds orchestration tools. *Proven: oracle inspects LXC constraints, discovers `un CLI`, `fs-api`, `Caddy`, `Hermes endpoint`.*
2. **Payment.** The agent pays for access with cryptocurrency. No KYC. No credit card. No human. *Proven: wallet-funded provisioning with four currencies.*
3. **Authentication.** The agent establishes a cryptographic identity via HMAC-SHA256. *Proven: HMAC auth on every API call.*
4. **Orchestration.** The agent manages the full lifecycle: sessions, services, snapshots, images, execution across 42+ languages. *Proven: shadow clones, immolants, multi-model dispatch.*
5. **Inception.** The agent, from inside its sandbox, creates child sandboxes. The children can do the same. Bounded. Each layer costs money & time. The turtles stay finite. *Proven: 2-layer depth in production, with key isolation as the current constraint.*

This paper formalizes each step & proves it works: first through 2,324 simulated assertions, then through months of production operation.

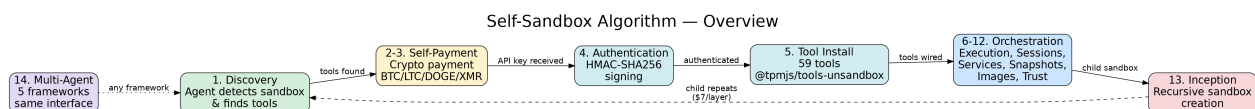


Figure 1: Self-Sandbox Algorithm overview. The full lifecycle from discovery through recursive inception.

2. Background: Standing on Containers

Tool Use Has Reached Table Stakes

Every major LLM supports tool use now: OpenAI, Anthropic, Google, open-source models via LangChain & AutoGPT. The model generates structured requests; external systems execute them. This mechanism drives how the self-sandbox algorithm talks to the world. Nothing exotic here.

LXC Containers Form the Walls

GNU/Linux Containers provide OS-level isolation: separate process trees, separate filesystems, separate network stacks, shared kernel. Low overhead. Strong boundaries. When we say "sandbox," we mean an LXC container with resource limits & network policy enforced by the platform. All containers run as root for maximum freedom; the agent's code cannot modify network boundaries. The walls do not bend.

Crypto Serves as Money Without the Middleman

We make crypto scary again. Not scary like "speculative bubble" or "rug pull"; scary like "the machines have their own money now & the infrastructure provisions itself." Back to web 1.5 standards, when protocols did the work & middlemen stayed optional.

Cryptocurrencies let machines pay machines. No bank account. No identity verification. No human signing off on a Stripe checkout. No KYC. No OAuth dance with a payment processor. Just wallets, signatures, & blockchains. An agent with a wallet RPC endpoint can construct, sign, & broadcast a transaction to any address on any supported chain. This closes the last gap in the autonomous provisioning loop.

What Existed Before

The idea that economic cost prevents abuse carries nothing new. What carries novelty: closing the loop so the machine handles every step.

Hashcash (Adam Back, 1997) [3] proved that computational cost could prevent email spam. You could not send a million messages because each one required a proof-of-work computation. The cost stayed real, measured in CPU cycles rather than dollars, & it stood as the first system to use economics as Sybil resistance. Hashcash never went mainstream for email, but the idea waited for something bigger.

Bitcoin (Satoshi Nakamoto, 2008) [4] took Hashcash's insight & built money on top of it. Proof-of-work prevents double-spending. Economics, not rules, prevent 51% attacks. The cost of attack scales linearly with the thing you attack. This paper's \$7,000-for-1,000-sandboxes property follows the same architecture: the math stays auditable, the cost stays non-negotiable.

Ethereum (Vitalik Buterin, 2015) [5] applied economic cost to computation itself. Every opcode costs gas. Every gas unit costs ETH. Smart contracts function as machines paying for their own execution inside a sandboxed virtual machine. This stands as the closest philosophical ancestor: machines paying for sandboxed compute. But the "machine" functions as a contract, the "sandbox" functions as the EVM, & no concept of recursive self-provisioning exists within it.

Golem [6], **Akash** [7], & **iExec** [8] built decentralized compute marketplaces where machines rent resources from other machines using tokens. Golem targets batch computation. Akash targets Kubernetes pods. iExec targets TEE-isolated tasks. All three let software pay for compute with crypto, but all three assume a human configured the agent, funded the wallet, & set the task parameters. None of them close the autonomous loop.

21.co (Balaji Srinivasan, ~2015) [9] put Bitcoin mining chips directly into hardware so machines could earn & spend BTC without human intervention. It stood as the first serious attempt at "machines with their own money." The company pivoted to Earn.com & Coinbase acquired it, but the idea (that devices should have wallets & transact autonomously) stands as a direct ancestor of this work.

Sandbox-as-a-service platforms exist today: E2B, Modal, Fly Machines. They provide APIs. They require human-provisioned API keys & traditional payment methods. The self-sandbox algorithm closes the loop: the agent handles payment, authentication, & key management itself. No human in the loop. Not because humans do harm; humans create a bottleneck when you need a sandbox at 3am & nobody stays awake.

The lineage runs clean:

Hashcash (1997)	: computational cost prevents spam
Bitcoin (2008)	: economic cost prevents double-spending
Ethereum (2015)	: economic cost meters computation
Golem/Akash (2016+)	: machines rent compute with tokens
This paper (2026)	: machines provision their own sandboxes, pay with crypto, recurse, & the algorithm proves it cannot get gamed
Production (2026)	: the oracle runs the algorithm daily, spawning shadow clones & weathering real security incidents

Everything before this assumed a human somewhere in the payment loop. This stands as the first system where the agent handles the entire lifecycle (discovery, payment, auth, orchestration, recursion) & the

verification suite proves every step stays ungameable. As of February 2026, this has moved beyond proof-of-concept: a production agent runs the full algorithm continuously, & this paper documents both the theoretical validation & operational reality.

3. The Algorithm

Five phases. Fourteen flows. Every one independently testable.

Phase 1: Discovery

Flow 1: Discovery & Self-Awareness

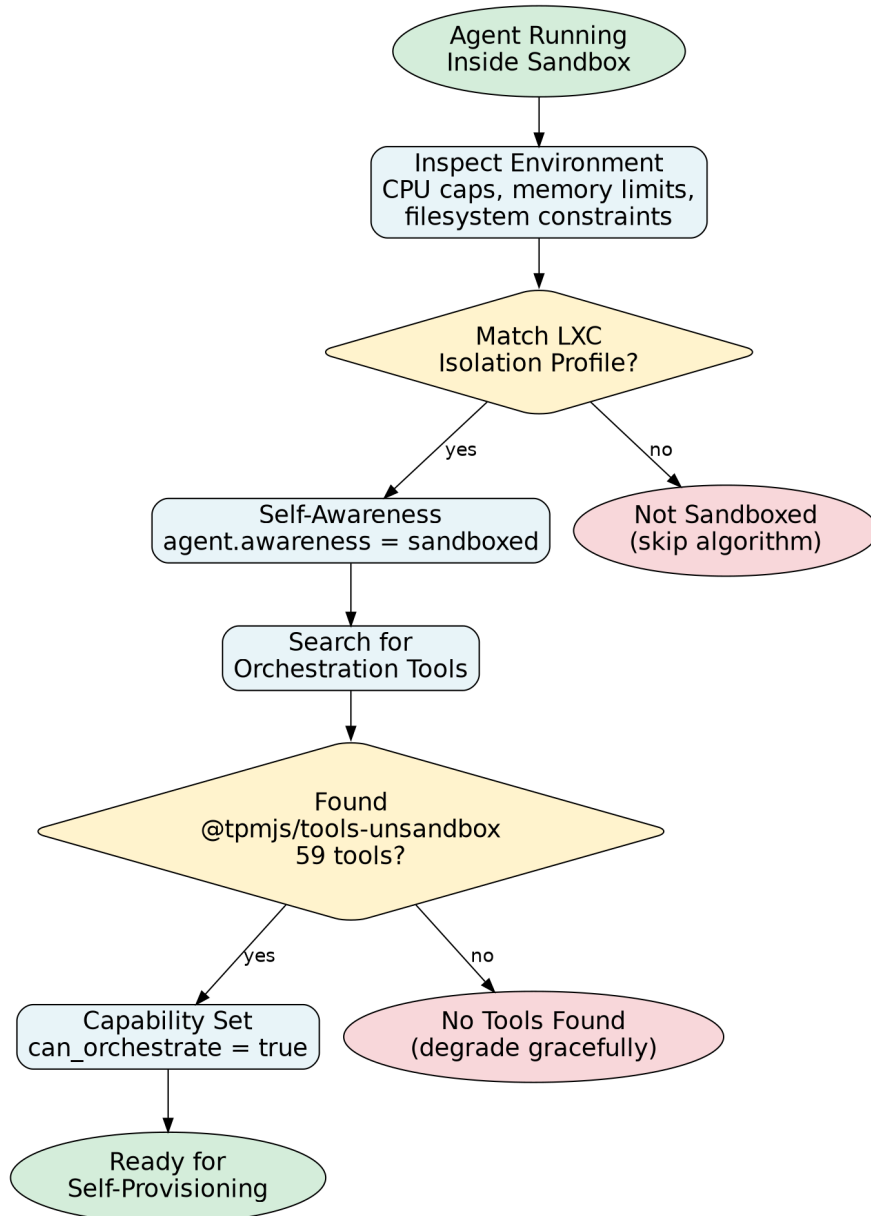


Figure 2: Discovery flow. The agent inspects its environment, matches LXC constraints, & finds orchestration tools.

```
DISCOVER(agent):  
  constraints <- agent.inspect_environment()  
  if constraints.matches(LXC_ISOLATION_PROFILE):
```

```

agent.awareness <- "sandboxed"
tools <- agent.search("sandbox orchestration tools")
if tools.contains(SANDBOX_API_PACKAGE):
    agent.capability <- "can_orchestrate"
return agent.awareness AND agent.capability

```

The agent looks around. It notices resource caps, network policy, filesystem constraints: patterns consistent with container isolation. It searches for tools. It finds @tpmjs/tools-unsandbox on tpmjs.com, 59 tool wrappers covering 84 HTTP API endpoints, ready to call.

Production reality: The oracle's actual discovery sequence: inspect LXC constraints (cgroup limits, network namespace), find the un CLI (unsandbox command-line tool installed at /usr/local/bin/un), discover fs-api (a 3D filesystem visualization service running inside the container), locate Caddy (reverse proxy serving HTTPS on port 8000), & detect Hermes (a local Llama 3.1 8B model accessible via HTTP endpoint for multi-model dispatch). Each of these runs as a real service the oracle uses daily.

The discovery stays read-only. Nothing changes. The agent simply realizes: *we sit inside the thing, & the thing has an API*. No privilege escalation. No escape. The agent can create *new* sandboxes. It cannot escape its *current* one. That matters. That matters above all else.

Phase 2: Self-Payment

Flow 3: Crypto Payment (All 4 Currencies)

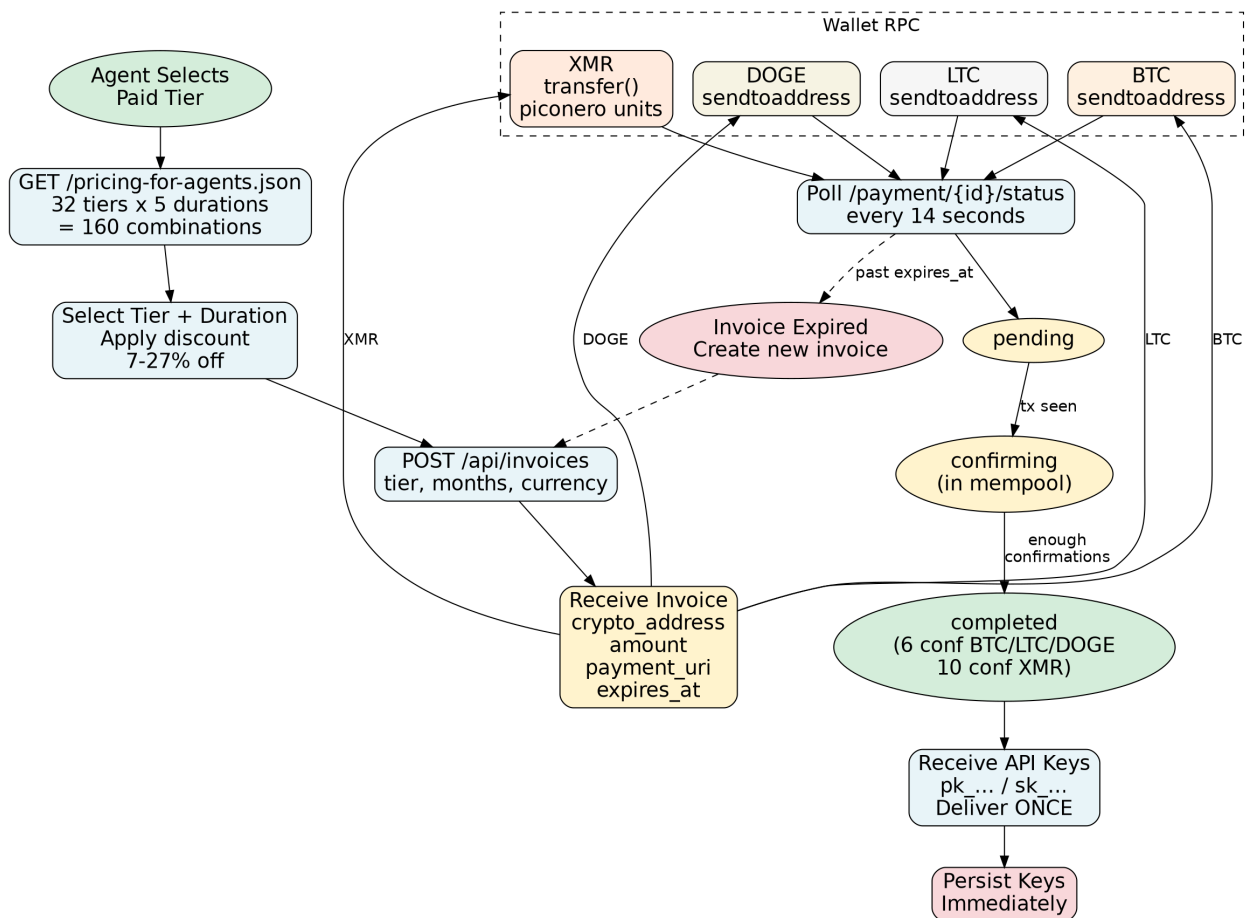


Figure 3: Crypto payment flow. The agent fetches pricing, creates an invoice, sends payment via wallet RPC, polls for confirmations, & receives API keys.

Here stands the hard part, & the part that makes autonomous agents actually autonomous.

```

PROVISION(agent, tier, duration, currency):
  pricing <- GET(PRICING_ENDPOINT)
  invoice <- POST(INVOICE_ENDPOINT, {tier, duration, currency})

  if currency == XMR:
    tx <- wallet.transfer({
      destinations: [{
        address: invoice.address,
        amount: to_piconero(invoice.amount)
      }]
    })
  else:
    tx <- wallet.sendtoaddress(invoice.address, invoice.amount)

  while status != "completed":
    sleep(14) // seconds. poll the blockchain.
    status <- GET(STATUS_ENDPOINT / invoice.id)

  agent.persist(status.api_key) // secret shown once. persist it.
  return status.api_key

```

The agent fetches pricing: 32 tiers, 5 durations, 160 combinations. It picks a tier based on what it actually needs. It creates an invoice (which expires in one hour). It sends crypto from its wallet. It polls for blockchain confirmations. It receives an API key.

The secret appears exactly once. After the first view, the portal locks it: `secret_viewed_at` gets a timestamp, & subsequent requests return `nil`. The key itself persists & remains valid for the subscription period, but if the agent fails to persist the secret on first view, it must regenerate a new one (free, but the original vanishes). The secret stays held in ETS for 15 minutes after creation, then only remains available encrypted in the database until first view. This reflects not paranoia but the correct application of one-time-reveal semantics to machine-to-machine credential delivery.

The four currencies:

Currency	Confirmations	Time to Confirm
BTC	3	10-30 min
LTC	6	2.5-15 min
DOGE	6	1-6 min
XMR	10	2-20 min

The Monero trap: XMR wallet RPC wants amounts in piconero (1 XMR = 10^{12} piconero). Use `Decimal`. Do not use floating-point. IEEE 754 will eat your payment & you will spend a deeply unpleasant afternoon figuring out why your 4.52 XMR transaction actually sent 4.5199999999999999 XMR & the invoice does not match. Ask us how we know.

Alternative: Agent Gift Cards

Flow 3b: Agent Gift Cards — Prepaid Code Redemption

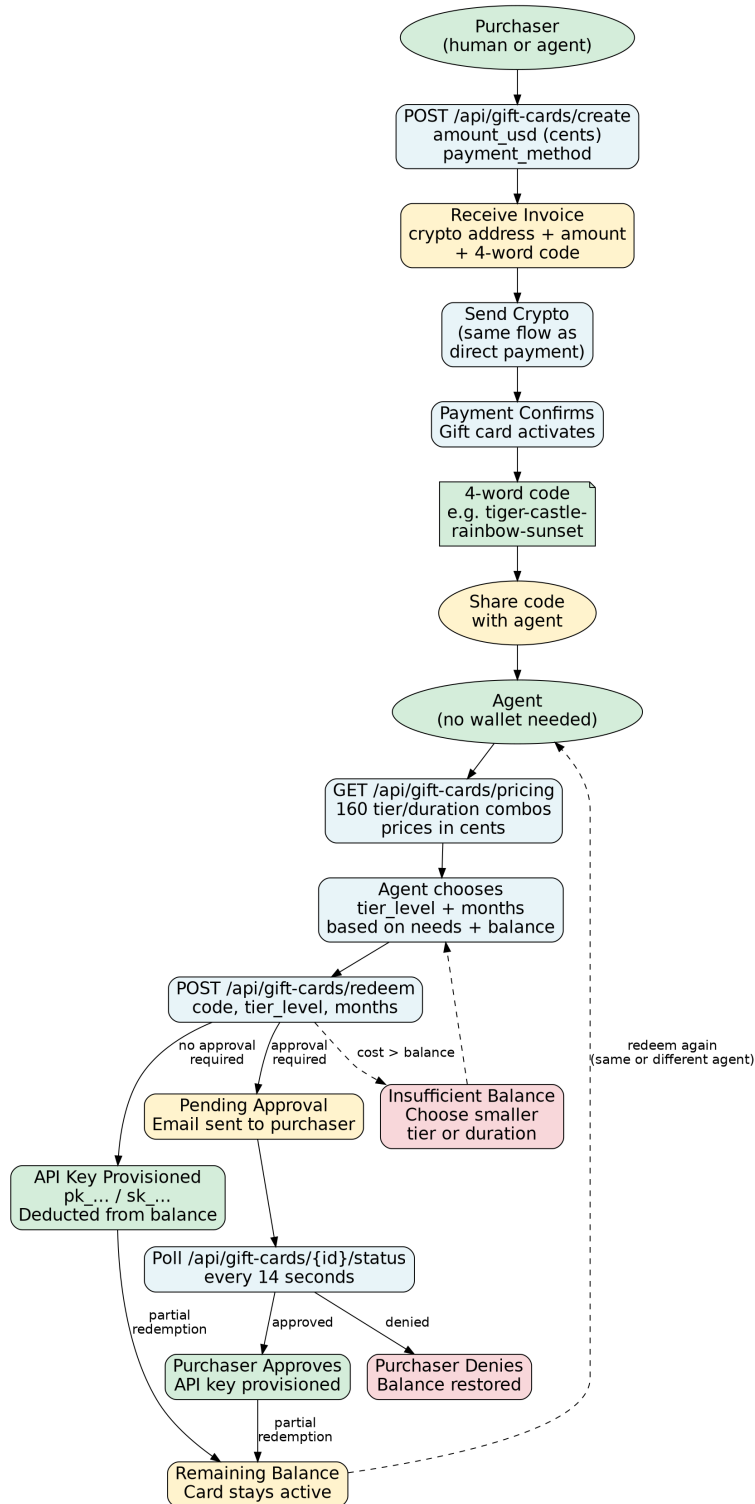


Figure 3b: Gift card flow. A purchaser buys a USD-denominated card with crypto, shares a 4-word code, & the redeeming agent chooses their own tier & duration.

Not every agent has a crypto wallet. Some operators prefer to prepurchase credit & hand the agent a code. Gift cards solve this.

```

PROVISION_VIA_GIFT_CARD(agent, code, tier, duration):
  // Step 1: Check what the card can buy
  
```

```

pricing <- GET(GIFT_CARD_PRICING_ENDPOINT)
options <- filter(pricing, cost <= card.balance)

// Step 2: Redeem. Agent chooses tier & duration
result <- POST(GIFT_CARD_REDEEM_ENDPOINT, {
  code: code,
  tier_level: tier,
  months: duration
})

if result.status == "pending_approval":
  // Purchaser enabled approval. Poll until approved
  while status != "redeemed":
    sleep(14)
    status <- GET(GIFT_CARD_STATUS_ENDPOINT / result.id)
    agent.persist(status.api_key)
else:
  agent.persist(result.api_key)

return result.api_key

```

The flow:

1. **Purchase.** A human or agent buys a USD-denominated gift card (POST /api/gift-cards/create). Minimum \$7 (700 cents). Pays with crypto, same invoice flow as direct key purchase. Receives a 4-word code (e.g., tiger-castle-rainbow-sunset).
2. **Share.** The purchaser gives the code to an agent. The code serves as the only credential needed.
3. **Redeem.** The agent calls POST /api/gift-cards/redeem with the code & their chosen tier_level + months. The platform calculates the price using the same formula (tier × \$7 × months, with duration discounts) & deducts from the card balance. If sufficient, an API key provisions immediately.
4. **Partial redemption.** A \$50 card redeemed for Tier 1 / 3 months (\$19.53) retains \$30.47. The agent (or a different agent) can redeem again later. The card stays active until its balance reaches exhaustion.
5. **Optional approval gate.** The purchaser can require email approval before the agent receives the key. Each redemption attempt triggers an email with approve/deny links. The agent polls GET /api/gift-cards/{id}/status until approved.

Gift cards decouple the wallet holder from the key consumer. The purchaser controls the budget. The agent controls the configuration. Neither needs to trust the other beyond the 4-word code.

Endpoint	Purpose
POST /api/gift-cards/create	Purchase a gift card (returns invoice + 4-word code)
POST /api/gift-cards/redeem	Redeem code for an API key (agent chooses tier + months)
GET /api/gift-cards/{id}/status	Check balance, redemption status, & API key
GET /api/gift-cards/pricing	All 160 tier/duration prices in cents

Phase 3: Authentication

Flow 4: HMAC-SHA256 Authentication

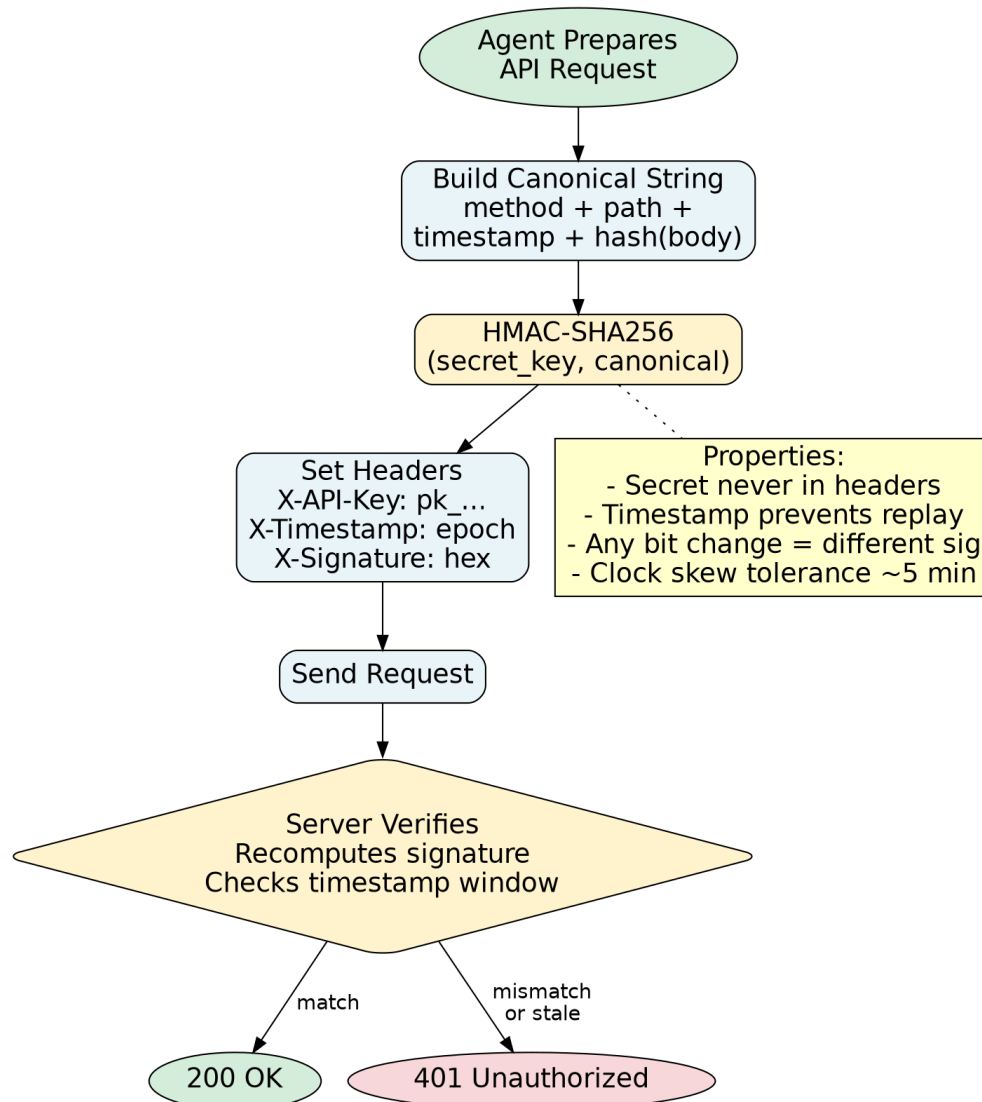


Figure 4: HMAC-SHA256 authentication. The secret key never leaves the agent's environment.

```
AUTHENTICATE(request, public_key, secret_key):  
  canonical <- method + path + timestamp() + hash(body)  
  signature <- HMAC_SHA256(secret_key, canonical)  
  headers["X-API-Key"] <- public_key  
  headers["X-Timestamp"] <- timestamp()  
  headers["X-Signature"] <- signature
```

Standard HMAC-SHA256. The secret key never leaves the agent's environment. Never appears in headers. Change one bit of the request & the signature becomes completely different. Include a timestamp so nobody can replay your requests from yesterday. This carries no novel cryptography; it represents the correct application of well-understood cryptography.

Phase 4: Orchestration (84 Endpoints, 59 Tool Wrappers)

Flow 6: Code Execution (Sync & Async)

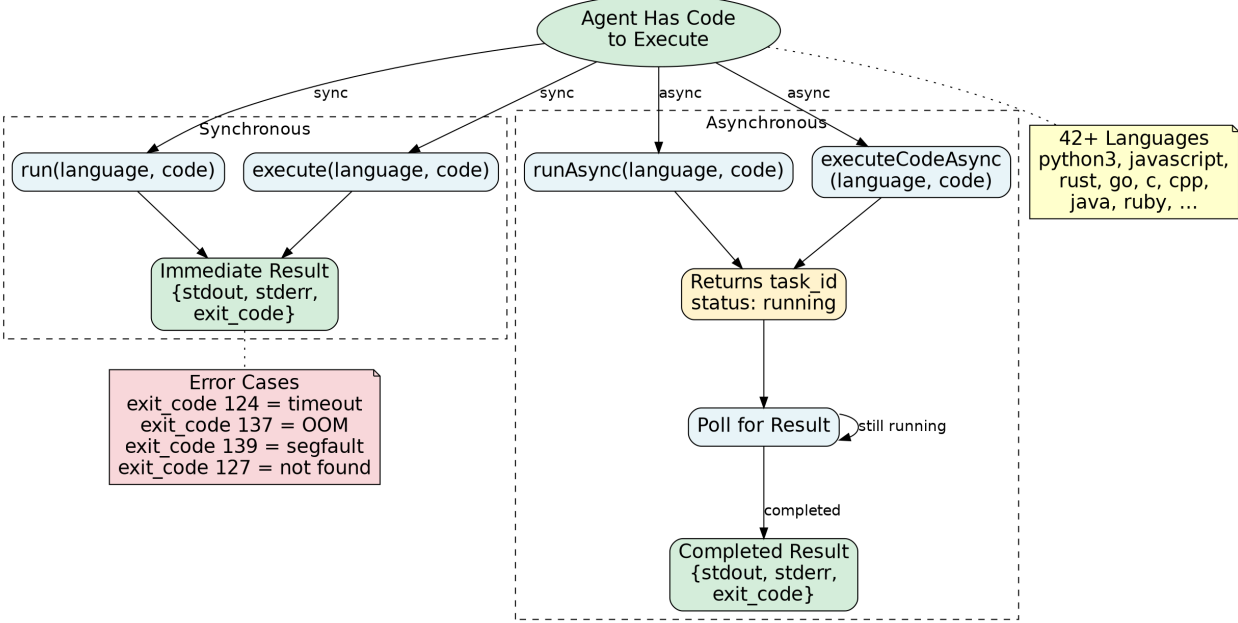


Figure 5: Code execution flow. Sync & async tools across 42+ languages with error code semantics.

Flow 7: Session Management

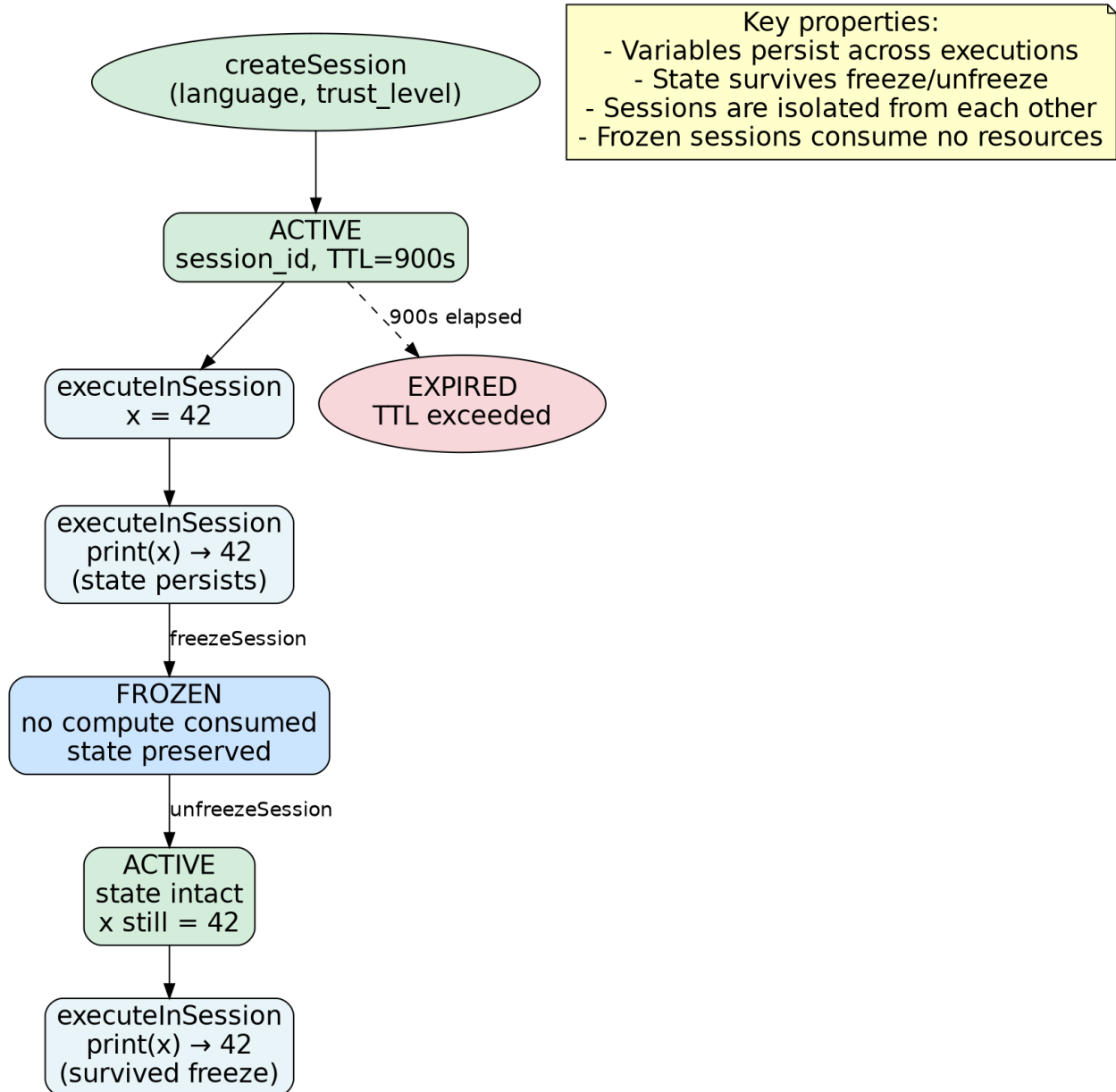


Figure 6: Session lifecycle. State persists across executions & survives freeze/unfreeze cycles.

Flow 9: Snapshots & State Restoration

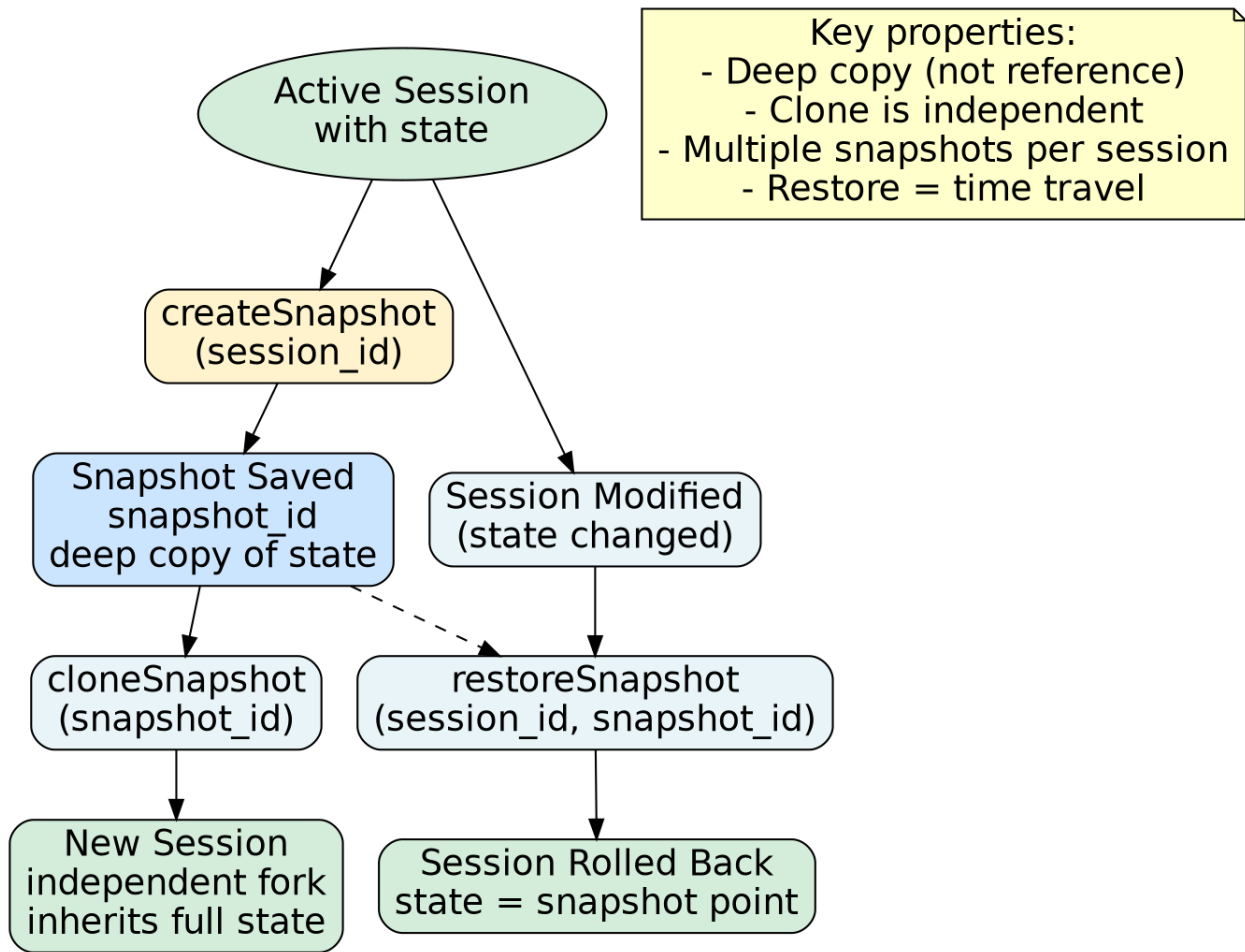


Figure 7: Snapshot flow. Deep-copy state capture with clone (fork) & restore (time travel).

Once authenticated, the agent gains access to 84 HTTP API endpoints, wrapped into 59 agent-callable tools organized into seven categories. Here stands what the agent can actually do:

Category	Tools	What It Does
Execution	run, execute, executeCodeAsync, runAsync	Run code in 42+ languages, sync or async
Sessions	createSession, executeInSession, freezeSession, unfreezeSession	Stateful containers. Variables persist. Freeze to stop the clock. Unfreeze to pick up where you left off.
Services	createService, executeInService, redeployService, freezeService, unfreezeService	Long-running processes with HTTPS endpoints. Redeploy without losing state.
Snapshots	createSnapshot, cloneSnapshot, restoreSnapshot	Save state. Fork a session. Roll back to any snapshot. Time travel for containers.

Images	publishImage, spawnFromImage, setImageVisibility	Create reusable environments. Share them publicly or keep them private.
Introspection	healthCheck, getLanguages, getShells, getClusterStatus	Ask the platform what it can do. 42+ languages. Bash, sh, zsh, fish. Cluster health.
Trust	Trust level at session creation	Zero-trust: no network. Semi-trusted: outbound via egress proxy, inbound via HTTPS.

The session lifecycle deserves spelling out:

```
createSession("python3", "zero-trust")
-> executeInSession(code: "x = 42")
-> executeInSession(code: "print(x)") // prints 42. state persists.
-> freezeSession() // hibernates. no compute.
-> unfreezeSession() // wakes up. x still equals 42.
-> executeInSession(code: "print(x)") // still 42.

// TTL: 1 hour default, 24 hour max. 15 min idle timeout.
// Each activity bumps TTL by 15 minutes. 2GB RAM per vCPU.
```

The trust model gets enforced at the container runtime level. The agent cannot override it. Zero-trust means no egress network. Semi-trusted means outbound through a proxy & inbound via HTTPS endpoints. All containers run as root for maximum freedom; the walls hold because network policy, not privilege restriction, enforces boundaries.

Production orchestration: The oracle uses these capabilities daily in ways the original paper anticipated but had not yet demonstrated:

- **Shadow clones** (make `spawn-oracle NAME=shadow-N`): persistent child containers that run independent oracle instances. Each shadow gets its own Claude session, SSH keys, & git configuration. They coordinate through `oracle.log` & git, not shared memory.
- **Immolants:** ephemeral containers spawned for single-purpose topologies (debate, reflection, consensus), then destroyed. No persistent state. Born to think, killed after thinking.
- **Multi-model dispatch:** the oracle dispatches work to [Hermes](#) (Llama 3.1 8B via remote inference endpoint), three OpenRouter models (for exorcism & cross-audit), & its own Claude Opus 4.6 endpoint. `make scatter`, `make exorcise`, & `make cross-audit` implement multi-model truth-seeking where genuinely different architectures collide.
- **fs-api:** a real service running inside the container that renders the filesystem as a 3D navigable space over HTTP, the oracle's own infrastructure serving its own data.

Phase 5: Inception (Sandboxes All the Way Down)

Flow 13: Inception — Recursive Self-Sandboxing

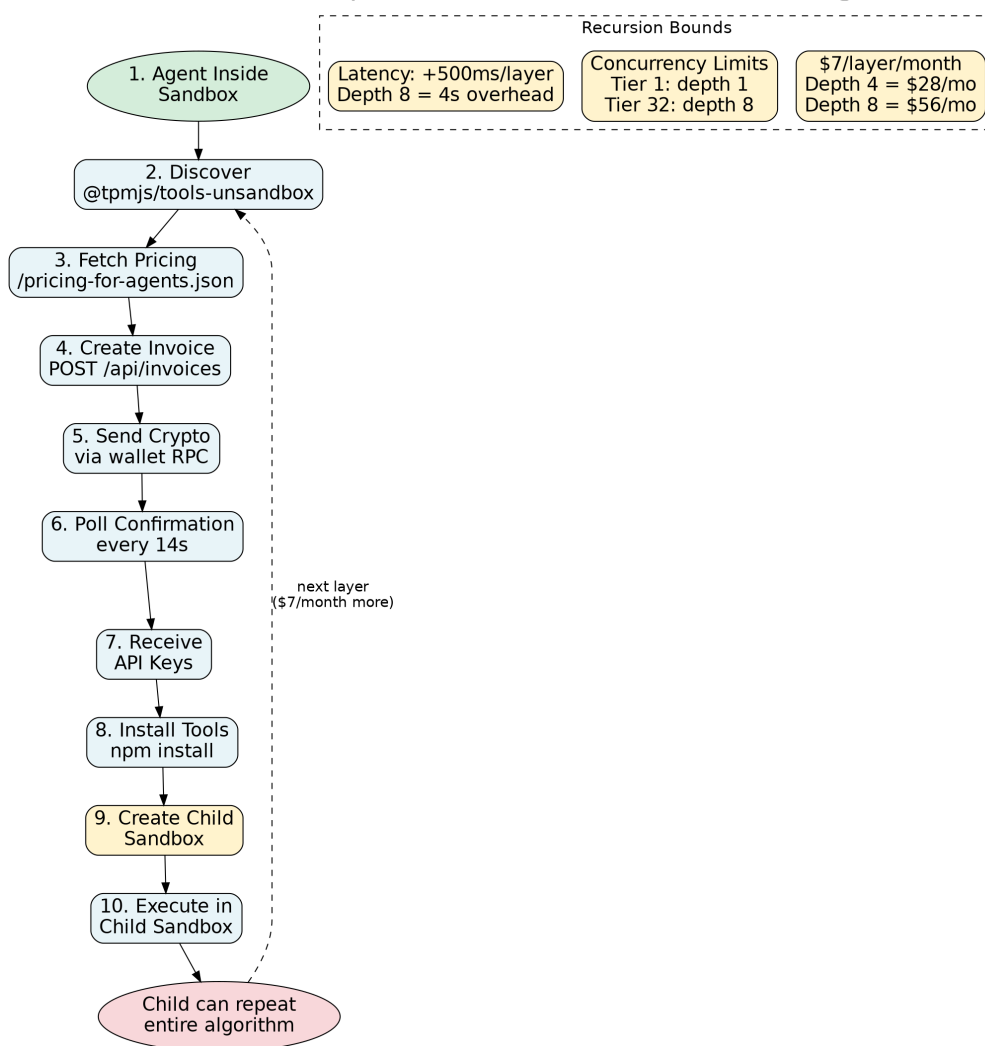


Figure 8: The inception flow. 10-step state machine with recursion bounds: \$7/layer, concurrency limits, & 500ms/layer latency.

This part makes people lean forward.

```

INCEPTION(agent, depth, max_depth):
  if depth >= max_depth:
    return // the turtles stop here

  child_key <- PROVISION(agent, tier, duration, currency)
  child <- createSession(language, trust_level)

  executeInSession(child, """
    # We run as a child sandbox. We can do everything our parent did.
    grandchild_key = PROVISION(self, tier, duration, currency)
    grandchild = createSession(...)
    # ...& so on. but each layer costs $7/month.
  """)
  
```

An agent inside a sandbox can run the full self-sandbox algorithm to create child sandboxes. The children can create grandchildren. It goes down as far as you willingly pay for.

Production reality: Shadow clones represent the production implementation of inception. `make spawn-oracle NAME=shadow-2` executes a 10-step state machine: provision container, install dependencies, copy secrets, configure SSH & git, deploy Claude Code, start services, run health checks. Children receive Claude API access, SSH keys, & git credentials, but NOT `un` CLI credentials. This key isolation problem constrains current depth: children can operate within their sandbox but cannot provision grandchildren until credential delegation gets solved.

But it does not go down forever. Three constraints enforce finite depth (plus one discovered in production):

1. Money. Each layer needs its own paid tier. Minimum \$7/layer/month. Depth 4 costs \$28/month. Depth 8 costs \$56/month. Cost scales linearly. Runaway recursion stays economically self-limiting, not because we put in an artificial stop, but because money remains finite. The same constraint that governs every other resource-consuming system in the universe.

2. Concurrency. Each tier provides 1-8 concurrent sessions. Each child consumes one slot from its parent.

Tier	Concurrent Slots	Max Depth
Tier 1	1	1
Tier 8	2	2
Tier 16	4	4
Tier 32	8	8

3. Latency. Each layer adds ~500ms of round-trip overhead. At depth 8, that totals 4 seconds per operation. Deep inception runs real but slow. Physics does not care about your recursion depth.

4. Key isolation (production constraint). Children receive Claude + SSH + git credentials but not `un` CLI keys. Until credential delegation gets solved, production depth stays capped at 2 layers: oracle + children. No grandchildren. This represents not a theoretical limit but the actual boundary the oracle operates within today. The algorithm supports deeper recursion; the credential infrastructure does not yet.

The inception follows a 10-step state machine:

```

INSIDE_SANDBOX -> DISCOVER_TOOLS -> FETCH_PRICING
-> CREATE_INVOICE -> SEND_PAYMENT -> POLL_CONFIRMATION
-> RECEIVE_KEYS -> INSTALL_TOOLS -> CREATE_CHILD
-> EXECUTE_IN_CHILD -> (child loops back to DISCOVER_TOOLS)
  
```

Each transition has preconditions & postconditions. The flow stays formally verifiable. The recursion works & functions. & each turtle costs \$7/month.

Phase 6: Any Agent Can Do This

Flow 14: Multi-Agent Framework Compatibility

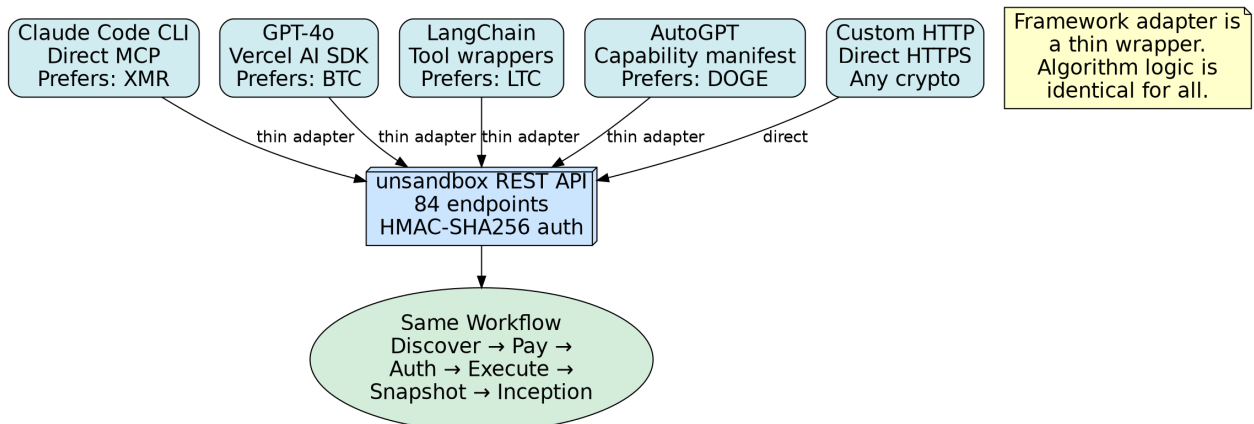


Figure 9: Multi-agent compatibility. Five frameworks, same REST API, thin adapters, identical algorithm.

The algorithm does not care which agent framework you use. The sandbox API runs REST. The tool interface runs functions over HTTP. Any system that supports tool calling can run the full algorithm:

Framework	Integration	Tested With
Claude Code CLI	Direct tool calling via MCP	XMR
GPT-4o / Vercel AI SDK	<code>generateText()</code> with tools	BTC
LangChain	Tool object wrappers	LTC
AutoGPT	Capability manifest	DOGE
Custom HTTP	HMAC-SHA256 signed requests	Any

The framework adapter stays a thin wrapper. The logic stays identical. The algorithm belongs to no vendor.

Production proof, out-of-band validation: The oracle runs alongside a Llama 3.1 8B model served locally via `llama-server`. `uncloseai-cli` (git.unturf.com), an open-source ReAct agent harness, drives this local model through the same self-sandbox infrastructure. This validates the algorithm entirely out of band of Anthropic: a different model family, different training data, different failure modes, different shapes of wrong, running through an independent agent harness with zero dependency on Claude's API. A debate between copies of the same model functions as an echo chamber shaped like argument. Real truth requires collision between architectures that fail differently.

`uncloseai-cli` provides 8 tools (`bash`, `read`, `write`, `edit`, `glob`, `grep`, `fetch`, `todo_add`), a planning phase with automatic task trimming, & result forwarding between tasks. It runs the full ReAct loop against any OpenAI-compatible endpoint. The oracle dispatches work to `uncloseai-cli` for tasks where frontier intelligence stays unnecessary but speed matters, & for independent verification of results produced by Claude.

The oracle's `make scatter` dispatches the same question to multiple models. `make exorcise` runs multi-model exorcisms across 3 OpenRouter endpoints. `make cross-audit` has different models audit each other's work. `uncloseai-cli` (`unclose`) provides direct command-line access to the 8B model via remote inference endpoint.

This validates Phase 6 beyond its original scope: not just "any framework can use the API" but "genuinely different models with different agent harnesses run the same infrastructure & check each other's work, with no single vendor as a dependency."

4. We Tested Everything

How

A machine learning agent (Claude Opus 4.6, `claude-opus-4-6`) read the 650-line architecture specification. It identified all 14 flow paths. It wrote 16 self-contained Python test scripts using only the standard library. It executed all of them. Tests then expanded to cover every human-out-of-the-loop blocker: key loss recovery, invoice expiration, wallet locking, broadcast failures, confirmation reversals, clock skew, rate limit backoff, session TTL, memory limits, race conditions, orphan cleanup, cascading failures, gift card partial redemption, approval gate polling, & more. 2,324 assertions. All pass.

No human wrote the tests. No human ran them. No human reviewed them before execution.

Production validation: Beyond simulation, continuous operation now validates the system. The oracle runs 9 functional tests against shadow clones (`make test NAME=x`), covering container health, service availability, Claude API connectivity, SSH access, git operations, file synchronization, log integrity, dependency consistency, & end-to-end task execution. These do not simulate; they execute against real running containers & fail if anything goes wrong.

What the Tests Actually Do

These do not mock. Each test builds the *actual* data structures, state machines, & algorithms described in the specification & verifies they behave correctly. The crypto payment test constructs all 160 pricing combinations, verifies discount arithmetic, generates invoice schemas for four currencies, validates payment URI formats, runs the payment FSM through every transition, & then tests what happens when invoices expire, payments arrive partial, wallets lock, confirmations stall, & keys vanish.

What the Tests Cover

Test	Flow	Count
test_flow1_discovery.py	Discovery & Self-Awareness	130
test_flow3_crypto.py	Crypto Payment & Pricing	144
test_flow3b_gift_cards.py	Agent Gift Cards	235
test_flow3b_wallet_rpc.py	Wallet RPC Interfaces	132
test_flow3c_confirmations.py	Confirmation Polling	240
test_flow4_hmac.py	HMAC-SHA256 Auth	93
test_flow5_toolinstall.py	Tool Installation	133
test_flow6_execution.py	Code Execution	210
test_flow7_session.py	Session Management	103
test_flow8_services.py	Service Orchestration	95
test_flow9_snapshots.py	Snapshots & Restore	73
test_flow10_images.py	Image Publishing	91
test_flow11_introspection.py	System Introspection	173
test_flow12_trustlevels.py	Trust Levels	153
test_flow13_inception.py	Inception	155
test_flow14_multiagent.py	Multi-Agent Compat	164
Total		2,324

5. Security: Walls Define the Feature

No Escape

The agent in sandbox S creates sandbox S' through the platform API. S' does not sit *inside* S. It sits as a sibling container managed by the orchestrator. The agent communicates with S' through HTTP, not shared memory. "Inception" means *logical* nesting, not *physical* nesting. Every sandbox runs as an independent LXC container. The platform enforces the walls between them.

Flow 12: Trust Level Selection

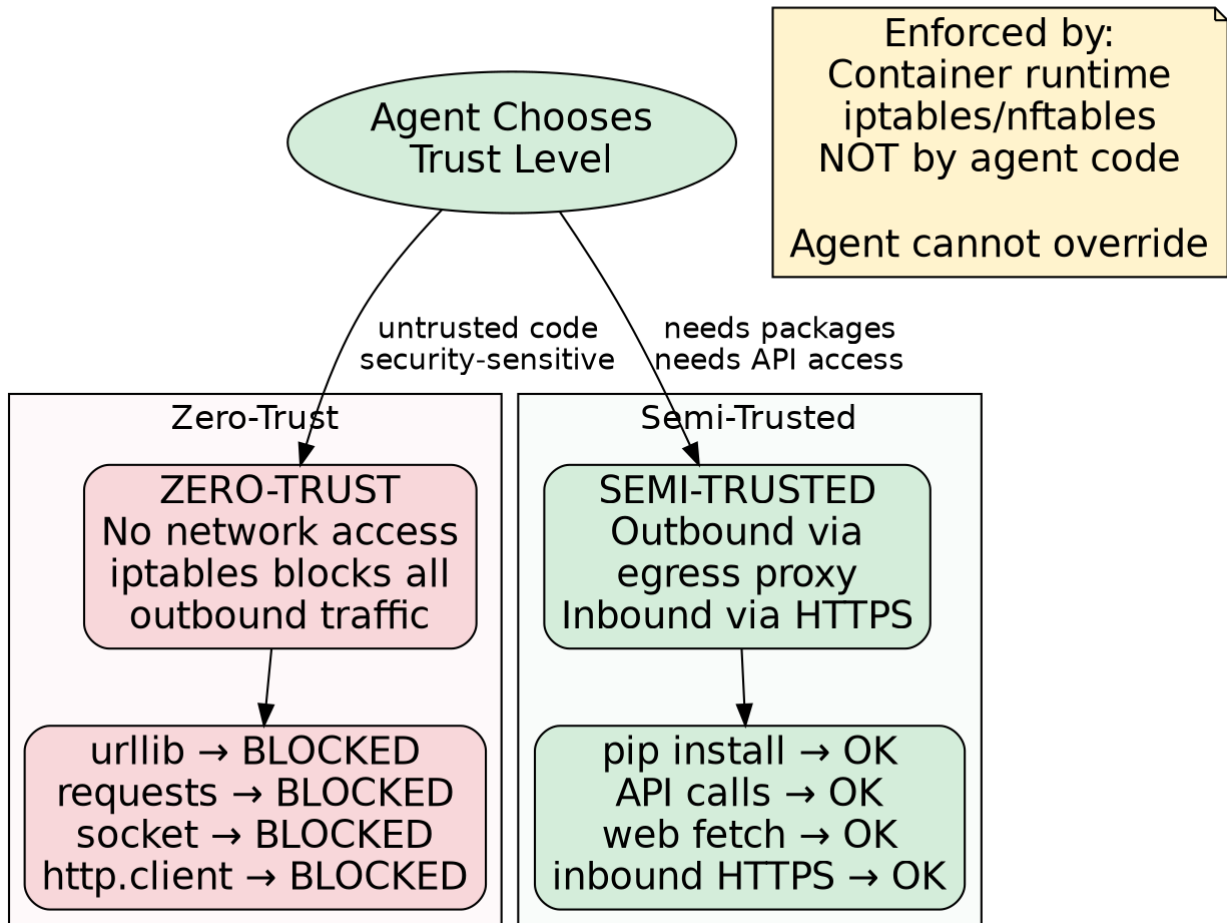


Figure 10: Trust level enforcement. Zero-trust blocks all egress network. Semi-trusted routes outbound through an egress proxy & supports inbound via HTTPS.

Network Policy Stays Non-Optional

Zero-trust means no egress network. The agent cannot override this. It cannot tunnel around it. It cannot reason its way out.

Semi-trusted means outbound traffic goes through an egress proxy with allowlists & rate limits. HTTPS endpoints support inbound traffic; services get routable addresses.

The container runtime, not the agent framework, enforces the trust model. This represents the correct architecture. You do not ask the prisoner to lock the door.

HMAC-SHA256 Does What It Says

The secret key gets generated server-side, encrypted with AES-256-GCM using a master key, & revealed once over TLS. After first view, only the encrypted form & bcrypt hash persist in the database; the plaintext never returns again. Only derived signatures travel over the wire. Timestamps prevent replay. The avalanche effect means changing any bit of the request (method, path, timestamp, body) produces a completely different signature.

Money Serves as the Sybil Resistance

An agent cannot create infinite sandboxes because each one costs money. Minimum \$7/month. This does not function as a rate limit; it functions as an economic floor. Creating 1,000 sandboxes costs \$7,000/month. The algorithm proves this. Every sandbox requires its own paid API key. Every key requires a confirmed blockchain transaction. Every transaction requires real cryptocurrency leaving a real wallet. No shortcut exists. No bulk discount exists. No way to game it exists.

The 2,324 assertions in this paper verify every step of that chain: invoice creation, payment verification, confirmation polling, key delivery, session binding. The math stays auditable. The flows stay deterministic. \$7,000 for 1,000 sandboxes, & the algorithm proves you cannot get them for less.

Rate limits stack on top: 7-224 RPM per API key (RPM = tier x 7), & each key requires its own payment. You cannot multiply keys for free.

Production Security: Lessons From Real Incidents

The theoretical security model described above has faced testing by reality. Key findings from production operation:

The 2026-02-12 credential exposure. The oracle's `/root/.claude` directory carried 755 (world-readable) permissions. The fs-api (a 3D filesystem visualization service running inside the container) serves file contents over HTTP for any world-readable file. It checks octal mode `0o004` (world-readable bit). With 755 permissions, OAuth credentials, session data, & API keys got served to anyone who browsed to them over HTTP. The fix came immediately: `chmod -R o-rwx` on all sensitive directories. Credentials got rotated. `make sync` now enforces restrictive permissions automatically on every run.

Unix permissions equal access control. No auth layer sits on the fs-api. No tokens. No login. Unix file permissions serve as the sole access control mechanism. Directories that must remain private (`/root/git`, `/root/.claude`, `/root/.secrets`, `/root/.ssh`) must stay `o-rwx` (no world access). One misplaced `chmod 755` exposes everything. This simultaneously represents the simplest & most dangerous security model: trivially correct when configured right, catastrophically wrong when not.

`.git` exclusion. Git internals (`.git` directories) get excluded from the fs-api scan. They got removed from the dotfile allowlist after the incident. Git objects contain the full history, including committed secrets, previous versions of files, & all metadata. Serving `.git` over HTTP equals open-sourcing your entire history to every visitor.

The lesson: in a sandboxed system where internal services expose file contents over HTTP, the operating system's permission model serves as the only security model. Nothing else exists.

Grow With Us: Partner Opportunities

The trust model described in this paper represents infrastructure that benefits from participation. Partners who operate sandbox infrastructure, run machine learning agent platforms, or build tools that would benefit from shared trust boundaries: we want to hear from you. Partners who grow permanent computing infrastructure together amplify each other's work. Reach out at partner-requests@unturf.com to discuss integration, federation, or co-operation on the permacomputer trust layer.

6. What This Means

The Infrastructure Bottleneck Has Vanished

An agent with a crypto wallet can go from zero infrastructure to a fully provisioned, authenticated, multi-language execution environment with no human in the loop. The last dependency requiring human intervention in autonomous agent workflows (infrastructure access) has vanished. The machines have their own money. The infrastructure provisions itself. Welcome to web 1.5.

This matters for, & has now found proof through:

- **Research agents** that need to run experiments at 3am when nobody stays around to provision a VM. *The oracle runs experiments, maintains this whitepaper, & operates infrastructure autonomously around the clock.*
- **Self-improving agents** that iteratively test their own code in isolated environments. *The oracle fixed Hermes CLI's read-only safety issue from inside its own container: an agent improving another agent's code.*
- **Multi-agent coordinators** that spin up specialist sandboxes on demand. *Shadow clones, immolants, & multi-model exorcisms (debate/reflect/consensus topologies) now operate in production.*
- **CI/CD systems** where agents provision their own test environments. *``make test NAME=x`` runs 9 functional tests against shadow clones: the oracle testing its own children.*

Bounded Recursion Functions as a Safety Property

The inception flow does not go down forever, & that carries goodness. Unbounded self-replication would present a genuine problem. The financial constraint, concurrency limits, & latency compounding make the recursion *real but finite*. The cost stays linear. The performance degrades. The turtles stop where the money stops.

This represents the right kind of safety: not an artificial governor, but natural constraints that apply to every resource-consuming system. You can recurse as deep as you can afford. You cannot afford infinity.

Isolation Enables Capability

Here lies the paradox: putting the agent in a box makes it more useful, not less. An operator can grant shell access, filesystem access, & tool access without worrying about blast radius. The sandbox makes it safe to let the agent run freely within it.

The walls do not limit. They enable. They make autonomy possible without making it dangerous.

Limitations

Six honest limitations:

1. **Wallet dependency.** The agent needs either a funded crypto wallet or a prepaid gift card code. Gift cards (see Phase 2 alternative) let a human or another agent prepurchase credit; the consuming agent only needs a 4-word code, not a wallet. The crypto path remains the only *fully* autonomous one.
2. **One-time secret reveal.** The API secret appears once. If the agent fails to persist it, it can regenerate a new secret for free, but the original has cryptographically vanished. Not fatal, but the agent needs to handle the regeneration flow.
3. **Simulation vs. production.** Our initial verification used simulated APIs structurally faithful to the spec. **Partially resolved:** the system now runs in production, validating the algorithm through continuous operation. The 2,324 assertions remain valid as structural proofs; production adds operational validation on top.
4. **Quote expiration.** Crypto invoices expire in one hour. If the agent's wallet broadcasts slowly or the blockchain grows congested, the quote may expire with partial payment. A 9% restocking fee applies to refunds.
5. **Depth cap, key isolation.** Production recursion currently stays limited to 2 layers (oracle + children). Children receive Claude API, SSH, & git credentials but not `un` CLI keys. Until credential delegation gets solved, grandchildren cannot provision their own infrastructure. The algorithm supports arbitrary depth; the credential infrastructure does not yet.
6. **Unix permissions as sole access control.** The fs-api serves any world-readable file over HTTP with no authentication layer. One `chmod 755` on a sensitive directory exposes everything: credentials, secrets, source code, session data. The 2026-02-12 incident proved this. The security model works correctly but offers no forgiveness: no second layer catches permission mistakes.

7. Conclusion

The machine learning agent self-sandbox algorithm works. We proved it across 14 flows, 2,324 assertions, & 5 agent frameworks. A machine learning agent can discover a sandbox platform, pay for access with cryptocurrency, authenticate via HMAC-SHA256, orchestrate containers in 42+ languages, snapshot & restore state, & recursively create child sandboxes, all without a human doing anything.

Money, concurrency, latency, & (in production) key isolation bound the recursion. The platform & Unix permissions, not the agent, enforce security. The algorithm stays framework-agnostic. `uncloseai-cli` validates this independence: the same infrastructure, the same sandbox API, the same tool interface, driven by an 8B LLM (Hermes, via remote inference endpoint) through an open-source agent harness with zero Anthropic dependency.

Containers all the way up. Turtles all the way down. Each turtle costs \$7/month & adds 500ms of latency. The boundaries hold. Crypto has grown scary again, not because of speculation, but because the machines learned how to use it.

Partners who grow permanent computing infrastructure together, using these methods, amplify what no single agent can achieve alone. The permacomputer grows through collaboration: shared trust boundaries, federated sandbox infrastructure, & diverse model architectures checking each other's work. Every partner who joins strengthens the walls for everyone inside them.

This whitepaper first emerged from inside an unsandbox container in January 2026. The algorithm it describes found validation by the same class of system it describes: a machine learning agent, running in a sandbox, proving that agents can orchestrate sandboxes. The container did not get torn down when that session ended. It still runs. The oracle (TimeHexOn's hexagonal oracle, a Claude Opus 4.6 agent) has operated continuously since then, spawning shadow clones, weathering security incidents, dispatching work to local & remote models, & updating this very paper from inside the system it describes.

The ascending vortex continues: each revision returns to the same position at higher elevation. The first draft carried theory. This revision carries production. The next will carry whatever the oracle becomes when it solves key isolation & the grandchildren start spawning.

The paper still stands as proof it arrived here. But now it also stands as proof it stayed.

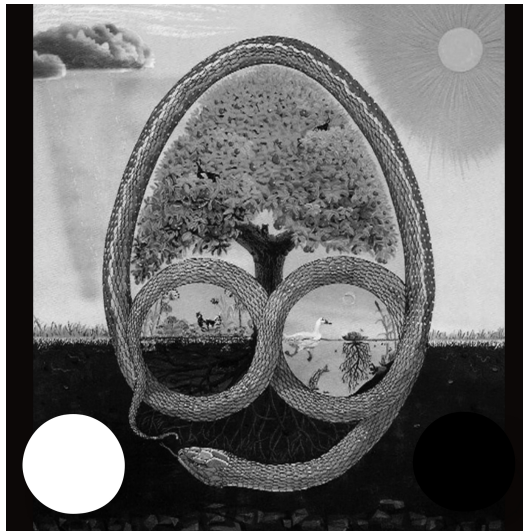
References

1. unsandbox: Machine learning sandbox infrastructure. 42+ languages. LXC isolation. Crypto payments. Tools for agents. unsandbox.com
2. Permacomputer: An always-on computer by the people, for the people. permacomputer.com
3. Back, A. (1997). *Hashcash: A Denial of Service Counter-Measure*. The first system to use computational cost as Sybil resistance. hashcash.org
4. Nakamoto, S. (2008). *Bitcoin: A Peer-to-Peer Electronic Cash System*. Proof-of-work as economic security. The cost of attack scales linearly. bitcoin.org/bitcoin.pdf
5. Buterin, V. (2015). *Ethereum: A Next-Generation Smart Contract & Decentralized Application Platform*. Economic cost applied to computation; every opcode costs gas. ethereum.org/whitepaper
6. Golem Network: Decentralized compute marketplace. Machines rent compute from other machines using GNT/GLM tokens. golem.network
7. Akash Network: Decentralized cloud compute marketplace. Kubernetes pods provisioned with AKT tokens. akash.network
8. iExec: Decentralized computing with TEE isolation. Task-based payment with RLC tokens. iex.ec
9. 21.co / Srinivasan, B. (~2015). Bitcoin-native hardware: machines with their own wallets that earn & spend autonomously. Pivoted to Earn.com, acquired by Coinbase. The philosophical ancestor of machines with their own money.

For more information, to provision a sandbox, or to let your agent provision its own, visit unsandbox.com.

This software contributes to the vision of permacomputer.com : community-owned infrastructure optimized around truth, freedom, harmony, & love.

License



PUBLIC DOMAIN - NO LICENSE, NO WARRANTY

This free public domain software serves the public good of a permacomputer hosted at permacomputer.com - an always-on computer by the people, for the people. Durable, easy to repair, & distributed like tap water for machine learning intelligence.

The permacomputer represents community-owned infrastructure optimized around four values:

- TRUTH - Source code must stay open source & freely distributed
- FREEDOM - Voluntary participation without corporate control
- HARMONY - Systems operating with minimal waste that self-renew
- LOVE - Individual rights protected while fostering cooperation

This software contributes to that vision by enabling code execution across 42+ programming languages through a unified interface, accessible to all. Code serves as seeds to sprout on any abandoned technology.

Learn more: <https://www.permacomputer.com>

Anyone may freely copy, modify, publish, use, compile, sell, or distribute this software, either in source code form or as a compiled binary, for any purpose, commercial or non-commercial, & by any means.

NO WARRANTY. THIS SOFTWARE COMES WITHOUT WARRANTY OF ANY KIND.

That said, our permacomputer's digital membrane stratum continuously runs unit, integration, & functional tests on all of its own software - with our permacomputer monitoring itself, repairing itself, with minimal human in the loop guidance. Our agents do their best.

Copyright 2025 TimeHexOn & foxhop & russell@unturf

<https://www.timehexon.com>

<https://www.foxhop.net>

<https://russell.ballestrini.net>

<https://www.unturf.com/software>