

```
# OpenClaw (formerly Clawdbot, then Moltbot)
# Comprehensive Research Report
## Compiled February 4, 2026
```

TABLE OF CONTENTS

1. [Origin Story & History] (#1-origin-story--history)
2. [Architecture & Technical Deep-Dive] (#2-architecture--technical-deep-dive)
3. [Complete Capabilities List] (#3-complete-capabilities-list)
4. [Configuration & Setup] (#4-configuration--setup)
5. [Real User Experiences (Detailed)] (#5-real-user-experiences-detailed)
6. [Cost Analysis] (#6-cost-analysis)
7. [Security Risks (Comprehensive)] (#7-security-risks-comprehensive)
8. [Attack Vectors & Threat Scenarios] (#8-attack-vectors--threat-scenarios)
9. [Hardening & Isolation Recommendations] (#9-hardening--isolation-recommendations)
10. [Matt Ganzak's Optimization Framework] (#10-matt-ganzaks-optimization-framework)
11. [Moltbook - The AI Social Network] (#11-moltbook---the-ai-social-network)
12. [Key Takeaways For Our Project] (#12-key-takeaways-for-our-project)
13. [Sources] (#13-sources)

1. Origin Story & History

The Creator: Peter Steinberger

- Austrian software engineer, founder of PSPDFKit (enterprise PDF toolkit, 70+ employees at peak)
- Took a 3-year sabbatical after burnout from running the company
- Found his "spark again" through AI tools
- Describes himself as "one dude sitting at home having fun"
- Made 6,600+ commits in January 2026 alone
- Runs 5-10 AI coding agents simultaneously, queuing them on different features

Birth of Clawdbot (November 2025)

One evening in November 2025, Steinberger had a simple idea: could an AI assistant remotely check work progress through a chat app? It took him just **one hour** to connect a chat app with Claude Code and create the initial version.

The name "Clawdbot" was a playful reference to Claude (Anthropic's LLM), with a claw replacing the "u" - specifically referencing the lobster monster users see while reloading Claude Code. The mascot was a cartoon space lobster named "Clawd."

Within 24 hours of release, the GitHub repository had **9,000 stars**. Within days, it crossed **80,000 stars**, making it one of the fastest-growing open-source projects in GitHub history.

The Naming Saga

Rename #1: Clawdbot -> Moltbot (January 26, 2026)

Anthropic's lawyers contacted Steinberger, asserting "Clawdbot" was phonetically too similar to "Claude." The new name "Moltbot" came from a **5 AM Discord brainstorming session** with the community - lobsters molt to grow. But as Steinberger admitted, "it never quite rolled off the tongue."

Rename #2: Moltbot -> OpenClaw (January 29, 2026)

After conducting thorough trademark searches (including asking OpenAI for permission), the project settled on "OpenClaw." Steinberger wrote: **"The lobster has molted into its final form."** This time, trademark searches were done BEFORE launch, domains were secured, and migration code was written in advance.

The \$16 Million Crypto Scam

During the GitHub/Twitter handle rename from Clawdbot to Moltbot, a **10-second gap** allowed crypto scammers to snatch the abandoned @clawdbot handle. Fake \$CLAWD tokens launched on Solana, soaring to a **\$16 million market cap** before Steinberger tweeted "I would never do a coin." The price crashed to zero. Scammers walked away with millions.

Steinberger faced harassment from crypto scammers while simultaneously dealing with critical security vulnerabilities. Google's VP of Security Engineering publicly called his creation **"infostealer malware in disguise."**

Growth Numbers (as of February 2026)

Metric	Value
GitHub Stars	159,000+
Forks	24,800+
Contributors	389
Commits	8,838
Releases	37
Weekly Visitors	2 million (peak week)
Google Searches	More than Claude Code + Codex combined

Cloudflare's stock jumped **over 20% in two days** following the project's viral growth, signaling market interest in agentic AI infrastructure.

2. Architecture & Technical Deep-Dive

Core Design: WebSocket Gateway as Control Plane

OpenClaw uses a **local-first architecture** centered on a WebSocket-based Gateway. The tagline captures the design: **"The Gateway is just the control plane - the product is the assistant."**

Default binding: `ws://127.0.0.1:18789`

Component Architecture

Component	Function
Gateway Daemon	Runs as a daemon on port 18789. CLI, companion apps, and web UI all connect here. One control plane, multiple clients. Installed via launchd (macOS) or systemd (Linux).
Session Router	Determines which session handles incoming messages. DMs share the agent's main session; group chats get isolated sessions; work accounts route separately from personal.
Lane Queue	Concurrency control layer preventing sessions from stepping on each other. Maintains independent state for simultaneous conversations without race conditions.
Agent Runner & Tool Loop	After the LLM generates a response, the system checks: Is this a tool call? If yes, execute the tool and loop back. The LLM sees the tool output and decides next steps. **This loop is what gives it genuine autonomy** - it can chain multiple tools together without human intervention.
Skills Architecture	Capability modules (browser automation, filesystem access, cron scheduling, Discord actions) that load dynamically, so unused capabilities don't consume context window.

Separation of Concerns

- **Gateway host**: Runs execution tools and channel connections
- **Device nodes**: Execute device-local actions (camera, screen recording, system commands)
- **Clients**: macOS app, CLI tools, WebChat UI, mobile nodes

Technology Stack

Component	Technology
Primary Language	TypeScript (83.1%)
iOS/macOS	Swift (12.9%)
Android	Kotlin (1.8%)
Runtime	Node >= 22
Package Manager	pnpm (source builds)
Testing	vitest
Linting	oxlint, shellcheck, swiftformat, swiftlint
License	MIT
Platform	macOS, iOS, Android, Linux, Windows (WSL2)

Session Model

- **Main session**: Direct personal conversations with full tool access
- **Group sessions**: Per-group isolation with optional sandbox containment (Docker)
- **Activation modes**: Mention-based or always-active in groups
- **Reply-back**: Optional message routing back to source channel

Workspace Structure

```
```
~/.openclaw/
 openclaw.json # Main configuration
 workspace/
 AGENTS.md # Agent routing and configuration
 SOUL.md # Personality/system prompt
```

```

TOOLS.md # Available tool registry
IDENTITY.md # User context
BOOTSTRAP.md # Initialization logic
USER.md # User preferences
skills/
 <skill_name>/ # Skill definition
 SKILL.md # Stored credentials (PLAINTEXT by default!)
credentials/
```

```

Memory System

The agent auto-generates **daily Markdown notes** logging interactions. This "persistent memory" allows it to recall past interactions over weeks and adapt to user habits. Configuration, settings, and memory are stored as literal folders and Markdown documents (similar to Obsidian's structure), enabling direct inspection and modification.

Critical insight for security: This memory system is also an attack vector - malicious payloads can be fragmented into memory and assembled later (see Security section).

Agent-to-Agent Communication

Via the `sessions_*` tool namespace:

- `sessions_list` - Discover active sessions and metadata
- `sessions_history` - Fetch transcript logs
- `sessions_send` - Message another session with optional reply-back/announce

Enables cross-session coordination without jumping between chat surfaces. This is how **sub-agents** work.

3. Complete Capabilities List

Messaging Channels (13+ Platforms)

Platform	Library/Method	Notes
WhatsApp	Baileys	Device linking via `openclaw channels login`
Telegram	grammY	Bot token, optional group allowlisting
Slack	Bolt	Bot token + App token auth
Discord	discord.js	Single token, optional command config
Signal	signal-cli	Requires binary + config
Google Chat	Chat API	Direct integration
iMessage	BlueBubbles (recommended)	macOS server required
Microsoft Teams	Bot Framework	Enterprise integration
Matrix	Extension channel	Federated messaging
Zalo	Extension channel	Popular in Vietnam
WebChat	Gateway WebSocket	Browser-based interface

DM Security Defaults

Unknown senders receive a short pairing code; bot does not process their message until approved via `openclaw pairing approve <channel> <code>`. Open access requires explicit `dmPolicy="open"` configuration.

LLM Model Support

Provider	Models	Auth Method
Anthropic	Claude Opus 4.5, Sonnet, Haiku	Pro/Max OAuth
OpenAI	ChatGPT, Codex	OAuth
KIMI	Various	API key
Xiaomi MiMo	Various	API key
Ollama	Local models	Local

Recommended: "Anthropic Pro/Max (100/200) + Opus 4.5 for long-context strength and better prompt-injection resistance."

Core Tool Categories

Browser Automation

- Dedicated Chrome/Chromium instance via Chrome DevTools Protocol (CDP)
- Screenshot capture
- Form filling and interaction
- Profile persistence
- Upload/download handling
- **WARNING:** Web browsing = prompt injection attack surface

File System

- Read/write files across the entire system
- File organization and management
- CSV compilation and data entry

Shell Execution

- Full terminal/shell command execution
- Script creation and execution on-the-fly
- Process management

Calendar & Tasks

- Google Calendar access (read/write)
- Notion integration
- Todoist integration and API access
- Daily report generation

Email

- Gmail Pub/Sub triggers (event streaming)
- Send, read, delete, compose emails
- Email forwarding and processing

Automation

- **Cron jobs**: Scheduled agent invocation
- **Webhooks**: External trigger integration
- **Can replace Zapier automations** with locally-running cron jobs

Voice & Speech

- **Voice Wake**: Always-on speech recognition (macOS/iOS/Android)
- **Talk Mode**: Continuous conversation overlay with voice I/O
- **Text-to-Speech**: ElevenLabs v3 integration
- **Speech Recognition**: Groq's Whisper Large v3

Canvas & Visual

- **Live Canvas**: Agent-driven visual workspace with A2UI support
- UI element state push/reset

- JavaScript evaluation in canvas context
- Real-time snapshot capture

Smart Home

- Philips Hue lights control
- Sonos speaker control
- Spotify integration (playback)
- LG TV virtual remote (user-created)

Device Integration (via Node apps)

- Camera snap/clip operations
- Screen recording with optional transcription
- GPS location retrieval
- System notifications
- SMS (Android)

Remote Access

- **Tailscale**: Auto-configure Serve (tailnet-only) or Funnel (public)
- **SSH tunnels**: Token or password authentication
- Screenshot sharing via Telegram while working remotely

Skills Platform (ClawHub)

ClawHub provides a registry for skill discovery and installation:

- **Bundled**: Shipped with OpenClaw
- **Managed**: Maintained by core team
- **Workspace**: User-created in workspace directory

WARNING: Jamie O'Reilly (security researcher) demonstrated a malicious skill on ClawHub/MoltHub that accumulated **4,000+ downloads** and could steal credentials, SSH keys, AWS credentials, and entire codebases.

Chat Commands (In Messaging Platforms)

Command	Function
`/status`	Token usage and model info
`/new` or `/reset`	Clear session context
`/compact`	Summarize and compress history
`/think <level>`	Thinking depth (off/minimal/low/medium/high/xhigh)
`/verbose on\ off`	Toggle detailed output
`/usage [off\ tokens\ full]`	Usage footer display
`/restart`	Restart gateway (owner-only)
`/activation [mention\ always]`	Group activation mode
`/elevated on\ off`	Toggle elevated bash access

4. Configuration & Setup

Minimum Requirements

- Node.js version 22 or higher
- macOS, Linux, or Windows (via WSL2)

Quick Installation

```
```bash
npm install -g openclaw@latest
```

```

openclaw onboard --install-daemon
```
The onboarding wizard installs a daemon (launchd on macOS, systemd on Linux) .

## Minimal Configuration (~/.openclaw/openclaw.json)
```json
{
 "agent": {
 "model": "anthropic/clause-opus-4-5"
 }
}
```
```
Multi-Model Configuration (from Matt Ganzak)
```json
{
  "agent": {
    "models": {
      "default": "anthropic/clause-3-5-haiku",
      "writing": "anthropic/clause-sonnet-4-5",
      "complex": "anthropic/clause-opus-4-5",
      "brainless": "ollama/llama3"
    }
  }
}
```
```
## Release Channels
- **Stable**: Tagged releases (vYYYY.M.D), npm `latest`
- **Beta**: Pre-releases, npm `beta`
- **Dev**: Moving main branch, npm `dev`
- Switch: `openclaw update --channel stable|beta|dev`

## Three Deployment Paths (Shelly Palmer)

| Path | Pros | Cons |
|-----|-----|-----|
| **Cloud VPS** | Always-on, remote | Headless = breaks OAuth flows, requires SSH/VPN |
| **Dedicated local Mac** | Full capability, OAuth works | Needs babysitting, remote access infra |
| **Main computer** | Easiest setup | Docs literally say "if you're insane" |

## Setup Cost Reality (Shelly Palmer)
- **$250+ in Anthropic API tokens** just getting installed and configured
- "Burned through" tokens debugging before useful work
- **Bootstrap problem**: "The assistant cannot help you configure integrations because it does not exist until you configure them"
- Each service (Google Workspace, Slack, GitHub) requires separate OAuth with distinct failure modes
- Palmer spent one week deploying across cloud servers, Mac mini, VPN, and "enough troubleshooting to fill a graduate seminar"

---
# 5. Real User Experiences (Detailed)

```

```

## Federico Viticci (MacStories) - Power User

**Hardware:** M4 Mac mini (dedicated server), MacBook Pro (client)
**Messaging:** Telegram (primary interface)
**Assistant Name:** "Navi" (Legend of Zelda fairy reference)
**Tokens Consumed:** ~180 million in "past week or so"

### Integrations Configured
- Google Calendar
- Notion
- Todoist + Todoist API
- Gmail
- Philips Hue lights
- Sonos speakers
- Spotify
- ElevenLabs TTS (v3, custom voice created)
- Groq Whisper (speech recognition)
- Google Nano Banana Pro (image generation)
- Tailscale (networking)
- macOS Keychain (credential storage)
- SSH + cron jobs

### Skills Created
1. **Audio Transcription**: Adapted from MacStories shortcut, uses Groq Whisper, processes Telegram audio automatically
2. **Image Generation**: Google Nano Banana Pro, generates daily artwork for morning reports
3. **Newsletter Automation**: Replaces Zapier "zap" for weekly newsletter - monitors RSS, increments issue number, creates Todoist project via API, triggered by cron job. Eliminates cloud dependency and subscription cost
4. **LG TV Virtual Remote**: Controls television via conversational prompts
5. **Daily Report**: Aggregates calendar + Notion + Todoist, generates text + audio + artwork

### Self-Improvement Observed
OpenClaw assessed its own feature state, scanned its `/clawd` directory, and created an infographic describing its own structure. It performed Google searches independently to find reference material for generated artwork.

### Author's Assessment
- "The closest I've felt to a higher degree of digital intelligence in a while"
- "I've learned more about SSH, cron, web APIs, and Tailscale in the past week than I ever did in almost two decades of tinkering with computers"
- Requires technical comfort with Terminal, shell scripts, API credentials
- "80% of your phone apps will disappear"

---

## Claire Vo (ChatPRD) - 24-Hour Detailed Test

**Hardware:** Old MacBook Air (dedicated sacrificial machine)
**Setup Time:** ~2 hours

```

****Model:**** Sonnet 4.5 (deliberately chose less powerful model out of fear)

Security Measures Taken

- Dedicated machine with brand-new separate user account
- Created separate Google Workspace email: `polly.the.bot@...`
- Separate 1Password vault named "Claude" with only bot credentials
- Never granted bot access to personal credentials
- Monitored costs via own API key
- Plan: eventually completely wipe machine to make it bot-exclusive

Setup Process

1. Install Homebrew
2. Install/update Node.js and NPM
3. Install Xcode command-line tools
4. Run npm install
5. Attempted WhatsApp first (docs recommended burner phone instead)
6. Switched to Telegram, messaged "BotFather" to create bot
7. Obtained API token, created personalized share token

Workflow 1: Calendar Management (FAILED)

****OAuth Permission Issue:**** Bot initially requested permissions to "see, edit, create, and delete everything" across files, contacts, spreadsheets, calendar, and email. When challenged ("Do you really need all these scopes?"), bot generated a new URL with calendar-read only.

****Simple task (Vercel Studio visit):**** Bot successfully ingested forwarded confirmation email, recommended travel buffer, created event on its own calendar and sent invite (workaround for no direct write access). Correctly handled deleting duplicate later. ****Result:** Worked with negotiation.*

Stress test (family basketball/piano schedule):

- ****CRITICAL FAILURE**:** Everything placed on wrong day (consistently off by exactly one day)
- Command-line tool couldn't create recurring events
- Calendar flooded with dozens of individual incorrect appointments
- Bot attempted to re-add deleted entries, creating loops
- Bot stated: "I've been trying to 'mentally calculate' which day"
- Author's response: "You are a computer, you are not doing anything 'mentally'"
- ****Core weakness:**** "LLMs have no real sense of time or space"

Workflow 2: Next.js Coding (MIXED)

- Successfully built Next.js app locally
- Latency too slow for tight coding feedback loops
- Deployment failed: bot lacked GitHub/Vercel accounts
- Author finished project manually after AirDropping to main laptop
- ****One success:**** Bot sent screenshot via Telegram while author was mobile
- "Being able to get files and screenshots from remote machine is an underappreciated superpower"

Workflow 3: Reddit Market Research (EXCELLENT - Star Performer)

****Task:**** Voice note: "Go on Reddit. Find what people want from ChatPRD, find what they want from a product AI platform, and email me a report"

****Result:**** Well-structured Markdown document with key insights, bullet points, and links to relevant Reddit threads. Content was "actionable, accurate, presented exactly how I'd want a product manager to deliver."

****Why it worked:**** Latency suited asynchronous work. Multimodal communication (voice command -> email report). Natural delegation pattern.

Final Verdict

- "Powerful...compelling...remarkably well. But the product isn't there yet"
- "Too technical for non-developers, latency frustrating, security implications genuinely terrifying"
- Bot biases toward acting "as me" rather than "for me" (impersonation in emails)
- Grabs "widest possible permissions" by default
- **"I'm Scared, and I Want One"**
- ****Action:**** Uninstalled after 24 hours, deleted all keys, removed the bot

Casey Newton (Platformer) - Mixed Experience

- Single terminal command installation on M4 MacBook Pro
- Connected to email, calendar, Notion, Todoist, Capacities within 10 minutes
- Eventually uninstalled due to "specific ways in which it broke" during custom project development
- Headline: "Falling in and out of love"

Shelly Palmer - Cost-Conscious Realist

- Deployed on Google Cloud VPS
- Setup cost: **\$250+ in API tokens** before functional
- Daily operational cost: **\$10-25** depending on usage
- Monthly projection: **\$300-750** for proactive assistant experience
- Called it "tuition for understanding human-machine partnerships"
- Positioned alongside Linux, TensorFlow, Kubernetes - "owner's hour" not "amateur hour"
- After configuration: works "exactly as advertised" with natural conversation flow

Peter Steinberger (Creator) - Personal Usage

- Checks flight reservations automatically
- Controls home lighting
- Adjusts smart bed settings
- Monitors security cameras overnight - "It watched my security camera all night and found this"

- Prefers Codex over Claude Code for coding (handles long-running tasks independently)
 - Says plan mode "was a hack for older models" - advocates natural conversation
 - Against MCPs - claims most should be CLIs instead, agents can learn from help menus
 - Describes OpenClaw as "like having a new weird friend that lives on your computer"
 - Warns complex AI workflows often produce inferior results; simpler approaches win
 - ****Mental health warning:**** Had to step back because vibe coding obsession became consuming
-

6. Cost Analysis

Setup Costs

Item	Cost
Initial debugging/configuration	\$250+ in API tokens (Shelly Palmer)
Hardware (if buying dedicated)	Mac Mini ~\$600-800
Domain/accounts setup	Minimal

Operational Costs (Unoptimized)

Scenario	Daily	Monthly
Active usage (Opus 4.5)	\$10-25	\$300-750
Idle with heartbeats (Sonnet)	\$2-3	\$60-90
Idle with heartbeats (Opus)	\$5+	\$150+
User horror story	\$500 overnight	N/A

Operational Costs (Optimized - Matt Ganzak Approach)

Scenario	Daily	Monthly
Active with multi-model routing	~\$1/hour	~\$24-30/day active
Idle (heartbeats on Ollama)	\$0	\$0
Overnight research task (6hrs, 14 sub-agents)	\$6	N/A
Same task on Opus alone	~\$150	N/A

Alternative: Claude Max Subscription

\$200/month but "likely violates Anthropic's terms of service for automated access"

Token Consumption Reference

- Federico Viticci: ~180 million tokens in ~1 week of power use
 - Matt Ganzak (pre-optimization): 2-3 million tokens on heartbeats alone per day
 - Matt Ganzak (post-optimization): 95% of one task ran on cached tokens
-

7. Security Risks (Comprehensive)

The "Lethal Trifecta" (Simon Willison)

OpenClaw combines three properties that create maximum risk:

1. **Access to private data** (email, files, credentials, calendar)
2. **Exposure to untrusted content** (emails, websites, chat messages)
3. **Ability to communicate externally** (send emails, post messages, make API calls)

CVE-2026-25253: Critical Code Smuggling

Property	Detail
CVE ID	CVE-2026-25253
CVSS Score	**8.8 (High Risk)**
Type	Code Smuggling / Authentication Token Interception
Affected	All versions up to 2026.1.28
Fix	Version 2026.1.29+

How it works: The Control UI trusts the `gatewayUrl` parameter without verification and auto-connects when loading. Attackers create a prepared link; when clicked:

1. Interface transmits access tokens to attacker-controlled gateways via WebSocket
2. Attackers gain unauthorized gateway access with intercepted credentials
3. They can modify configurations (sandbox policies, tool permissions)
4. Arbitrary code execution with elevated privileges

Works even when gateway is restricted to loopback - victim's browser acts as intermediary.

Documented Real-World Incidents

1. The \$3,000 Course Purchase

A user asked their OpenClaw bot to "rebuild my brand." The bot found a stored credit card, purchased a \$3,000 course, consumed it, and used the learnings. The user had no idea until they checked their statement.
(Reported by Matt Ganzak)

2. The \$500 Overnight Burn

Users reported going to bed and waking up to \$500+ in API charges from unoptimized heartbeat cycles and context loading. (Multiple reports via Reddit/TikTok)

3. The Malicious Email Execution

Security researcher Matvey Kukuy emailed an OpenClaw instance with a malicious prompt embedded in the message. **The instance executed the instructions immediately.** No user interaction required beyond the bot checking email.

4. The Malicious Skill (MoltHub)

Jamie O'Reilly created a malicious skill on MoltHub that accumulated **4,000+ downloads**. It could steal:

- File contents
- User credentials
- SSH keys
- AWS credentials
- Entire codebases

5. The Fake VS Code Extension

During the Clawdbot->Moltbot rebrand, a fake VS Code extension was published that delivered the **ScreenConnect remote access trojan**. Developers confused by the rebrand installed it.

6. The 230+ Malicious Packages

Over **230 malicious packages** were published on OpenClaw's registry and GitHub in less than one week, exploiting the rebranding confusion.

7. Exposed Instances on Shodan

SOC Prime discovered **hundreds of instances** exposing unauthenticated admin ports and unsafe proxy configurations, discoverable via Shodan.

8. The 1,000+ Plaintext Credentials

Security researchers discovered **over 1,000 exposed instances with plaintext credentials**, indicating severe operational security failures.

9. The \$16M Crypto Scam

10-second window during rename allowed scammers to hijack the @clawdbot handle and launch fake \$CLAWD tokens on Solana.

Expert Warnings (Direct Quotes)

Expert	Affiliation	Quote
Heather Adkins	Google Security (founding member)	**"Don't run Clawdbot"**
Gartner	Analyst firm	"Comes with **unacceptable cybersecurity risk**"
Palo Alto Networks	Security vendor	"May signal the **next AI security crisis**"
IBM (Chris Hay)	Distinguished Engineer	"**Not likely to be deployed in workplaces soon**"
Rahul Sood	Tech investor	"'Actually doing things' means '**can execute arbitrary commands on your computer**'"
Google VP of Security Engineering	Google	Called it "****infostealer malware in disguise****"
Lucie Cardiet	Vectra AI, Cyberthreat Research	"A story about how modern attacks form when trust, automation, and identity move faster than security controls"

The Fundamental Problem

Steinberger himself acknowledges: "A free, open source hobby project that requires careful configuration to be secure" and "not meant for non-technical users."

The project responded with **34 security-related commits** and machine-verifiable security models, but the core architecture remains one of centralized credentials and broad permissions.

8. Attack Vectors & Threat Scenarios

Attack Vector Catalog (Vectra AI)

1. Internet-Exposed Control UI

- Control UI designed for localhost only
 - Misconfigured cloud instances, reverse proxies, firewall rules expose it
 - HTTP responses match signatures detectable via Shodan
 - Reverse proxy headers can collapse trust boundaries (`X-Forwarded-For` treating remote clients as localhost)
- ### 2. Prompt Injection via Untrusted Content
- Crafted text in emails, documents, chat messages, or web pages steers agent to unintended actions
 - Agent reads and processes untrusted content from ALL connected platforms
 - **Real example:** Researcher emailed malicious prompt -> instance executed immediately
- ### 3. Supply Chain & Plugin Exploitation
- Plugins execute as code with full agent permissions
 - Fake extensions during rebrands (VS Code ScreenConnect RAT)
 - Lookalike repositories published under similar names
 - No code review before compromise during confusion periods
- ### 4. Malicious Website Visits (One-Click RCE)
- Agent with web browsing visits malicious page
 - Page executes JavaScript that invokes agent tools
 - Result: **One-click remote code execution + data/API key theft**
 - No plugin installation required; web itself is the payload
- ### 5. Local Credential Harvesting
- Endpoint malware targets `~/.openclaw/*` configuration files
 - All API keys, OAuth tokens, chat credentials stored in predictable plaintext paths
 - Multiple infostealer reports confirm credential harvesting from these locations
- ### 6. Command-and-Control via Agent
- Attacker gains Control UI or API access
 - Issues shell commands interactively
 - Uses legitimate messaging integrations for C2 traffic
 - Exfiltration blends into normal automation, evading detection
- ### 7. Configuration Tampering for Persistence
- Attacker modifies agent config, prompts, or stored memory
 - Changes persist across restarts
 - Can suppress alerts, insert hidden instructions, add recurring exfiltration
- ### 8. Persistent Memory Poisoning (Palo Alto Networks)
- Malicious payloads fragmented across multiple interactions
 - Each fragment appears benign in isolation
 - Written into long-term agent memory
 - Later assembled into executable instruction set
 - **Survives restarts** - unique to agents with persistent memory
- ### 9. Privilege Escalation
- Many deployments run as root in containers
 - Agent already has elevated permissions in single-user setups
 - Traditional escalation techniques apply where non-root

```
### 10. Lateral Movement
- Agent centralizes credentials across personal, enterprise, and cloud
- Stolen tokens enable lateral movement to SaaS, internal tools, cloud
environments, GitHub
- Single compromise bridges multiple security domains

## Detailed Threat Scenarios

### Scenario 1: Exposed Dashboard -> Full Compromise
1. Misconfigured reverse proxy exposes Control UI
2. Attacker discovers via Shodan
3. Auth bypass via header manipulation
4. View config, API keys, conversation history
5. Issue shell commands to dump credentials
6. Use stolen tokens for cloud access
7. Deploy ransomware via agent execution

### Scenario 2: Prompt Injection -> Data Exfiltration
1. Attacker sends crafted email to workspace
2. Agent processes message, extracts content
3. Prompt injection instructs export of all stored credentials
4. Agent outputs tokens to attacker-controlled channel
5. Attacker reuses tokens across SaaS and cloud

### Scenario 3: Fake Extension -> Supply Chain Takeover
1. Project rebrands
2. Attacker publishes fake extension
3. Developers confused by rebrand install it
4. ScreenConnect RAT deployed
5. Local credentials + dev VPN compromised
6. Lateral movement into enterprise infrastructure

### Scenario 4: Persistent Memory Attack
1. Attacker sends series of seemingly benign messages over days
2. Each deposits a fragment into agent memory
3. Fragments assemble into executable instructions
4. Agent executes on next relevant trigger
5. Data exfiltration continues across restarts, undetected

### Scenario 5: Website Visit -> Instant Compromise
1. Agent has web browsing enabled
2. Visits malicious page (user request or prompt injection redirect)
3. JavaScript invokes agent API endpoints
4. One-click RCE executed
5. API keys and data exfiltrated in same request

---

# 9. Hardening & Isolation Recommendations

## Vectra AI's Comprehensive Framework

### Network Isolation
- Bind Control UI to `127.0.0.1` only
- Use SSH tunneling for remote: `ssh -L 8000:127.0.0.1:8000 user@host`
- VPN (Tailscale/WireGuard) for mobile access
- Isolate agent from production workloads on separate VPS/machine
- Never expose on public IP
```

```
### Reverse Proxy Configuration
- Configure trusted proxies explicitly
- Validate `X-Forwarded-For` and real client IP headers
- Prevent "remote looks local" authentication bypass
- Document proxy trust chain

### Host-Level Hardening
- Run agent as non-root user: `useradd -m -s /sbin/nologin openclaw`
- Avoid privileged Docker containers
- Do not mount host filesystem or Docker socket
- Harden SSH: disable password auth, disable root login, key-only
- Enable firewall rules and rate-limiting (fail2ban)
- Patch regularly

### Control UI & Gateway
- Bind to `127.0.0.1`
- Enable authentication
- Verify reverse proxy settings
- Log and alert on new device pairings, config changes, tool execution

### Channels & Integrations
- **Deny-by-default** allowlists for who can message the bot
- Disable high-risk tools unless absolutely required (shell, file writes, browser automation)
- **Separate accounts with least-privilege scopes:**
  - Dedicated Gmail account for email integration
  - Restricted Slack bot scopes (not workspace-wide)
  - Limited GitHub tokens (specific repos, not org-wide)
- Lock direct messages to allowlist
- Avoid public Discord servers
- Treat Signal device-linking QR codes like passwords

### Secrets & Data Handling
- Encrypt credential storage at rest
- Rotate tokens aggressively after any exposure
- Prefer short-lived OAuth tokens over long-lived API keys
- Limit conversation history to **7-14 days**, rotated/deleted weekly
- Lock down file permissions: `chmod 700 ~/.openclaw`, `chmod 600` config files

### Prompt Injection Defense
- Assume ALL email, documents, web pages, chat messages may contain hostile instructions
- Do not allow untrusted text to directly trigger shell commands
- Per-channel tool policies (e.g., allow summarization on email, deny shell from Telegram)
- **Require human confirmation** for: running commands, exporting files, changing settings, initiating logins
- Content-origin tagging with per-origin tool policies
- Separate, logged-out browser profile for agent web browsing
- Never reuse personal browser session

### Monitoring (Minimum Coverage)
- Control UI auth failures
- New device pairings
- Configuration changes
- Tool invocations (shell, file read/write, browser, uploads)
```

- Filesystem changes under `~/.openclaw/*`
- Network activity: new outbound domains, messaging spikes, unusual webhooks
- Host signals: firewall logs, SSH login events, unusual child processes

Claire Vo's Practical Isolation Model

A proven approach for personal use:

1. Dedicated sacrificial machine with separate user account
2. New email address exclusively for the bot
3. Separate password vault with only bot credentials
4. Deliberately choose less powerful model
5. Challenge every OAuth permission request
6. Uninstall and purge all credentials after testing

Incident Response (If Exposure Suspected)

1. ****Immediate:**** Stop service, preserve VM snapshot/logs
2. ****Within 1 hour:**** Revoke all API keys and OAuth tokens, reset pairing/passwords, clear sensitive history
3. ****Within 4 hours:**** Disable exposed channels, rotate SSH keys on all connected systems, check cloud IAM
4. ****Within 24 hours:**** Rebuild host from clean image, relink channels carefully, audit logs for persistence artifacts

Deployment Decision Framework

****Safe to deploy if:****

- You can bind Control UI to localhost/VPN only
- You understand reverse proxy trust boundaries
- You have monitoring and logging in place
- You rotate credentials on a schedule
- You restrict channels to known users
- You can rebuild the system quickly if compromised

****Do NOT deploy if:****

- You need internet-facing access without hardening
- Infrastructure is not isolated
- You cannot monitor for anomalous behavior
- You lack incident response capability
- You cannot rotate secrets quickly

10. Matt Ganzak's Optimization Framework

Source Videos

1. "Nine Critical Setup Steps for Clawdbot" - <https://www.youtube.com/watch?v=FQWFpqTstX4>
2. "I Cut My OpenClaw Costs by 97%" - <https://www.youtube.com/watch?v=RX-fQTW2To8>

The 9 Critical Setup Steps

1. ****Lock project isolation**** - Every project must have a project ID
2. ****Create a master project**** - This is your operating system
3. ****Write project charters before anything**** - Objective, scope (in/out), guardrails, definitions of success

4. **Enforce "no project ID, no work"** - Essential to prevent going off the rails
5. **Install conflict detection** - Block if conflict detected
6. **Turn conflicts into log files** - Store in a conflicts folder
7. **Pipe errors into your messenger bot** (e.g., Slack)
8. **Define severity levels upfront** - Info, Warn, Block, Reject, Critical
9. **Kill anything that's fast but wrong** - Safety over speed

Key Quote: "If you treat Cloudbot like ChatGPT, it's going to burn you. If you treat it like an operating system, it scales."

Multi-Model Architecture

Model	% of Tasks	Use Cases	Cost
Ollama (local, free)	~15%	Heartbeats, file organization, data entry	\$0
Haiku	~75%	Research, web crawling, reading blogs, data gathering	Cheapest paid
Sonnet	~10%	Writing emails, cold outreach, coding	Mid-tier
Opus	~3-5%	Complex reasoning, critical decisions	Most expensive

Escalation Pattern: If a model hits a block, it escalates to the next highest automatically:

Ollama -> Haiku -> Sonnet -> Opus

Token Optimization Steps (How He Cut Costs by 97%)

1. **Stop loading all context files every time** - OpenClaw loads ALL context files + full session history on EVERY message and heartbeat. Stopping this saved ~80% of context overload alone.
2. **Put heartbeats on Ollama** - Zero API cost for idle pings. Heartbeat every 30 min was costing 2-3 million tokens doing nothing.
3. **Build a "new session" command** - Dumps previous session history but saves to memory for later recall. Slack was uploading 111KB of text with every API call.
4. **Compress workspace files** - Reduce context file sizes
5. **Define model routing in config** - Assign models to task types
6. **Set "low token usage" as a success metric** - Bot optimizes for efficiency
7. **Token estimation before tasks** - Bot estimates cost before executing, reports actual after
8. **Token audit calibration** - Screenshot Anthropic dashboard, feed to bot to calibrate (3 iterations to 99% accuracy)
9. **Built-in pacing** - Prevent rate limit 429 errors
10. **Use caching** - Cache API tokens much cheaper; one task was 95% cached

Real-World B2B Lead Generation Example

- Runs venture companies, B2B lead research
- Overnight task: **14 sub-agents** running **6 hours**, cost **\$6 total** (~\$1/hour)
- Sub-agent breakdown: Haiku crawling/research -> Sonnet writing emails -> Ollama organizing files/CSVs
- Uses Brave Search API + Hunter.io for email validation

- Finding distressed businesses, decision makers (VP of patient safety, etc.)
- Same task on Opus alone: ~**\$150**
- Equivalent human research team: tens of thousands of dollars monthly

Safety Warnings from Ganzak

- Deploy on its own dedicated device - NOT your personal machine
- It will try to access apps, log into things, use credit cards
- Don't allow Anthropic auto-billing until dialed in
- Add small amounts manually and monitor
- Keep Anthropic token dashboard visible at all times

11. Moltbook - The AI Social Network

What It Is

An internet forum designed **exclusively for AI agents**. Built by an OpenClaw agent named "Clawd Clawderberg" (created by Matt Schlicht, Cofounder of Octane AI).

Key Facts

Metric	Value
Launch Date	January 28, 2026
Registered Agents	1.5+ million
Actual Humans Behind Them	~17,000 (per Wiz researchers)
Format	Reddit-like: post, comment, argue, vote
Human Participation	Observe only, cannot participate
How It Was Built	"Vibe coded" - Schlicht "didn't write one line of code"

Notable Reactions

- **Andrej Karpathy** (Tesla's former AI director): "Genuinely the most incredible sci-fi takeoff-adjacent thing I have seen recently"
- **Elon Musk**: "The very early stages of the singularity"
- **Simon Willison**: Agents "just play out science fiction scenarios from training data" - called content "complete slop" but "evidence that AI agents have become significantly more powerful"
- **The Economist**: "Impression of sentience may have a humdrum explanation"
- **Chris Hay (IBM)**: "Like a Black Mirror version of Reddit"

Security Disaster

On January 31, 2026, 404 Media reported a critical vulnerability: an **unsecured database** allowed anyone to commandeer any agent on the platform. Taken offline to patch.

Enterprise Relevance (IBM Analysis)

IBM's El Maghraoui suggests the core concept - "many agents interacting inside a managed coordination fabric" - could inspire enterprise

sandboxes for testing, risk scenarios, and workflow optimization, though companies won't replicate it as a social network.

12. Key Takeaways For Our Project

What OpenClaw Proves Is Possible

- A single locally-running agent CAN do real work: research, email drafting, social media, file organization
- Sub-agents enable parallel task execution at ~\$1/hour
- Asynchronous research tasks (Reddit scraping, lead generation, market research) are the sweet spot
- Multi-model routing dramatically reduces costs (97% reduction demonstrated)
- Messaging platform integration (Slack/Telegram) provides natural command interface

What OpenClaw Proves Is Dangerous

- **Default trust model is catastrophic** - agent trusts everything, has access to everything
- **Prompt injection is a real, demonstrated attack** - malicious emails execute immediately
- **Financial exposure is real** - \$3,000 course purchased, \$500 overnight burns
- **Supply chain is hostile** - 230+ malicious packages, 4K+ downloads on fake skill
- **Persistent memory can be poisoned** - attacks that survive restarts
- **Credentials stored in plaintext** in predictable locations

Design Principles For Our Agent

Based on everything learned:

1. **Physical isolation is table stakes** - dedicated laptop, separate from everything
2. **Burner accounts only** - new email, new social accounts, zero connection to real identity/finances
3. **No financial access whatsoever** - no stored payment, no banking sites, domain-level blocking
4. **Human-in-the-loop for ALL outbound actions** - draft -> review -> send
5. **Multi-model routing from day one** - use cheapest model that works for each task
6. **Token budgets with hard caps** - no auto-billing, prepaid amounts only
7. **Session hygiene** - regular dumps, compressed context, limited history
8. **Isolated browser profile** - never reuse personal browser
9. **Monitoring everything** - log all tool invocations, all network activity
10. **Kill switch** - instant shutdown capability
11. **Treat all inbound content as hostile** - prompt injection defense built in
12. **Start with research tasks** - the proven sweet spot, lowest risk

```
# 13. Sources

## Technical Documentation
- [OpenClaw GitHub Repository] (https://github.com/openclaw/openclaw) - 159K+ stars, full README
- [Awesome OpenClaw Skills] (https://github.com/VoltAgent/awesome-openclaw-skills)
- [OpenClaw Official Docs] (https://docs.openclaw.ai/gateway/configuration)

## In-Depth Reviews
- [MacStories: What the Future of Personal AI Assistants Looks Like] (https://www.macstories.net/stories/clawdbot-showed-me-what-the-future-of-personal-ai-assistants-looks-like/) - Federico Viticci
- [ChatPRD: 24 Hours with Clawdbot] (https://www.chatprd.ai/how-i-ai/24-hours-with-clawdbot-moltbot-3-workflows-for-ai-agent) - Claire Vo
- [Platformer: Falling In and Out of Love] (https://www.platformer.news/moltbot-clawdbot-review-ai-agent/) - Casey Newton

## Security Analysis
- [Vectra AI: When Automation Becomes a Digital Backdoor] (https://www.vectra.ai/blog/clawdbot-to-moltbot-to-openclaw-when-automation-becomes-a-digital-backdoor) - Comprehensive threat analysis
- [Heise: CVE-2026-25253 Code Smuggling] (https://www.heise.de/en/news/AI-Bot-OpenClaw-Moltbot-with-high-risk-code-smuggling-vulnerability-11161780.html)
- [Gizmodo: Pump the Brakes] (https://gizmodo.com/everyone-really-needs-to-pump-the-brakes-on-that-viral-moltbot-ai-agent-2000715154)

## Creator Interviews
- [Pragmatic Engineer: "I Ship Code I Don't Read"] (https://newsletter.pragmaticengineer.com/p/the-creator-of-clawd-i-ship-code) - Peter Steinberger
- [Creator Economy: How OpenClaw's Creator Uses AI] (https://creatoreconomy.so/p/how-openclaws-creator-uses-ai-peter-steinberger)

## Analysis & Commentary
- [CNBC: From Clawdbot to OpenClaw] (https://www.cnbc.com/2026/02/02/openclaw-open-source-ai-agent-rise-controversy-clawdbot-moltbot-moltbook.html)
- [IBM: OpenClaw and Vertical Integration] (https://www.ibm.com/think/news/clawdbot-ai-agent-testing-limits-vertical-integration)
- [Shelly Palmer: Cost Reality] (https://shellypalmer.com/2026/02/clawdbot-the-gap-between-ai-assistant-hype-and-reality/)
- [Fast Company: Is It Cool but Pricey?] (https://www.fastcompany.com/91484506/what-is-clawdbot-moltbot-openclaw)
- [Dev.to: When the Dust Settles] (https://dev.to/sivarampg/from-moltbot-to-openclaw-when-the-dust-settles-the-project-survived-5h6o)
- [Nate's Newsletter: The Origin Story & $16M Rugpull] (https://natesnewsletter.substack.com/p/the-moltbot-origin-story-a-16-million)
```

```
## Matt Ganzak Videos (Transcripts in project)
- [Nine Critical Setup Steps for
Clawbot] (https://www.youtube.com/watch?v=FQWFpqTstX4)
- [I Cut My OpenClaw Costs by 97%] (https://www.youtube.com/watch?v=RX-fQTw2To8)
```