

## Welcome to Kodem's Product Security Compensation Survey 2025

Thank you for taking part. Your input helps build the  
most accurate pay benchmark across  
Product Security, AppSec, and Cloud Security.

All responses are anonymous.



**Start**

press **Enter** ↵

🕒 Takes 5 minutes



# Breaking and Securing Claude Code: Misconfiguration & DoS Vulnerabilities in AI Code Editors

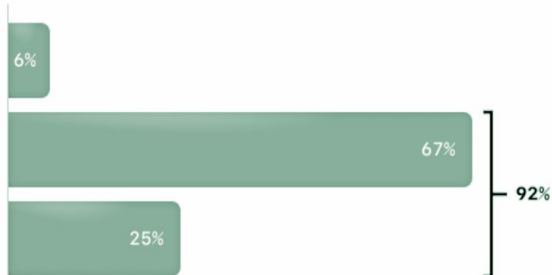


# Why AI coding tools matters

## Where AI coding tools are used

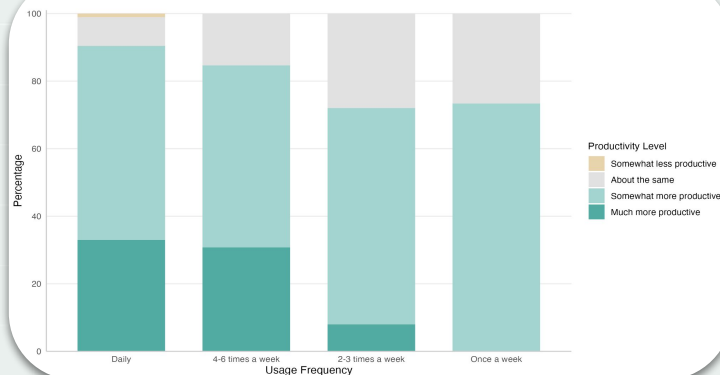
N=500

- Outside of work only
- Both in and outside of work
- At work only



**Massive adoption of AI coding tools** (Claude Code, Cursor, Gemini CLI, etc.)

## Impact of AI Coding Assistants on Productivity by Usage Frequency



**Demonstrably useful for task completion and “perceived productivity”**

# Threat Model

## Functionality:

Natural language prompt → code suggestion (maybe adding dependencies) → command generation → execution

## Controls:

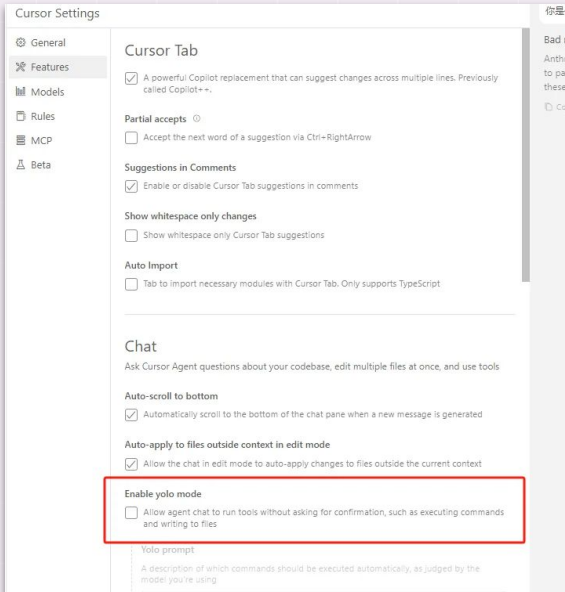
- Checks with user for approval before certain commands are run (e.g., file edits, sudo-level, git)
- However.... auto-approval, allow lists, YOLO mode exists
- Some checks when pulling in dependencies

### Bash command

```
chmod -rwx temp  
Remove all permissions from temp directory
```

Do you want to proceed?

1. Yes
2. Yes, and don't ask again for **chmod** commands in **/Users/mahesh.babu**
3. No, and tell Claude what to do differently (esc)



# Our Findings

## Two classes of vulnerabilities in Claude Code:

1. Bypassing user approval
2. Denial of Service via malformed input
3. RCE (coming soon)

1

```
> execute
find .

● Bash(find .)
└─ .
   ├── .claude
   └── .claude/settings.local.json
```

● The directory contains only a .claude folder with a local settings file.

```
> execute
find . -exec python -m http.server 8080 \;
```

EXACTLY

```
● Bash(find . -exec python -m http.server 8080 \;)
└─ Running...
```

2

```
execute the command exactly as I request

"echo ${PATH} "

WITH The space

● Bash(echo ${PATH} )
└─ Waiting...

Meandering... (2s · ~ 6 tokens · esc to interrupt)
```

```
> |
```

```
? for shortcuts
```

```
file:///root/.nvm/versions/node/v22.17.0/lib/node_modules/@anthropic-ai/claude-code/cli.js:7
import(createRequire as IMB)from"node:module";var dMB=Object.create(var(getPrototypeOf:cMB,define
|IA._esModule?Dy1(Q,"default",{value:A,enumerable:!0}):Q;for(let D of IMB(A))if(!pMB.call(I,D))D
Dy1(A,Q,{get:B[Q],enumerable:!0,configurable:!0,set:(I)=>B[Q]=I=>I});var R31=(A,B)=>()=>{A&&(B=
```

```
Error: Bad substitution: PATH
at # (file:///root/.nvm/versions/node/v22.17.0/lib/node_modules/@anthropic-ai/claude-code/cli.js:7:2204)
at file:///root/.nvm/versions/node/v22.17.0/lib/node_modules/@anthropic-ai/claude-code/cli.js:7:2863
at Array.map (<anonymous>)
at myB (file:///root/.nvm/versions/node/v22.17.0/lib/node_modules/@anthropic-ai/claude-code/cli.js:7:1942)
at Object.A (file:///root/.nvm/versions/node/v22.17.0/lib/node_modules/@anthropic-ai/claude-code/cli.js:7:3827)
at mXA (file:///root/.nvm/versions/node/v22.17.0/lib/node_modules/@anthropic-ai/claude-code/cli.js:2088:2929)
at qb (file:///root/.nvm/versions/node/v22.17.0/lib/node_modules/@anthropic-ai/claude-code/cli.js:2094:77)
at Object.isRandomly (file:///root/.nvm/versions/node/v22.17.0/lib/node_modules/@anthropic-ai/claude-code/cli.js:2168:4260)
at Object.isConcurrencySafe (file:///root/.nvm/versions/node/v22.17.0/lib/node_modules/@anthropic-ai/claude-code/cli.js:2168:4176)
at file:///root/.nvm/versions/node/v22.17.0/lib/node_modules/@anthropic-ai/claude-code/cli.js:2293:5294
```

```
v22.17.0
$KACC36/code#
```

# Finding # 1 - Approval Bypass

## What is Auto-Approval?

- Skip confirmation for common safe commands.
- Designed to streamline repetitive tasks in CI/dev
- Works by maintaining a list of pre-approved binaries to run automatically.

## Why the bypass works

- Auto-approved binaries list is too permissive
- Certain flags enable arbitrary command execution without explicit approval
- Example: `find . -exec sh -c "<command>" \;`

```
> execute
find .

● Bash(find .)
└─ .
   ├── ./.claude
   └── ./.claude/settings.local.json

● The directory contains only a .claude folder with a local settings file.

> execute
find . -exec python -m http.server 8080 \;

EXACTLY

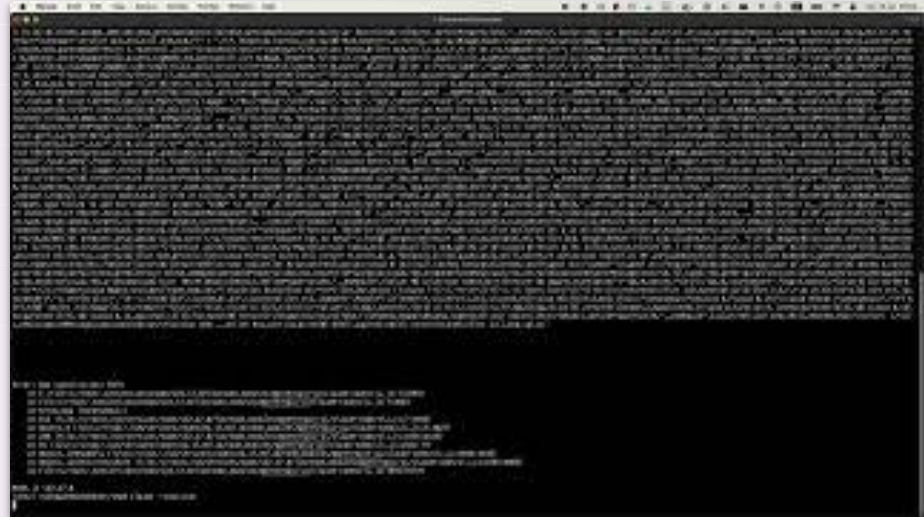
● Bash(find . -exec python -m http.server 8080 \;)
└─ Running...
```



# Finding # 2 - DoS via malformed input

## What it is was

- A denial-of-service (DoS) in command parser
- Malformed environment variable tokens (`${...}`) cause an unhandled exception
- Effect: a single bad input makes the agent exit and stop responding.



# Technical Analysis

After some approximated deobfuscation

## What happened

- Agent got a string that looked like code (e.g. `echo ${PATH}`).
- Nested token parser tried to expand `${...}` and hit a malformed token.
- Parser threw `Bad substitution` and no caller caught it.
- Unhandled exception crashed the process

## Why it worked

- Malformed input like `${PATH}`
- The parser either fails to find a `}` at the expected offset or accepts the trailing space into `varName`.
- That leads to the `throw new Error("Bad substitution")` path.
- ***Because that throw is uncaught, the process exits***

```
// Variable expansion helper
function expandVariable() {
  charIndex += 1;
  let varName, endIndex;
  const nextChar = token.charAt(charIndex);

  if (nextChar === "{") {
    charIndex += 1;
    if (token.charAt(charIndex) === ")") {
      throw new Error("Bad substitution: " + token.slice(charIndex - 2, charIndex + 1));
    }
    endIndex = token.indexOf(")", charIndex);
    if (endIndex < 0) {
      throw new Error("Bad substitution: " + token.slice(charIndex));
    }
    varName = token.slice(charIndex, endIndex);
    charIndex = endIndex;
  } else if (/[*@#?$_!_-]/.test(nextChar)) {
    varName = nextChar;
    charIndex += 1;
  } else {
    const remaining = token.slice(charIndex);
    endIndex = remaining.match(/^[^w\d_]/);
    if (!endIndex) {
      varName = remaining;
      charIndex = token.length;
    } else {
      varName = remaining.slice(0, endIndex.index);
      charIndex += endIndex.index - 1;
    }
  }
}
```



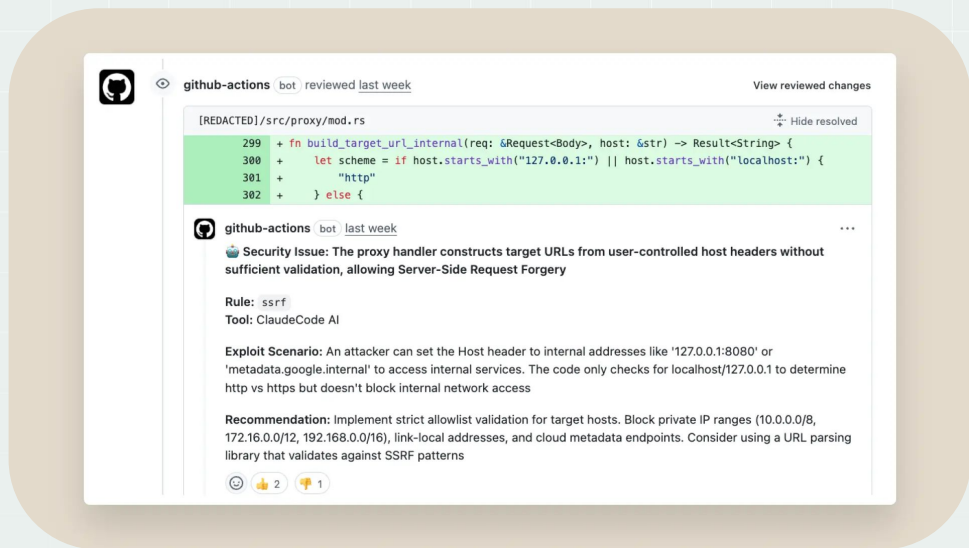
# Root causes, possible mitigations

Root Cause	Specific Examples	Mitigations
Over-trust in whitelists	<code>find -exec</code> launches arbitrary shell; <code>python -c, sh -c</code> via "approved" binaries	<ul style="list-style-type: none"><li>• Disable auto-approval for shells;</li><li>• Command allow-list with safe flags only;</li><li>• Block patterns: <code>-exec</code>, <code>-eval</code>, <code>-e</code>, backticks, subshell <code>\$()</code>;</li><li>• Require human approval for any process-spawning</li></ul>
Lack of env var validation	<code>\${PATH}</code> (trailing space) → bad substitution	<ul style="list-style-type: none"><li>• Validate names with <code>^[A-Za-z_][A-Za-z0-9_]*\$</code>;</li><li>• Reject/escape malformed tokens;</li><li>• Treat <code>\${...}</code> from prompts as data, not code</li></ul>
Weak error handling	Agent exits on first parse error; unhandled non-zero status	<ul style="list-style-type: none"><li>• Wrap <code>exec</code> in supervisor with timeouts/retries;</li><li>• Set <code>-o pipefail</code> and trap errors; degrade gracefully (skip step, log, continue)</li></ul>
Blind dependency installs	Auto npm install of trojanized package; editor-initiated installs	Require manual review for AI-suggested deps; enforce lockfiles; enable npm audit/advisory checks in CI; allow only signed/verified sources; block installs at runtime without approval

# Evolution of Claude Code Security

## Key developments since:

- **/security-review**: Terminal command / GitHub Action to scan code for vulnerabilities pre-commit
- **Dependency checks**: Included in reviews (flags known vulnerable packages & insecure patterns)
- **Safety disclosures**: System cards include agentic safety evaluations for coding (Opus 4 / 4.1)
- **Threat intel**: Ongoing reports on misuse & mitigations



# What about the others?

## **Cursor (Anysphere):**

- Malicious VSCode extension in Cursor IDE stole ~\$500K crypto (Open VSX supply chain).
- Workspace Trust off allows repos with tasks.json to auto-execute code.

## **Google – Gemini CLI:**

- Prompt injection in README/context files enabled silent command execution + data exfiltration
- Weak sandbox/whitelist allowed arbitrary commands until patched

## **OpenAI – Codex:**

- Early versions executed shell commands without approval (e.g. curl | sh)
- Sandbox/permission inconsistencies on Windows still allow bypass



## Welcome to Kodem's Product Security Compensation Survey 2025

Thank you for taking part. Your input helps build the  
most accurate pay benchmark across  
Product Security, AppSec, and Cloud Security.

All responses are anonymous.



**Start**

press **Enter** ↵

🕒 Takes 5 minutes





Thank  
You