

# Security Architecture — claude-config-template

## Honest Risk Assessment

This system gives an AI agent significant access to your machine. That is inherently dangerous.

However, the risk profile varies dramatically by deployment model:

	Claude Code + this config	OpenClaw (self-hosted)
<b>Runs as</b>	Interactive CLI, user present	24/7 daemon, unattended
<b>Human in the loop</b>	Yes — user approves destructive ops	No — fully autonomous
<b>Attack surface</b>	Local machine, user's permissions	Server + all connected channels (Slack, Teams, WhatsApp)
<b>Credential exposure</b>	.mcp.json on local disk	API keys in running service, accessible from network
<b>Blast radius of compromise</b>	One user's files	All connected channels + server + all users who interact
<b>Sandbox</b>	bubblewrap (Linux namespaces)	None (runs as service user)
<b>Permission gating</b>	Explicit allow-list per tool	No tool-level permission model
<b>Prompt injection risk</b>	Low — input is the user typing	High — any message in Slack/Teams can be crafted input
<b>Audit trail</b>	Git history + session context	Application logs only
<b>Recovery from bad action</b>	git revert, files are local	Messages sent, API calls made — not reversible

**Bottom line:** Claude Code with this config is a powerful tool with real risks, but the human-in-the-loop model + sandbox + permission gating + git audit trail make it **categorically safer** than any unattended AI agent. OpenClaw's always-on, multi-channel exposure is a fundamentally different (and higher) risk class.

## Residual risks even with this setup

- Claude runs with **your user permissions** — it can read/write anything you can
- Write(\*) is pre-approved — a hallucinating agent could overwrite important files (git history is your safety net)
- MCP servers extend reach to external services (GitHub, Twitter) — a confused agent could post or commit unwanted content
- docker is in the allow list — container operations have their own blast radius
- The config repo itself is a high-value target — anyone with push access controls Claude's behavior

This is not “safe.” It is “managed risk with recovery options.”

---

## Sandboxing (Claude Code Built-in)

Layer	Mechanism	What it does
<b>Process sandbox</b>	bubblewrap (bwrap)	Linux namespace isolation — Bash commands run in a restricted container with limited filesystem/network access
<b>IPC proxy</b>	socat	Proxies communication between Claude Code and sandboxed processes — prevents direct system access
<b>Platform</b>	Linux namespaces + seccomp	Kernel-level syscall filtering. No privilege escalation from within the sandbox

**Required packages:** socat + bubblewrap (installed by setup/install-base.sh and setup/bootstrap-fedora.sh)

Without these, **sandboxed tool calls fail silently** — Claude Code falls back to unsandboxed execution or auto-denies.

---

## Permission Model

### How it works

Claude Code uses a grant-based permission system in `settings.json`:

```
permissions.allow = [ "ToolName(glob_pattern)", ... ]
```

### What's pre-approved (shipped in template)

Category	Allowed	Examples
<b>Read-only tools</b>	All	Read(*), Glob(*), Grep(*), WebSearch, WebFetch(*)
<b>File modification</b>	All	Write(*), Edit(*)
<b>Safe Bash commands</b>	40+ patterns	git, npm, node, python3, docker, gh, curl, mkdir, cp, mv
<b>Orchestration</b>	Skill only	Skill(orchestration)

### What's NOT pre-approved (prompts user)

Command	Why blocked
rm / rm -rf	File deletion — must confirm
sudo	Privilege escalation — must confirm
chmod (partially)	chmod is allowed but chown is not
Arbitrary Bash	Commands not matching any allow pattern trigger a prompt

### Subagent permissions

Background subagents **cannot prompt the user**. If a tool isn't in the allow list, the call is **auto-denied silently**. This is the #1 cause of mysterious subagent failures. Fix: add the minimal matching pattern to `permissions.allow`.

---

## Credential & Secret Management

Layer	Mechanism
<b>API keys</b>	Environment variables loaded from <code>secrets.env</code> (never committed)
<b>MCP credentials</b>	Stored in <code>.mcp.json</code> (per-user, not in repo)
<b>Git authentication</b>	<code>git-credential-mcp</code> helper reads PAT from <code>.mcp.json</code> — single source of truth, no embedded tokens
<b>VPS secrets</b>	<code>secrets.env.template</code> tracked, <code>secrets.env.gitignore</code>

### `.gitignore` protection

```
.env, .env.* , *.key, *.pem, secrets/, credentials/
```

Both repos enforce this. The template ships with a conservative `gitignore`.

---

## Trust Boundaries

Layer	Role	Key controls
<b>USER</b> (interactive)	Approver	Can approve/deny any tool call in real-time
<b>CLAUDE CODE</b> (main process)	Executor	Runs with user's filesystem permissions. Bash sandboxed via bubblewrap. Permission model gates destructive ops.
<b>SUBAGENTS</b> (background)	Workers	Same sandbox. CANNOT prompt user – auto-denied if not in allow list. Isolated context windows.
<b>MCP SERVERS</b> (external tools)	Integrations	Each server has its own credentials in .mcp.json.
<b>VPS</b> (remote, headless)	Infrastructure	Scoped access. No cross-server credential sharing. secrets.env loaded at bootstrap only. SSH tunnel for services. Let's Encrypt TLS.

## What's NOT Covered (Gaps)

Gap	Risk	Mitigation
<b>No dedicated security doc</b> **Write(*) is fully open**	Security rules scattered across 10+ files  Claude can overwrite any file the user owns	This document consolidates them  Sandbox + git history provide recovery
<b>No network egress filtering</b>	Sandboxed processes may still make outbound connections	MCP servers are the primary external access; Bash curl is allowed
<b>No file integrity monitoring</b>	Modified files not detected between sessions	Git diff at session start (hook) shows changes
<b>Subagent blast radius</b>	A subagent with write(*) can modify any file	Git-based audit trail; session hooks commit state
<b>MCP token scope</b>	GitHub PAT with repo scope = full repo access	Use fine-grained PATs where possible

## Platform-Specific Security

### WSL (Windows Subsystem for Linux)

Feature	Detail
Sandbox deps Windows Defender	sudo apt install socat bubblewrap Exclude WSL paths to avoid scanning overhead (see wsl-environment.md)
Line endings Performance boundary	core.autocrlf input prevents CRLF injection /mnt/c/ is slow — keeping files in /home/ also reduces Windows-side exposure

### Native Linux (Fedora, etc.)

Feature	Detail
Sandbox deps	<code>dnf install socat bubblewrap</code> (via <code>bootstrap-fedora.sh</code> ) bubblewrap uses kernel namespaces natively
No additional hardening needed	

---

## Audit Trail

Mechanism	What it captures
<b>Git history</b>	Every file change, every session, every machine
<b>session-context.md</b>	What was done, when, by which session
<b>Session hooks</b>	Auto-commit + push at session end; auto-pull at session start
<b>Registry</b>	Which projects exist, where, what status

The combination of git + session hooks means **every action is version-controlled and recoverable**.

---

## Zone Identifiers / Windows Security Tools

**Not currently documented in the config template.**

On Windows/WSL, two security mechanisms may apply to downloaded files:

Tool	What it does
<b>NTFS Alternate Data Streams (ADS)</b>	Windows attaches a :Zone.Identifier stream to downloaded files marking their origin (Internet, Intranet, etc.) — the “Mark of the Web” (MOTW)
<b>Unblock-File (PowerShell)</b>	Removes the Zone.Identifier ADS, telling Windows the file is trusted

These are **Windows-side** mechanisms. Within WSL’s Linux filesystem (/home/), NTFS ADS don’t exist. They only matter for files accessed via /mnt/c/ or when executing scripts from the Windows filesystem.

**If you were using tools that left Zone Identifiers during WSL setup, those were likely Windows-side script downloads.** They don’t affect the Claude Code sandbox (which is pure Linux bubblewrap).

---

Generated 2026-02-25. Source: *claude-config-template* + *claude-config analysis*.