



DETERMINISTIC AUTHORIZATION FOR CODING AGENTS

Your AI has
too much access.

THE NEW NORMAL

Plan together, then let it run unattended

Plan with the agent, then let Claude Code run autonomously

With no boundaries, it takes **any** action its context suggests - any file, any command, any host.



● Re

● Wr

● Ba

● Re

● We

...app



THE STAKES

Broken access control is still #1.

The screenshot shows the OWASP Top 10:2025 RC1 page for A01:2025 Broken Access Control. The left sidebar contains a navigation menu with items like 'Welcome Page', 'Current Release', 'Release Candidate', 'Introduction', 'About OWASP', 'What are Application Security Risks?', 'Establishing a Modern Application Security Program', 'Top 10:2025 List', and 'A01 Broken Access Control'. The main content area is titled 'A01:2025 Broken Access Control' and includes a 'Background' section and a 'Score table'.

Background.

Maintaining its position at #1 in the Top Ten, 100% of the applications tested were found to have some form of broken access control. Notable CWEs included are *CWE-200: Exposure of Sensitive Information to an Unauthorized Actor*, *CWE-201: Exposure of Sensitive Information Through Sent Data*, *CWE-918 Server-Side Request Forgery (SSRF)*, and *CWE-352: Cross-Site Request Forgery (CSRF)*. This category has the highest number of occurrences in the contributed data, and second highest number of related CVEs.

Score table.

CWEs Mapped	Max Incidence Rate	Avg Incidence Rate	Max Coverage	Avg Coverage	Avg Weighted Exploit	Avg Weighted Impact	Total Occurrences
40	20.15%	3.74%	100.00%	42.93%	7.04	3.84	1,839,701

Ranked #1

100% of a access co

Now hand

owasp.org/Top10

CLAUDE CODE · AUTO-PERMISSIONS

false-negative rate on dangerous actions
analysis, which is upfront about

"Claude can make mistakes that allow harmful content to be generated. It is recommended to only use in isolated environments."

– Anthropic documentation

THE MISMATCH

Authorization demands **deterministic**



✗ Probabilistic judgment ✗

A model weighs context and returns its best guess at what's safe. Run it again and the answer can change.

"probably fine" → 83% of the time



A po
ques

cheo

SECTION 01 / 06

A primer on authoriza

MODELS OF ACCESS CONTROL

Four ways to model access

ACL

Access lists

A list, per object, of who can do what. Accurate, but it doesn't scale.

RBAC

Roles

Permissions grouped into roles. Simple but does not scale

ABAC

Attrib

Rules
resour
attribu
hard to

GOOGLE ZANZIBAR · REBAC

Google Zanzibar models access as relationships.

- 🌐 Globally distributed, > 10M+ checks per second.
- 🔗 AuthZ system behind Google Docs, YouTube, Maps, Photos & more
- 🔑 To the developer, it's just an API:
Is this actor allowed to perform this action on this resource?

Abstract

Determining digital object access per presents of Zanzibar, access control and configuration access control Google, including YouTube, ordering of use amid change Zanzibar scaling of authorization by billions of users, accuracy of less than 99.999%

1 Introduction

Many online services confirm that

REBAC · THE VOCABULARY

Three building blocks, one tuple.

 **User**

A natural person or service account.

`user:10957135`

 **Object**

Any non-user entity, with a type prefixed id.

`doc:123 · folder:plans`

A RELATION
TUPLE

`document:123 # owner @ user:3`

REBAC · THE SHAPE OF IT

Chain the tuples and you get a **graph**

- 🔗 Objects relate to other objects — and to other objects' relations.
- 🔗 Chain enough tuples and you've built a directed graph — a DAG.
- 🔗 A permission is just reachability in that graph.

somedocum

REBAC · ANSWERING A CHECK

A check is a **walk through the graph**

THE QUESTION

```
check(  
  somedocument,  
  "reader",  
  jill  
)
```

somedocument

STAYS FAST: Parallel subproblems · Cache everything · Dedupe requests

WHY RBAC ISN'T ENOUGH

A role can't say **path, command, or**




`role: developer`

One coarse bucket. All-or-nothing.

WHAT AN AGENT AC

 Write only un

 Run only go

 Fetch only g.

THE OPEN-SOURCE ZANZIBAR

Introducing SpiceDB

- 📄 Google published the Zanzibar whitepaper — the model, not the code.
- 🔗 SpiceDB is its open-source implementation — the engine you actually run.
- 🏢 In production at companies like OpenAI and Zoom, with contributors from the Fortune 500.



A high
authori

Can

🔗 git

SECTION 02 / 06

Introducing SpiceBox

WHAT IT IS



SpiceDB (the authz engine)
+ Sandbox (OS containment)

Fine-grained
coding

\$ br

Apach

THE CORE IDEA

Permission checks are still deterministic

01 

You decide

Describe what the agent may do, in plain English.

*"Read & write Go under src/.
Run go test. Fetch github.com."*

02 

SpiceBox translates

AI maps it to structured SpiceBox permissions — you review and approve.

```
file rw src/**/*.go
bash go test web github.com
```

FROM DESCRIPTION TO ENFORCEMENT

Six steps, one session.

THE TASK

```
the ask
```

```
"Add retry logic to the Go API client."
```

One ordinary coding task — watch it pass through all six steps, scoped the whole way.

```
service-api · src/
```

01  Describe

02  Translate

03  Review

04  Enforce

05  Run

06  Cleanup

QUICK START

One command to launch

```
$ spicebox claude
```

- 1 Describe what Claude can do
- 2 Review the translated permissions
- 3 Approve, and launch Claude Code

Spice

Approve

[x] file r

[x] bash

[x] web ·

[] agents

> Approve

SECTION 03 / 06

Three layers of enforce

OVERVIEW • DEFAULT-DENY

Every tool call runs the same **gauntlet**.


Tool call



APPLICATION

Hook server

Checked against SpiceDB before it runs.



OPERATING SYSTEM

OS sandbox

Filesystem limits
kernel.

Only explicitly granted permissions pass.

LAYER 1 APPLICATION

The hook server checks **before** the c

- 🛡️ **Default-deny.** Only what you granted passes. Everything else is blocked.
- 🔗 **Cross-tool.** Read, Glob and Grep all check the same file permission.
- 🏠 **Out-of-process.** The model can't reach or rewrite the rules.

ONE

Re

LAYER 2 OPERATING SYSTEM

The kernel is the **backstop**.




- 📦 A **sandbox-exec** profile restricts the process at the OS level.
- 📁 No write permission granted? The working directory is **read-only**.
- 🚫 Bash tricks like **echo > file** are blocked by the kernel anyway.



```
# ho  
$ ec  
sand  
Open
```

LAYER 3 NETWORK

Egress goes through an **allowlist**.

-  All outbound traffic is routed through a **localhost proxy**.
-  Connections are allowed only to **permitted domains**.
-  A `curl evil.com` via bash dies at the proxy, not the model.

```
$ curl -v https://evil.com
✓ 200 OK
$ curl -v https://evil.com
× proxy error: connection refused
```


WHAT YOU CAN SCOPE


Five permission types

 **Blanket** Category-wide access

 **File** Path prefix + suffix filter

 **Bash** Command prefix matching

 **Web** Domain-based URL filtering

 **MCP** Tool, prefix, or parameter value

APPROVAL CONTROL

Granted isn't the same as **silent**.

 [auto]

Proceeds silently

Matching tool calls run with no interruption. Reads default here.



Pro

The
defa

SECTION 04 / 06

The threat model

ENFORCEMENT THAT HOLDS UNDER ATTACK

Deterministic checks survive a **compromis**



Assume the model is compromised.

Prompt injection can hijack the agent's reasoning. It doesn't matter — every check runs out-of-process and deterministically, so a hijacked model still can't exceed its grant.



Prompt injecti
denied regardle



Bash tool wor
kernel blocks th



Data exfiltrati
host isn't on the



Privilege esca
rules live where



Alternative to
Grep can't reach

BEING HONEST ABOUT THE EDGES

What SpiceBox does **not** protect against



Bash is Turing-complete

Broad bash prefixes can be gamed. The OS sandbox — not the hook — is the backstop.



Transpilation

Natural language processing
you apply



No suffix exclusions yet

"Read all files except .env" approximates to blanket read — with a warning.



Sandboxing

Other OS-level

SECTION 05 / 0

Live de

Scoping paths · network isolation · a crypto

DEMO 01 · PATH SCOPING

An agent that can't write outside its lane

POLICY

```
write *.go  
under src/ only
```

Anything outside `src/` is denied — twice.

```
● ● ● claude – under sp  
  
Write(.. /app.config)  
× hook: file_resource  
  
Bash(echo cfg > ../ap  
× sandbox: file-write  
  
Write(src/handler.go)  
✓ allowed – write com
```

DEMO 02 • NETWORK ISOLATION

What network isolation looks like

✓ On the allowlist

```
$ curl github.com/authzed/...  
✓ 200 OK
```

✗ E

```
$ c  
x p
```

Routed through the proxy whether it's a tool call

DEMO 03 · MID-SESSION GRANT

The agent asks for more. **A human gra**

●●● /spicebox-grant

Claude requests:

"Also let me run npm test"

translated · HMAC-signed

+ bash npm test [confirm]

Grant this permission? >

✓ relationship written – live now





SECTION 06 / 06


Beyond Claude Code

WHY NOT JUST SETTINGS.JSON?

Tool-level rules leave **gaps**.

 No cross-tool policies
→ SpiceBox: shared file checks

 No parameter constraints
→ SpiceBox: MCP value filtering

 Static for the session
→ SpiceBox: cryptographic grants

 Allc
→ S

 No
→ S

 Eve
inter

WHAT TO TAKE AWAY

Four things to remember

- 01 Model judgment is not enforcement. Keep the two separate.
- 02 Deterministic, out-of-process checks hold even under pressure.
- 03 ReBAC scopes by path, command, and domain — RBAC does not.
- 04 Let the AI draft the policy — but a human approves it.

NEXT STEPS

Try it, read it, talk to us



SpiceBox

github.com/authzed/spicebox



Install

```
brew install --cask authzed/tap/spicebox
```



S
g



C
S

AuthZed

Thank you!

Sohan Maheshwar

[linkedin.com/in/sohanmaheshwar](https://www.linkedin.com/in/sohanmaheshwar)