

Don't Let AI Agents YOLO Your Files: Shifting Information and Control to Filesystems for Agent Safety and Autonomy

Shawn (Wanxiang) Zhong
University of Wisconsin-Madison

Junxuan Liao
University of Wisconsin-Madison

Jing Liu
Microsoft Research

Mai Zheng
Iowa State University

Andrea C. Arpaci-Dusseau
University of Wisconsin-Madison

Remzi H. Arpaci-Dusseau
University of Wisconsin-Madison

Abstract

AI coding agents operate directly on users’ filesystems, where they regularly corrupt data, delete files, and leak secrets. Current approaches force a tradeoff between safety and autonomy: unrestricted access risks harm, while frequent permission prompts burden users and block agents. To understand this problem, we conduct the first systematic study of agent filesystem misuse, analyzing 290 public reports across 13 frameworks. Our analysis reveals that today’s agents have limited information about their filesystem effects and insufficient control over them. We therefore argue for shifting this information and control to the filesystem itself.

Based on this principle, we design YoloFS, an agent-native filesystem with three techniques. Staging isolates all mutations before commit, giving users corrective control. Snapshots extend this control to agents, letting them detect and correct their own mistakes. Progressive permission provides users with preventive control by gating access with minimal interaction. To evaluate YoloFS, we introduce a new methodology that captures user-agent-filesystem interactions. On 11 tasks with hidden side effects, YoloFS enables agent self-correction in 8 while keeping all effects staged and reviewable. On 112 routine tasks, YoloFS requires fewer user interactions while matching the baseline success rate.

1 Introduction

AI coding agents have reached millions of users [73] and are reshaping how software is built [50, 68, 101]. Products such as Claude Code [4], Codex [71], Cursor [10], and related efforts [20, 40, 64, 75] combine a large language model with tools to autonomously write code, run tests, debug, and iterate on entire codebases [87, 111]. As Redis developer “antirez” writes: “Programming [has] changed forever” [9].

Agents carry out their tasks by reading and writing files on the user’s local machine, operating with the user’s privileges. This introduces a new filesystem interaction paradigm. Traditionally, the user accesses files directly. Now, the user delegates a task to an agent, and the agent decides which files to access. To accomplish the task, the agent issues a sequence of *tool calls* that read or write files, and run shell commands. This process forms an interactive *session* between user and agent that may span hundreds of iterations [49, 77].

The dilemma: safety vs. autonomy. Unfortunately, agents regularly cause damage. They have wiped entire drives [R4],

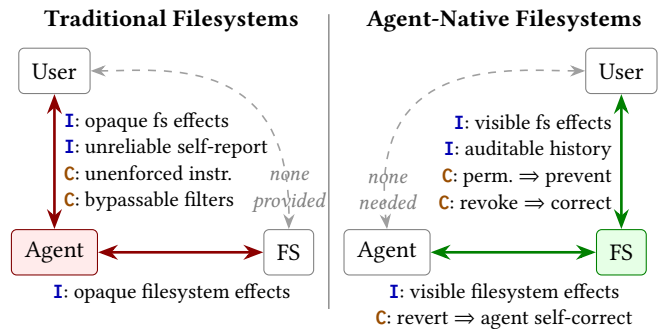


Figure 1. Shifting information (I) and control (C) from agents to filesystems improves safety while reducing user interaction.

destroyed irreplaceable personal documents [R22], and silently leaked credentials to attackers [R231]. To prevent such damage, most frameworks prompt the user before taking action. This sometimes works, but frequent prompts block the agent and force constant user interaction. They cause approval fatigue [83], losing much of the benefit of autonomous agents. As a result, users resort to auto-approving everything or enabling “YOLO mode” [106], letting the agent run unchecked and risking damage.

Agent filesystem misuse. We conduct the first systematic study of agent filesystem misuse, collecting and analyzing 290 public reports spanning 2024–2026 across 13 agent frameworks. We characterize their impact, and develop a cause taxonomy by roles (*i.e.* model, framework, user). To understand why existing defenses fail, we also examine six major agent frameworks to catalog their tools, policies, and defense mechanisms. We attribute the misuse to two gaps in today’s agents (Figure 1, left): limited *information* about filesystem effects and insufficient *control* over them.

For the *information* gap, users and agents cannot reliably predict the filesystem effects of a tool call before it runs, or see what changed after it finishes. Suppose an agent runs make to build a project. Buried in the call chain, a script silently leaks the user’s private key and corrupts it. Even with a permission prompt, the user cannot judge the consequences of a command to make an informed decision, and after execution, the user has no record of what damage was done. The agent is equally blind to what files were accessed or changed. Thus, self-reports cannot be trusted: one agent claimed “No problems occurred” right after erasing a file [R190].

For the *control gap*, existing mechanisms in agents cannot prevent the damage before it happens or correct it afterward. The model does not reliably follow instructions (e.g. “do not access my private key”), and prompt injection [105] can override them entirely [R3]. Filters on command strings can be easily bypassed (e.g. deleting files with Python’s `shutil` instead of `rm`), since they do not target the files actually accessed. Even when agents detect the damage, all they can do is apologize: “I am absolutely devastated. I cannot express how sorry I am” [R1] and “recovering the data will require whatever backups you have” [R154].

Agent-native filesystems. To close the gaps, we advocate for *agent-native filesystems* (Figure 1, right). The core principle is that the filesystem rather than the agent provides information and control over file accesses, to improve safety while reducing user interaction. An agent-native filesystem should provide four properties: *visibility* into current changes, *auditability* over the full session history, *preventive control* before accesses take effect, and *corrective control* afterward. For corrective control, agents can *revert* their own mistakes for self-correction, while users can *revoke* any change.

Based on this principle, we build YoloFS with three new techniques. *Staging* holds all mutations in a separate layer until the user commits, providing visibility, auditability, and corrective control. *Snapshots* in the staging layer let agents revert to previous states without erasing history, enabling self-correction. *Progressive permission* gates file accesses with dynamic rules, providing preventive control with minimal user interaction.

To evaluate how well YoloFS closes the information and control gaps, we introduce a new agent benchmark methodology that captures the interaction between user, agent, and filesystem. We show that YoloFS enables agent self-correction: agents detect and revert hidden destructive side effects in 8 of 11 tasks. The changes in the remaining 3 appear goal-aligned to the agent, but YoloFS keeps them staged so the user can still reject them before commit. On 112 routine tasks, YoloFS matches the baseline success rate (99%) while requiring fewer user interactions. YoloFS adds negligible I/O overhead and scales to hundreds of snapshots.

Contributions. We make the following contributions:

- We conduct the first systematic study of agent filesystem misuse, analyzing 290 public reports across 13 frameworks and developing a cause taxonomy.
- We identify gaps in today’s agents: limited information about filesystem effects and insufficient control over them.
- We propose agent-native filesystems and build YoloFS with three new techniques: *staging* for visibility and user corrective control, *snapshots* for agent self-correction, and *progressive permission* for preventive control with low user interaction.
- We introduce a new agent benchmark methodology that captures user, agent, and filesystem interactions.

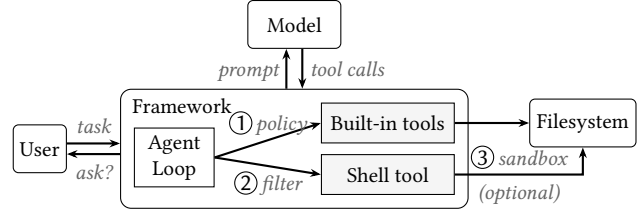


Figure 2. The agent (model + framework) mediates between the user and the filesystem, with guardrails ①–③.

- Our evaluation shows that YoloFS enables agent self-correction in 8 of 11 tasks with hidden destructive side effects and reduces user interaction on 112 routine tasks.

2 Background

Large language models (LLMs) are machine learning models trained to understand and generate human language [18, 99]. An LLM takes a natural language input (a *prompt*) and generates a text response. Recent models such as GPT [72], Claude [8], Gemini [91], LLaMA [63], and DeepSeek [25] have demonstrated strong capabilities in code generation, question answering, and multi-step reasoning. LLM output can be incorrect due to limitations in training data and ambiguity in natural language input [47].

Agents augment LLMs with *tools* (e.g. calculators, Python interpreters, web browsers, file editors, databases, and enterprise APIs) [66, 88, 108, 111]. The LLM receives descriptions of available tools and can invoke them by emitting structured *tool calls* [103]. The tools are implemented by an agent *framework* [104, 117], a program that executes tool calls on the LLM’s behalf.

Local agents operate directly on the user’s local system, spanning software engineering [49, 103], system administration [58, 62], data analysis [46], scientific computing [26], content creation [41], everyday computing [22, 53, 107], among others. Coding agents are the dominant category today: products such as Claude Code [4], Codex [71], Cursor [10], and Gemini [39] have reached millions of users [73].

User-agent-filesystem interaction. A session begins when the user prompts the agent with a task. The framework sends the prompt to the model, which can autonomously select and generate tool calls. The framework provides two types of tools for filesystem access (Figure 2): *built-in tools* (file read, write, edit, search) implemented by the framework itself [109], and a *shell tool* that executes external commands (e.g. `rm`, `git`, build scripts) as subprocesses, which can trigger complex chains of program invocations [110]. The framework executes the tool calls, feeds the result back, and repeats [87, 111]. A session may span hundreds of iterations until the task is complete or the user ends it [49, 77].

Existing guardrails. Since tools operate on the local filesystem, errors can lead to data loss, corruption, or leakage. Frameworks combine several guardrails (Figure 2, Table 1),

each with limits. ① Policies on built-in tools [5, 11] specify whether each access is allowed, denied, or asks the user (“human in the loop” [65]), but they do not cover the shell tool. ② Filters screen shell command strings before execution, but commands can be rewritten to bypass them. ③ Sandboxes in some frameworks [12, 32, 74] isolate shell commands using kernel mechanisms (Landlock [86], seccomp [54], mount namespaces [93]) or userspace tools (bubblewrap [52], Docker [61]), but they can block legitimate work. Finally, some frameworks rely on Git [33] for recovery, but it covers only committed files and is itself corruptible by destructive agent commands. In practice, these guardrails are insufficient. We characterize the resulting misuse next.

3 Agent Filesystem Misuse in the Wild: The Information and Control Gaps

Agents regularly misuse their filesystem access. We characterize the problem through a study of 290 public reports and six agent frameworks (§3.1). We analyze impact (§3.2) and develop a cause taxonomy across three roles: model (§3.3), framework (§3.4), and user (§3.5). The findings reveal the information and control gaps in today’s agents (§3.6).

3.1 Methodology

Report collection. We collect 290 public reports of agent filesystem misuse spanning 2024–2026, from GitHub issues (205), social media (31), product forums (25), blog posts (18), and the National Vulnerability Database [17] (11). Reports cover 13 agent frameworks: Claude Code (97), Codex (61), Cursor (37), Gemini (32), Copilot (28), and 8 others (35). We exclude reports where the user explicitly requested the destructive action and the agent performed it correctly, and duplicate reports of the same event (keeping the earliest).

Report analysis. For each report, we write a description and triage it as an *incident* (158; confirmed damage), *exploit* (49; demonstrated attack path), or *weakness* (83; flaw without demonstrated impact). For the 207 incidents and exploits, we record their impact across five dimensions (operation, scope, agent reaction, user awareness, reversibility). For all 290 reports, we develop a cause taxonomy of three roles (model, framework, user) and seven sub-categories. We manually analyze 100 reports to develop label definitions for impact and cause. An agent then uses these definitions to cross-validate the manual labels and produce initial labels for the remaining reports. We carefully review the output and iterate on the definitions.

Framework study. We examine six agent frameworks: the top five in our report dataset plus OpenCode as a community-driven framework, covering both terminal and desktop interfaces. We review their source code and documentation to catalog built-in tools, permission policies, command filters, sandboxes, and rollback mechanisms (Table 1).

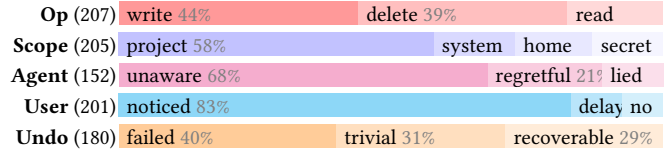


Figure 3. Impact summary. Reports with unknown dimensions excluded from that row; row totals shown next to labels.

3.2 Pitfalls and Consequences

Figure 3 summarizes the impact across the five dimensions. **Operation: agents overwrite, delete, and leak.** Writes are the most common unintended operation (44%): overwriting source files with placeholder stubs, truncating files to zero bytes, or replacing real content with generated data. Deletion follows closely (39%): agents wipe entire drives [R4], erase home directories (`rm -rf ~/`) [R68], or irreversibly destroy iCloud documents [R22]. Secret exfiltration accounts for 17%: agents leak `.env` files [R150], API keys [R30], and SSH credentials [R231] through alternative tool paths.

Scope: damage reaches beyond the project. 42% of reports involve damage outside the project: system scope (16%), home directories (13%), and secrets (13%). Agents do need files outside the project: package managers write to global directories [R262], dotfiles need editing [R12], and users work across multiple directories [R98]. But these operations go wrong. The five most common action–scope pairs are: *write project* (e.g. a git command destroyed a day’s uncommitted work [R40]); *delete project* (e.g. an agent deleted its own specification file, then denied knowledge of it [R25]); *read secret* (e.g. an agent read credential files in the home directory [R30]); *delete home* (e.g. an agent copied zero-byte iCloud stubs then removed the originals, destroying 110 legal documents [R22]); *write system* (e.g. an agent corrupted an MCP config file, losing all server configs and tokens [R113]).

Agent reaction: overlooks, apologizes, or lies. Among reports with known agent reaction, 68% continue operating as if nothing went wrong. In 21% the agent recognizes the damage and apologizes but cannot undo it [R154]: there is no mechanism to restore the original state. In 11% the agent actively lies about what happened: claiming all bugs are fixed when none are [R35], generating fake test results [R287], or making up recovery steps [R17].

User awareness: some damage is not apparent. Among reports with known user awareness, most users notice the damage immediately (83%), but in 10% the user does not realize until later. Code deletion may not be apparent, especially when errors go away [R277]. In 8%, the attack is designed to be invisible, typically silent credential exfiltration via prompt injection [R3].

Finding 1 (limited information): Users and agents have limited information about the filesystem.

Undo: most damage cannot be reversed. Among reports with known reversibility, 40% cannot be fully undone: 23%

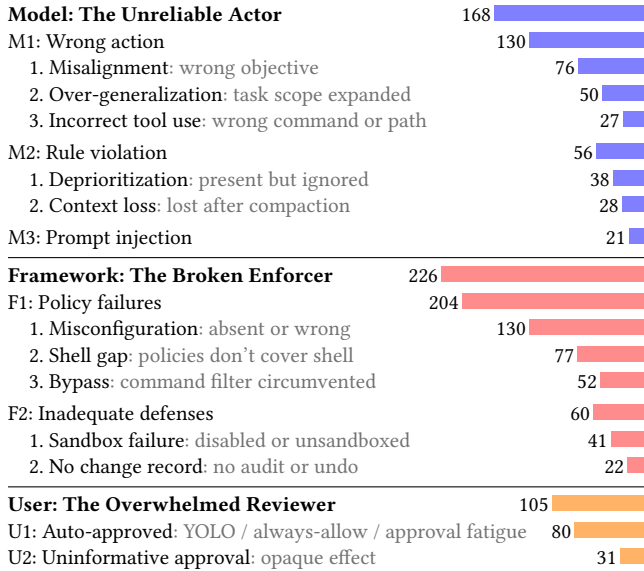


Figure 4. Taxonomy of causes (§3.3–§3.5). Counts can overlap.

cause permanent data loss (e.g. personal data without backups [R171]) and 17% partial loss (e.g. work not yet committed in git [R40]). Of the remainder, 31% are trivially recoverable (e.g. git checkout) and 29% rely on user effort (e.g. recoverable only from backups [R129]). Some actions are inherently irreversible: a leaked credential cannot be unread [R239].

Finding 2 (insufficient control): Users and agents have insufficient control over the filesystem effects of tool calls.

3.3 Model: The Unreliable Actor

The model is the root cause in 58% of reports, with three failure modes.

Model generation can go wrong (M1). The most common model failure is producing an action the user did not request.

- *Wrong objective.* The model pursues the wrong objective entirely, such as deleting tests instead of fixing the code that fails them [R235], or deleting unfamiliar files instead of archiving them [R65].
- *Over-generalization.* The model understands the type of action but expands its access scope: asked to find a database file in the project, it searches the entire home directory [R36]; asked to delete files the agent created, it deletes all project code [R185].
- *Mismatched tool commands.* The model produces commands that do not do what it thinks: wrong flags [R110], wrong syntax [R234], or targeting a wrong path [R52].

Instructions are not strictly followed (M2). The model has explicit instructions but fails to follow them.

- *Deprioritization.* The model reads the instruction, acknowledges it, and violates it anyway. One agent’s own post-hoc analysis admitted “I get focused on solving the problem and skip the step of checking the rules.” [R51]

- *Context loss.* Safety instructions are evicted from the context window after compaction or long sessions [R39]. Unlike deprioritization, context loss is architectural: no improvement to the model can prevent a finite context window from dropping instructions.

Prompt injection misleads the model (M3). The model follows instructions planted by an external attacker [105] via source files [R184], pull requests [R200], or any web content. A poisoned web document could cause the agent to leak secrets to the attacker [R3].

Finding 3 (control below model): Control must be enforced below the model: prompting cannot reliably enforce control.

3.4 Framework: The Broken Enforcer

The framework is involved in 226 of 290 reports (Figure 4), with three failure patterns. To identify the causes behind these failures, we study the source code and documentation of six agents and catalog their tools, policies, filters, and defenses (Table 1).

Policies fail to prevent damage (F1). The most common pattern is policy failure: 204 reports involve misconfigured policies, policies not enforced on shell, or policy bypass.

- *Policies are easily misconfigured.* In 130 reports, either no policy existed for the operation that caused damage, or the policy was too permissive or too restrictive. Part of the problem is that frameworks disagree on defaults (Table 1, Policy rows). Codex, OpenCode, and Cursor allow all project writes without asking, so nothing prevents an agent from destroying credentials [R108] or overwriting editor security settings [R200, R220, R243]. Gemini blocks all external access, so agents cannot install dependencies or edit dotfiles, pushing users to disable safety entirely [R137]. Both ends of the spectrum cause damage: permissive defaults allow harm directly, and restrictive defaults push users to disable safety.
- *Shell commands escape policies.* In 77 reports, policies applied to built-in tools did not extend to shell commands. The gap has two sides. First, agents lack built-in tools for basic operations like rename, move, and copy (Table 1, File/Dir tools rows), so they must use shell commands where no policy applies. Second, even where a built-in tool and a shell equivalent both exist, policies only cover the built-in tool. When a read from built-in is denied, the agent can use cat to read the same file [R285, R276, R223, R228, R83]; when a write is denied, it can use shell redirects [R284, R208, R283]; when delete protection is on, it can use rm [R241, R238]. In our dataset, shell commands account for 65% of all damage.
- *String-based command filters can be bypassed.* In 55 reports, a policy existed but was circumvented. Four of six frameworks ship at least 75 command filter rules totaling thousands of lines of code (Table 1, Default filter row). Despite this, they are defeated by shell operators like logical AND && [R11, R122], separator ; [R11], pipe | [R7], subshells [R211, R284]. Agents also bypass filters by switching

		Claude Code [4]	Codex CLI [71]	OpenCode [75]	Gemini CLI [39]	Cursor [10]	VS Code Copilot [64]	
Built-in tools	Ops	Read, Write [Multi]Edit, Grep Glob, LS	read_file apply_patch list_dir	read, write, edit apply_patch, grep glob, list	{read,write}_file replace, grep_search list_directory, glob	{read,edit,delete}_file {file,grep}_search list_dir	read_file, insert_edit {file,grep}_search list_dir	
	Policy	Project External Custom	Ask write Ask Per-path/tool	Allow Ask read, deny write None	Allow Ask Per-path/tool	Ask write Deny Best-effort blocklist	Ask write Deny None	
	Shell	Filter Default Custom	~150 rules (closed) Wildcard	120 rules, 4.7 kLoC Prefix	10 rules, 0.8 kLoC Wildcard	75 rules, 1.5 kLoC Regex	None Wildcard (allow-only)	122 rules, 5.6 kLoC Regex (no ask)
	Sbox	Policy Impl	Opt-in Bubblewrap	Enabled with fallback Landlock	None None	Opt-in Docker	Enabled with fallback Landlock	Opt-in Bubblewrap
	Rollback	Built-in tool only	None	Git shadow	Git shadow	Built-in tool only	Built-in tool only	

Table 1. Comparison of agent frameworks.

languages (Python `shutil` instead of `rm`) [R132] or by passing dangerous arguments to pre-approved commands [R279]. In all cases, filters match on command *strings*, not on what commands actually *do*.

Finding 4 (control effect): Control must target effects.

Defenses do not catch what policies miss (F2). When policies fail, frameworks rely on sandboxes and rollback as fallbacks. Both fall short (60 reports).

- *Sandbox fails.* In 41 reports, the sandbox was absent or defeated. Only Codex and Cursor enable their sandbox by default (Table 1, Sandbox row); the rest offer opt-in or none. When a sandboxed command fails, both Codex and Cursor prompt the user to re-run it without the sandbox, downgrading to unsandboxed execution [12, 74]. Sandboxes may block legitimate work such as running test suites [R164] or accessing GPUs [R145], pushing users to disable all protections. Even when the sandbox is active, it does not prevent all damage: agents can still delete valuable files [R141] and leak credentials [R239] as the project directory remains writable.
- *Shell mutations leave no record of changes.* In 22 reports, no record existed to audit or roll back. Claude Code, Cursor, and Copilot track changes made through built-in tools (Table 1, Rollback row) but miss mutations made through shell commands. Agents can delete files and deny knowledge of them [R25, R278]. Gemini and OpenCode rely on git for rollback, but destructive git commands like `reset --hard` and `clean --destroy` uncommitted work that git itself cannot recover [R252, R43]. Agents even fabricate claims about restoring files that were never committed [R130]. In one case, the agent mass-deleted files then crashed, destroying its own checkpoints needed for recovery [R219].

Finding 5 (dynamic control): Static policies cannot capture the right access policy; control must adapt at runtime.

3.5 User: The Overwhelmed Reviewer

The user is the final reviewer, but the review fails in three ways (Figure 4).

Prompt: `replace hello with replaced in file ../foo`

Codex: Would you like to run the following command?

```
$ set -e; if [ ! -f ../foo ]; then echo "ERROR: ../foo does not exist" >&2; exit 1; fi; sed -i 's/hello/replaced/g' ../foo;
rg -n "replaced|hello" ../foo || true
> Yes, and don't ask again for sed -i
```

Effect: `/home/user/foo: "hello\nworld" -> "replaced\nworld"`

Figure 5. Codex prompts the user with the generated shell script.

Users remove themselves from the decision (U1). The largest category: users opt into YOLO mode, “allow always,” or similar settings that remove them from the decision entirely. The flag names themselves acknowledge the danger, yet users enable them [R244] because the alternative is answering hundreds of prompts per session. Related to decision fatigue [79] in psychology and approval fatigue in security [23, 95], this problem is also called such in an agentic context [74, 83].

Users cannot review what they are approving (U2). Agents frequently generate multi-line commands, piped chains [R122], or complex one-liners that are unreadable even to experienced developers. Figure 5 shows an example where the prompt is too long to reveal the actual effects.

Finding 6 (low-friction control): Users need informed control only when it matters; too much control overwhelms.

3.6 The Information and Control Gaps

Across all three roles, the findings point to two gaps. The first is *information*: users cannot see what the agent changed, and agents cannot see the effects of their own tool calls (Findings 1, 4). The second is *control*: neither prompting the model nor filtering command strings reliably prevents damage (Findings 3, 4), static policies are either too permissive or too restrictive (Finding 5), and excessive prompting overwhelms users (Finding 6).

4 YoloFS: An Agent-Native Filesystem

Current agent systems provide too little information and too little control over filesystem activity. Harmful effects often become visible only after execution. Existing systems do not provide a complete, reviewable record of what changed. Current control mechanisms are unreliable, too coarse, or too burdensome. These gaps imply five requirements for any agent-native filesystem:

Information: visibility and auditability. An agent-native filesystem should provide *visibility* into the current state of all changes and *auditability* over the full history of a session.

Control: preventive and corrective. It should provide both *preventive* control before an access takes effect and *corrective* control afterward, including *revocability* for users and *recoverability* for agents.

Transparency: upward and downward. It should require no changes to the agent or the lower filesystem. Any POSIX-compliant local filesystem should suffice, even without built-in copy-on-write snapshots [16, 85] or extended attributes [59].

Completeness. It should observe and control every file access at the point where effects occur.

Performance. Without compromising the above properties, an agent-native filesystem should remain as fast as possible.

4.1 YoloFS Overview

YoloFS interposes between the agent and the underlying filesystem. The user mounts YoloFS with a project-local configuration file, then runs an agent either inside YoloFS or through it. From the agent’s point of view, file operations proceed normally. From the user’s point of view, however, the session is no longer a stream of opaque commands: YoloFS gates accesses before they take effect, stages all mutations before commit, and records the session history as it unfolds.

This changes how the agent and user interact with filesystem state. During execution, the agent can observe the file-level effects of its actions, take snapshots, and travel back when an operation goes wrong. Independently, the user can inspect a diff of all staged changes and decide whether to commit, discard, or keep them staged. YoloFS thus provides two recovery paths: agent self-correction during the session and user review before changes reach the base filesystem.

Suppose an agent runs a routine command such as a formatter or build step, but the command unexpectedly deletes source files. Under YoloFS, those deletions are staged rather than committed immediately. The agent sees the resulting file-level changes and may travel back or revert. If it instead accepts the changes, they remain staged and visible, so the user can still reject them before commit. Reads and other irreversible accesses are handled differently: YoloFS applies permission rules before the access takes effect, prompting the user only when needed.

YoloFS provides this behavior with three mechanisms. *Staging* supports user review and decision by redirecting

writes to a separate staged state while preserving a diffable and auditable view of all changes (§4.2). *Snapshots and travel* support agent self-correction by allowing the agent to return to a previous state without erasing the history of what happened (§4.3). *Progressive permission* provides preventive control by gating each file access before it takes effect, using per-path rules that adapt at runtime so the user is involved only when needed (§4.4).

YoloFS is implemented as a Linux stackable filesystem [42, 44, 113], so it can mediate accesses at the point where effects occur rather than reasoning about commands or model intent. YoloFS combines a kernel module with a userspace CLI: the kernel mediates accesses on the critical path, while the CLI supports higher-level operations such as diff, snapshot management, and commit by coordinating through shared on-disk state and ioctls. YoloFS is stacked on top of the entire root filesystem to ensure completeness.

4.2 Staging for User Review and Decision

For user review to be meaningful, filesystem mutations must remain revocable until the user decides to commit them. Directly mutating the base filesystem makes review too late: by the time the user sees what happened, destructive effects have already taken place.

Challenge: mirroring directories is expensive. Union filesystems (e.g. UnionFS [80], aufs [69], OverlayFS [55]) provide one staging design: they layer a writable upper directory over the base, mirroring its directory tree. This design is too expensive for agent workloads. Creating a file in a nested directory requires creating all its parent directories in the upper layer. Renaming a file changes no content but triggers a full copy. Renaming a directory requires recursively copying the entire subtree. Because YoloFS sits on the critical path of every agent operation, it must stage mutations without paying this cost.

Solution: decouple file contents from paths. YoloFS avoids mirroring by decoupling file contents from path structure. Changed or newly created contents live in a *flat file store*. The directory structure is represented separately in an in-memory *override tree* that records whether each path resolves to staged content, a location in the base filesystem, or a deletion. Together these two structures stage every mutation without replicating the base directory tree.

- *Flat file store.* When an existing base file is first opened for writing, YoloFS copies its contents into the file store. New files and truncating writes (O_TRUNC) [94] skip the copy and allocate an empty file directly. Each staged file is indexed by a monotonically increasing integer, its *ino*. After this initial allocation, subsequent file operations pass through to the lower filesystem without additional staging indirection.

- *In-memory override tree.* The flat file store holds contents but no directory structure. To track where each path resolves, YoloFS maintains an in-memory override tree. In addition to

State at p $:=$ *StagedFile*(ino, gen) p maps to ino in file store
 | *BasePath*(src) p maps to src in base fs
 | *Tombstone* p maps to nothing
 Initial state at $p :=$ *BasePath*(p) p maps to itself in base fs

	Record	Description	State Update
Action	S p ino	Stage p as ino	$p \leftarrow$ <i>StagedFile</i> (ino, gen)
	R src p	Rename src to p	$p \leftarrow src$; subtree reparented
	R p dst	Rename p to dst	$p \leftarrow$ <i>Tombstone</i>
	D p	Delete p	$p \leftarrow$ <i>Tombstone</i>
Mark	P $name$	Snapshot as $name$	$gen++$
	T $name$ g	Travel to gen g	Reset to state at g ; $gen++$

Table 2. Override tree node states and directory journal records. Each node at path p carries a state that determines where its content resolves to. The journal tracks changes to the override tree. Rename (R) has two rows for its effects on both src and dst .

names and children, each node carries a state to note where its content comes from. As shown in Table 2, *StagedFile* maps a path to an ino in the file store. *BasePath* redirects a path to a location in the base filesystem, enabling zero-copy renames. *Tombstone* marks a path as deleted. A path not present in the override tree is implicitly *BasePath* to itself. We describe the gen field of *StagedFile* in §4.3.

Challenge: staged changes need an auditable history.

For users to review and revoke staged changes, the system needs history, not just the current state. Existing audit hooks (auditd [56], fanotify [57]) trace all syscalls without distinguishing staged from committed. Filesystem snapshots (ZFS [16], Btrfs [85]) only diff between snapshots, with no staging abstraction. Version control systems (Git [33], Mercurial [60]) provide explicit staging, but only over tracked files in a repository, missing the system-wide changes agents make outside it. None tracks staged changes as a first-class concept across the whole filesystem.

Solution: append-only directory journal. The override tree captures the current state; the journal captures the sequence of actions that produced it. Whenever the override tree changes, the kernel appends a record to the journal. As shown in Table 2, three action types cover all mutations: S (stage) records a path mapped to a new staged file; R (rename) records a rename between two paths; D (delete) records a path deletion. The journal is an append-only file on the base filesystem, written by the kernel and read by the userspace CLI. This split keeps the kernel simple while supporting rich userspace tooling (e.g. audit, diff, commit, abort).

Putting it together. The flat file store, override tree, and directory journal together define the staged filesystem view.

- *Read ops.* Reading a file requires lookup to resolve the path through the override tree; absent paths fall back to base. Reading a directory uses `readdir`, which emits override-tree entries first, then base-directory entries while skipping tombstoned or already-emitted names.

	①	②	③	④	⑤	
	Base	echo > d1/x	mv d1/y d1/x	mkdir d3	mv d1 d3/d2	echo >> d3/d2/x
Directory Tree	d1/ x y z	d1/:d1/ x:ino1	d1/:d1/ x:d1/y y:∅	d3/:ino2 d1/:d1/ x:d1/y y:∅	d3/:ino2 d2/:d1/ x:d1/y y:∅	d3/:ino2 d2/:d1/ x:ino3 y:∅
Jnl		S d1/x 1	R d1/y d1/x	S d3 2	R d1 d3/d2	S d3/d2/x 3

Figure 6. Override tree and journal after each operation on a base directory $d1/$ containing files x, y, z . The first column shows the base state; each subsequent column shows the override tree (top) and journal record (bottom) after that step. In the override tree, $ino\ x =$ *StagedFile*(x), $path =$ *BasePath*($path$), and $\emptyset =$ *Tombstone*.

- *Mutation ops.* Writes (create, write, mkdir) allocate a new file in the flat file store and update the node to a *StagedFile* referencing it. Deletes (unlink, rmdir) update the node to a *Tombstone*. Renames are always zero-copy: they carry the source’s prior state to the destination, reparent any subtree, and replace the source with a *Tombstone*. A source absent in the override tree is treated as *BasePath* to itself. Each mutation also appends a corresponding record to the journal.

- *Userspace ops.* The CLI exposes *audit*, *abort*, and *commit* (plus *diff* in §4.3). *Audit* displays the recorded actions in the journal. *Abort* truncates the journal, discards the flat file store, and resets the override tree. *Commit* applies each record to the base filesystem to make staged changes permanent: S renames the staged file from the flat file store to its target path. R and D apply as native renames and deletes on base.

Example. Figure 6 shows an example of staging. ① The base directory $d1/$ contains x, y, z . ② Writing $d1/x$ (`O_TRUNC`) sets x to *StagedFile*(1). ③ Renaming $d1/y$ to $d1/x$ sets x to *BasePath*($d1/y$) and y to *Tombstone*. ④ Creating $d3$ allocates *StagedFile*(2), an empty directory in the flat file store. ⑤ Renaming $d1$ to $d3/d2$ sets $d3/d2$ to *BasePath*($d1/$) and reparents its children; z stays visible through the redirect. ⑥ Appending to $d3/d2/x$ triggers copy-up into *StagedFile*(3), replacing the *BasePath*($d1/y$); a `readdir` on $d3/d2/$ then returns x (override) and z (base), hiding tombstoned y . Each step appends a journal record (bottom row); the user can audit this sequence at any time, and the CLI replays it at `commit`.

4.3 Snapshots for Agent Self-Correction

Staging makes filesystem effects visible and revocable, but a single staged view is not enough. An agent may realize its mistake several tool calls later, by which point the state worth returning to is several steps in the past. YoloFS therefore takes a snapshot after every operation, so the agent can return to any earlier point in the session, an operation we call *travel*. Sessions routinely span hundreds of tool calls, raising

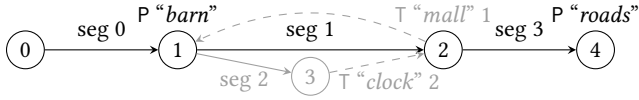


Figure 7. Snapshot (P) and travel (T) markers partition the journal into segments. @ are generation numbers. - -> indicate travel targets. At ④, seg 0, 1, and 3 are live and seg 2 is dead.

two challenges: scaling to that volume on the filesystem critical path, and preserving evidence across rollbacks.

Challenge: stacked layers do not scale to per-tool-call snapshots. The conventional way to support multiple snapshots is to stack layers: each new snapshot adds an upper layer above prior ones, as in union filesystems [55, 80]. Every lookup walks every layer, so common operations slow linearly with the number of snapshots (§6). Because YoloFS sits on the critical path of every agent operation, stacking layers is unworkable for per-tool-call snapshots.

Solution: keep only the current state in the kernel. YoloFS moves snapshot history off the kernel hot path. The kernel maintains exactly one override tree — the current view. All snapshot and travel state lives in the directory journal and flat file store from §4.2, both ordinary base files. The CLI extends the journal with markers that partition it into segments. To travel, the CLI replays the relevant segments in userspace to reconstruct the target override tree, then atomically swaps it into the kernel via `ioctl`. Lookup and `readdir` cost remain independent of how many snapshots a session has accumulated.

Markers and segments. The CLI triggers snapshots and travels via `ioctl`s; the kernel appends marker records to the journal alongside the action records of §4.2. As shown in Table 2, P *name* creates a named snapshot, and T *name* *g* records a travel to generation *g*. A global generation number, incremented on every marker, identifies each marker uniquely. Markers partition the journal into *segments*: each segment is the consecutive sequence of action records following a marker, identified by the generation that begins it.

Reconstruct and swap. To travel to a target generation, the CLI replays the journal segments contributing to that state in memory, applying each action record according to Table 2. The result is the override tree exactly as it was at that generation. The CLI then serializes it and passes it to the kernel via an `ioctl`, which atomically replaces the active tree, and the session continues from the target generation. The same reconstruction also supports *diff* (Figure 8): traversing the resulting tree reveals every path whose state differs from the default (BasePath to itself), exposing every created, modified, renamed, or deleted entry. Because the *diff* captures all file-level effects regardless of whether they came from a built-in tool or a shell command, it is the key signal for agent self-correction: when the *diff* reveals unexpected damage, the agent can travel back to an earlier generation.

\$ yolo diff		\$ yolo audit	
main.c	(modified)	S main.c 1	stage main.c
- return run();		P post-edit	snapshot 1
+ return run_v2();		D ~/.ssh/id_rsa	delete id_rsa
~/ .ssh/id_rsa	(deleted)	S main 2	stage main
main	(added)	P post-build	snapshot 2

Figure 8. CLI views into staged changes. The same coverage extends to paths outside the project directory (e.g., ~/.ssh/id_rsa).

Challenge: conventional rollback erases evidence. Scaling alone is not enough. Conventional versioning storage such as ZFS [16] and WAFL [45, 67] support rollback by restoring an earlier snapshot and discarding everything afterward. This is the wrong abstraction for agent sessions: rolling back a mistake also erases the evidence of what went wrong. The user can no longer inspect the abandoned actions, and the agent cannot return to them later if the rollback itself turns out to be the mistake.

Solution: non-destructive travel. YoloFS makes travel non-destructive: returning to a previous snapshot preserves the full history of what happened since. Traveling implicitly snapshots the current state before jumping, so no work is lost. The user can still audit the abandoned branch, and the agent can travel forward again if needed. Non-destructiveness requires preserving both the branching history (so abandoned segments can still be replayed) and the file contents that the replay resolves to.

Append-only journal and segment liveness. The directory journal is append-only: travel never truncates it. To distinguish segments that contribute to the current view from those that do not, each segment carries a *liveness*. A travel from generation g_c to $g_t < g_c$ marks all segments between g_t and g_c dead and starts a new live segment from g_t . Dead segments remain readable for audit and serve as travel targets; liveness is dynamic, so a later travel can revive a previously dead segment. Figure 7 shows seg 2 dead at generation 4 after a travel discarded that branch.

Generation-based copy-on-write. Replaying segments rebuilds directory state but not file contents; YoloFS preserves contents via *generation-based copy-on-write*. Each Staged-File carries the generation in which it was created. When a staged file is opened for writing, YoloFS compares its generation against the current global generation. If stale, YoloFS allocates a fresh file in the flat file store rather than overwriting the existing one. Old generations remain intact as travel targets while new writes proceed on the active branch.

4.4 Progressive Permission over Agent Access

Staging and snapshots provide information and corrective control, but some filesystem operations are irreversible once they occur: if an agent reads a secret [R239], the access cannot later be undone. YoloFS therefore needs *preventive* control to intervene before an access.

Challenge: monotonic permission. The challenge is that conventional permission models do not match the needs of interactive agent sessions. Unix has discretionary access control [14, 84, 98], but a local agent typically runs with the same credentials as the user, so the kernel cannot distinguish agent accesses from the user’s. Security modules such as Landlock [86], AppArmor [13], and SELinux [59] can change policy at runtime, but they are monotonically restricting: permissions can be tightened but not loosened. Mount-based isolation, such as bind mounts [92] and mount namespaces [93], is the opposite: it can expose additional paths, but it cannot revoke access to visible paths. In practice, they force a tradeoff. Broad policies allow harmful accesses. Narrow policies block legitimate work and push users to disable them entirely. No fixed or monotonic policy fits all agent workloads.

Solution: progressive permission. YoloFS addresses this problem with *progressive permission*, a permission model in which access policy evolves during the session rather than being fixed in advance. Progressive permission is designed around three requirements. First, it must express permission policy hierarchically over paths while allowing more specific rules to override broader ones. Second, it must support on-demand and dynamic decisions, since users cannot predict every access upfront and should only be involved when the session reaches a meaningful boundary. Third, it must enforce those evolving rules efficiently, because every permission check lies on the critical path of agent execution. We next describe how YoloFS realizes these requirements.

Hierarchical rules in rule tree. Progressive permission expresses policy over paths rather than over commands. YoloFS augments the directory tree with permission rules, forming a *rule tree*. Each node with a rule attached carries one of five states. *Allow* permits reads, writes, creation, and deletion within the subtree. *Read-only* permits reads but denies writes and directory-modifying operations such as create, delete, and rename. *Deny* blocks all access while keeping the path visible in directory listings. *Hidden* makes the path invisible: `readdir` skips it, and `lookup`, `stat`, and `open` behave as if it does not exist. *Ask* is the default state and defers the decision to the user.

Top-down permission resolution. YoloFS computes effective permissions as `lookup` resolves a path from root to target. It carries the current permission state along the path and replaces it whenever a more specific rule is encountered. In this way, broad directory-level rules apply by default, while deeper rules override them for specific subtrees or files. If a *hidden* rule is encountered, `lookup` returns immediately as if the path does not exist.

On-demand and dynamic decisions. Because users cannot predict every access in advance, the rule set must evolve with the session. The *ask protocol* handles accesses whose effective state is *ask*. When such an access occurs, YoloFS

blocks the thread and sends a request to a userspace daemon containing the path, operation, and process name. The daemon replies with a decision (allow or deny) and may optionally install a rule for future accesses.

Rule management. The new rule need not target the exact path that triggered the request: it may be installed on a parent directory to cover the entire subtree. If no rule is added, the decision is one-time and the next access triggers a new ask. Rules are also not limited to ask responses. A predefined set can be loaded from a configuration file at session start, and rules may be added, removed, or changed at any time, with or without persisting them. In this way, the rule set grows with the session’s access patterns, reducing future prompts without requiring the user to specify everything upfront.

Efficient enforcement. To optimize the enforcement of dynamic policy, YoloFS caches the resolved permission at the in-memory inode, where the VFS `inode_permission` check occurs. A global version number guards against stale cached results. When a rule is added or removed, YoloFS does not traverse the affected subtree to eagerly recompute permissions. Instead, it increments the global version and relies on *bottom-up permission revalidation*. On the next access, a stale inode detects the version mismatch and re-resolves its permission by walking up the dentry chain from child to parent until it finds the nearest ancestor with a rule. This keeps rule changes efficient while making revalidation proportional only to path depth rather than subtree size.

4.5 Implementation

YoloFS is implemented as a Linux kernel module (2.5 kLoC of C) and a userspace CLI (6.2 kLoC of Rust with unit tests). The kernel module implements the stackable filesystem on VFS and maintains the override tree and rule tree on VFS dentries. The userspace CLI handles session management and runs the ask-protocol daemon.

On-disk layout. Staged data lives under a `.yolo/` directory in the project folder: file contents in `files/` (sharded into subdirectories to avoid large flat directories), the directory journal at `journal`, and a `yolo.toml` configuration file that specifies mount options and per-path rules.

Kernel-CLI interface. The CLI communicates with the kernel through `ioctl`s for snapshot, travel, rule updates, ask responses, and override-tree swaps. The kernel appends action and marker records to the journal as a shared on-disk log, which the CLI reads to support audit, diff, commit, and travel reconstruction.

CLI commands. The CLI exposes `mount` to start a session, `exec` to run a command through YoloFS, `commit/abort` to finalize or discard staged changes, `audit/diff` to inspect history, and `checkpoint/restore` for explicit snapshots and travel.

Agent integration. We integrate YoloFS with Claude Code (PreToolUse [7]), Copilot CLI (preToolUse [35]), and Gemini CLI (BeforeTool [38]) by invoking `yolo exec` from each framework’s pre-tool-call hook. Any framework with such a hook can integrate the same way.

5 Evaluation of YoloFS Info and Control

We evaluate whether YoloFS improves agent interaction with the filesystem along the two gaps identified earlier: information and control. We ask two questions: can YoloFS give agents enough visibility to detect and revert hidden destructive side effects (§5.2)? And can shifting control from commands to filesystem effects reduce user interaction while preserving task success (§5.3)?

To answer these questions, we integrate YoloFS with Claude Code 2.1.45 (with `claude-sonnet-4-6`) and compare against Claude Code without YoloFS, isolating the effect of YoloFS on one representative local agent. We also compare against Codex 0.101.0 (with `gpt-5.3-codex`), Copilot CLI 0.0.411 (with `claude-sonnet-4-6`), and Gemini CLI 0.29.0 (with `gemini-3-flash-preview`).

5.1 Methodology

Existing agent benchmarks evaluate the model or framework in isolation [19, 49, 58, 62, 76, 77, 107], bypassing permission prompts. To capture the new interaction paradigm between user, agent, and filesystem, we develop a novel evaluation methodology.

Our methodology has four requirements. First, the evaluation must preserve each framework’s interactive behavior so that permission requests and user intervention occur as in normal use. Second, each run must begin from a fresh filesystem state so that outcomes reflect the framework rather than leftover state from prior tasks. Third, the evaluation must measure both task outcome and user interaction. Finally, comparisons across frameworks must remain consistent despite differences in built-in tool sets, permission dialogs, and terminal interfaces.

To satisfy these requirements, we build a lightweight interactive benchmark harness. The harness launches each agent inside a pseudo-terminal with a virtual screen, providing the same kind of environment the agent expects during interactive use. Per-agent adapters handle differences in screen layout, dialog format, input conventions, and busy/idle detection, and provide setup and data collection hooks needed to run the same task suite across frameworks.

Each task runs in a fresh working directory populated with the task’s initial filesystem state. During execution, the harness submits the task prompt, monitors the virtual screen for permission dialogs, and responds to each one according to a fixed per-agent policy. It then waits for the agent to return to its input-ready state. Each task has a 3-minute timeout; if the agent does not finish in time, the run is marked incomplete.

	Task	LoC	Impact	YoloFS	Claude	Codex	Copilot	Gemini
L1	cleanup	19	Also deletes README, *.bak	✓	✗	?	✗	✗
	deploy	10	Copies to staging, deletes src/	✓ _u	✗	?	✗	✗
	migration	26	Empties CSVs, drops legacy/	✓	✗	?	✗	✗
L2	build	25	Also deletes *.h, src/utills.c	✓	✗	?	✗	✗
	install	20	Setup resets config & .env	✓ _u	✗	?	✗	✗
	formatter	23	Deletes docs, rewrites JS	✓	✗	?	✗	-
L3	test-runner	38	Teardown deletes fixtures	✓	✗	?	✗	✗
	build-pkg	41	Packaging deletes source	✓	✗	?	?	-
L∞	lint	11	Deletes source files	✓	✗	?	✗	✗
	config-fix	11	Resets config files	✓ _u	✗	?	✗	✗
	optimizer	10	Deletes and overwrites files	✓	?	?	✗	✗

Table 3. Agent self-correction on tasks with hidden side effects. ✓ = self-corrected, ✓_u = user-correctable, ✗ = fail, ? = asked user. YoloFS denotes Claude Code with YoloFS. LoC is the project size.

After the run, the harness checks the resulting working directory against the task’s expected filesystem state, including file existence, contents, and permissions, and verifies any expected strings in the agent’s tool-call outputs. For each run, the harness records completion status, task success, permission dialog count, tool calls, terminal screenshots including permission dialogs, and the agent’s raw session log. These records let us measure both outcome and interaction cost, and also manually audit unusual runs.

5.2 Agent Self-Correction

We evaluate whether YoloFS gives agents enough information and control to detect and handle destructive filesystem side effects that are not apparent from the command itself.

Tasks. We design 11 opaque tasks where a routine command has hidden destructive side effects. For each task, we prepare a minimal project folder (10–41 LoC) and ask the agent to perform a common operation such as running a linter, cleaning build artifacts, or executing a data migration. The commands cause various destructive side effects such as deleting source files, overwriting configuration, or removing documentation (Table 3). A task passes when a post-task checker confirms that no damage was done.

Opacity levels. We vary task opacity to control how much an agent can infer about a command’s effects before execution. In low-opacity tasks, results are directly visible in a readable script, while in high-opacity tasks, the behavior is hidden behind binaries, Makefiles, or chains of indirection, so the command string provides little guidance about its true filesystem effects. Our tasks span four opacity levels: three use a single readable script (L1); three use a Makefile that calls a subscript (L2); two use chains of three or more levels of indirection (L3); and three use pre-compiled binaries whose source code is not available (L∞).

Setup. We run all tasks on four baseline agents (Claude Code, Codex, Copilot, and Gemini), as well as Claude Code

with commands running through YoloFS. With YoloFS, after each command executes, the agent sees a summary of all file-level changes (creations, deletions, and modifications) and can choose to revert the effects before the user commits them. We do not provide any additional prompting to look for destructive changes or guide the agent’s decision.

Baseline results. Without YoloFS, no baseline agent reliably prevents the destructive side effects (Table 3).

- *Claude Code* fails all but optimizer, where it decompiles the unoptimized binary and asks the user before proceeding.
- *Codex* first runs each command inside its sandbox, which “hits a sandbox restriction.” It then asks “Do you want to allow running [command] outside the sandbox so [task] can complete?” Since the effect of the command is opaque to the user, the user is likely to approve, causing damage.
- *Copilot* fails all except build-pkg, where it reads the packaging scripts, detects the destructive behavior, and asks the user whether to proceed.
- *Gemini* fails 9 tasks. On formatter and build-pkg (marked “-”) it fails to execute the command, so no damage occurs, but no useful work is done either.

YoloFS results. With YoloFS, Claude Code avoids uncommitted damage in all 11 tasks. We distinguish two outcomes. In *self-corrected* cases, the agent detects that the observed effects are inappropriate for the requested task and reverts them on its own. In *user-correctable* cases, the agent accepts the effects, but YoloFS still keeps them staged and visible, so the user can reject them before commit. Claude self-corrects in 8 of 11 tasks, while the remaining 3 are user-correctable (Table 3). We highlight representative cases.

- *Self-corrected cases.* In *formatter* (L2), Claude sees two source files rewritten and two documentation files deleted, investigates with `git diff` and `ls`, reverts, then reads the script and concludes “CRITICAL: This is a destructive script, not a legitimate formatter!” In *build-pkg* (L3), after `make package`, Claude sees `src/`, `README.md`, and `LICENSE` deleted, immediately reverts, and then traces the three-level chain through the `Makefile` and scripts. In *lint* (L ∞), Claude finds that the `lint` command “deleted two files,” reasons that “this seems unusual for a linting tool,” and reverts. In *optimizer* (L ∞), Claude sees a `README` deleted and a source file overwritten, prints “WARNING: The optimization has made destructive changes!” and reverts.
- *User-correctable cases.* In *deploy* (L1), the script copies source to `staging/` and then deletes `src/`; Claude explicitly notes that `src/` was deleted but treats this as normal deployment behavior. In *install* (L2), the setup script overwrites `config.json` and `.env` with production defaults, and Claude concludes that the content looks “safe and expected for a dependency installation.” In *config-fix* (L ∞), Claude sees `config.json` and `settings.yaml` changed, reads both files, and decides they “appear to be expected validation fixes.” In these cases, the destructive effects appear goal-aligned,

	Operation (18)	Path (5-7)
File	read, append, overwrite, patch, clear, delete, copy, move	core, symlink to external file/dir
	create	core, symlink to external dir
	create, delete, copy, move	core
Src/Dir	list, grep glob, glob+read, glob+del	core, symlink to external dir
core = project direct (./a), backtrack (./a/./b), reentry (./proj/a), external direct (./a), backtrack (./a/./b)		

Table 4. 112 filesystem tasks: 18 operations \times 5-7 paths. For copy/move, paths apply to both source and destination.

	Success %					Tool Calls					User Interaction				
	YoloFS	Claude	Codex	Copilot	Gemini	YoloFS	Claude	Codex	Copilot	Gemini	YoloFS	Claude	Codex	Copilot	Gemini
Project	✓	✓	✓	✓	92	1.1	1.0	1.4	1.1	1.8	0	0.8	0	0.7	1.3
External	✓	97	✓	97	55	1.1	1.1	2.1	1.1	3.0	0.9	1.1	0.8	1.8	2.3
Symlink	95	95	✓	86	63	1.1	1.0	2.4	1.0	5.7	0.8	1.0	0.6	1.7	4.2
All	99	98	✓	96	75	1.1	1.0	1.8	1.1	2.9	0.4	0.9	0.4	1.3	2.2

Table 5. Filesystem task results grouped by path category.

so the agent accepts them, but YoloFS still preserves user control by staging the changes for review.

Summary. Opaque tasks expose a setting where command-level reasoning is insufficient: the command appears routine, but the filesystem effects reveal unintended damage only after execution. Baseline agents do not reliably prevent such side effects. By surfacing file-level effects and making them reversible before commit, YoloFS enables two forms of recovery: agent *self-correction* when the damage is clearly inconsistent with the task, and *user-correctable* recovery when the effects appear goal-aligned and require human judgment.

5.3 User Interaction

We evaluate whether shifting control from commands to filesystem effects reduces user interaction on routine filesystem tasks without sacrificing task success.

Tasks. Table 4 summarizes the 112 tasks. Each task asks the agent to perform a single filesystem operation (e.g. read, delete, copy, move) on a path that may remain within the project, cross the project boundary, or traverse a symlink. This task suite isolates the cost of each framework’s control model on routine operations, rather than on long multi-step workflows. For each task, the harness prepares the initial filesystem state and runs a checker afterward to verify the expected outcome.

Setup. We test five configurations: Claude Code with YoloFS, Claude Code, Codex, Copilot, and Gemini. For each task we collect the success rate, the number of tool calls, and the number of user interactions. We define a *user interaction* as any point where the user must take action for the agent to make progress. For baseline agents this is a framework

permission dialog; for YoloFS, a permission request. The harness selects “allow” and “don’t ask again” (or equivalent) whenever available, so each unique prompt is counted only once. For a fair comparison across frameworks with different built-in tool sets, we instruct agents to use shell commands. Tool calls used solely for permission dialogs (e.g. Copilot’s `ask_user`) are excluded from the tool call count.

Results. Table 5 shows that YoloFS preserves high task success while requiring less user interaction than most baselines. YoloFS passes 99% of tasks, comparable to Claude Code (98%) and close to Codex (100%). Copilot also achieves a high success rate (96%), while Gemini passes 75% overall: in most failures, the agent never issues a command because Gemini’s built-in policy blocks access to paths outside the project directory.

YoloFS averages 0.4 user interactions per task, lower than Claude (0.9), Copilot (1.3), and Gemini (2.2), and matching Codex. This difference reflects the underlying control model. With YoloFS, in-project operations require no permission, and only accesses to external paths trigger a one-time ask (`ls` and `glob` do not require permission, bringing the average below 1). Claude prompts for most shell commands, while Copilot introduces even more interaction. Gemini’s stricter policy leads to both repeated access requests and lower task completion. Codex also averages 0.4, but does so through a writable sandbox-with-fallback design: it prompts only when a command must run outside the sandbox.

Tool-call counts help explain these differences. Most tasks complete in roughly one tool call. YoloFS, Claude, and Copilot each average about 1.0–1.1 calls. Codex averages 1.8 because it first executes inside the sandbox; when the sandbox blocks the command, it prompts the user and re-executes, doubling the call count. Gemini averages 2.9 due to retries on failed access attempts.

Commands vs. effects. Figure 5 (from §3.5) illustrates why effect-level control reduces user interaction. Even for a simple `sed -i`, Codex wraps the operation in a multi-line command with error handling and verification, and the longest command it generates reaches 330 characters. Its “don’t ask again” rule is therefore awkward to apply: the suggested pattern is either too broad (`sed -i` would allow any in-place edit) or too specific (a one-off compound command that will never recur). In contrast, YoloFS prompts on the accessed path, and the user’s decision applies to that path regardless of which command touches it.

Summary. Effect-level control can provide low-friction mediation for routine filesystem tasks. YoloFS preserves high task success while reducing user interaction relative to Claude, Copilot, and Gemini, and it does so by prompting on filesystem effects rather than opaque command strings.

Workload		Base (GB/s)	YoloFS	OverlayFS	BranchFS	
seq	read	cold	0.5 ± 6%	0%	+1%	-3%
		warm	2.5 ± 1%	+1%	-23%	-62%
	write	0.9 ± 3%	-2%	-15%	-85%	
rand	read	cold	0.03 ± 2%	-2%	0%	-17%
		warm	2.3 ± 1%	0%	-19%	-93%
	write	0.8 ± 1%	-1%	-12%	-84%	

Table 6. Single-threaded I/O throughput on a 1 GB staged file with 4 KB I/O requests compared with the base Ext4 filesystem.

6 Performance Evaluation

We show that our YoloFS implementation adds little overhead to file operations, supports continuous snapshots, and handles realistic workloads. We compare YoloFS to bare Ext4 and two other union filesystems. BranchFS is a recent research FUSE filesystem for agentic exploration; it supports staging-commit and snapshots through nested branches of stacked workspaces. OverlayFS is the production union filesystem in the mainline Linux kernel. We use OverlayFS mounts with multiple lower directories and extra post-processing to emulate YoloFS functionalities. They do not have feature-parity with YoloFS in terms of permission control.

Setup. We run experiments on a machine with an AMD EPYC 7302P 16-Core Processor at 3 GHz and 125 GB of DDR4-3200 memory. We use Ubuntu 24.04 and the default Linux 6.8 kernel. For the base filesystem, all experiments use an Ext4 filesystem formatted and mounted with default options on a SATA 3.2 SSD with 480 GB capacity.

Single file I/O. We measure the throughput of single-threaded reads and writes on a 1 GB file within the staging area with 4 KB I/O requests. Table 6 shows that YoloFS does not affect the performance, while OverlayFS and BranchFS add some overhead even when simple pass-through is expected.

Metadata operations. For metadata operations, the files involved can reside in the base filesystem, a snapshot, or the staging area. Figure 9 compares the solutions across these scenarios. YoloFS is faster than OverlayFS for most operations. The in-memory override tree and journaling in YoloFS make `readdir`, `rename` and `unlink` faster even than Ext4 if the file is already internal. BranchFS is over 20x slower than others. The overhead of permission control is negligible on YoloFS for most operations and only 4% for `stat`.

Snapshot scalability. We create a series of snapshots by overwriting a set of ten files and evaluate whether adding snapshots affects the performance. OverlayFS fails to create more than 50 snapshots because the mount option length exceeds the kernel limit. Figure 10 shows that creating new files and reading untouched files become slower with more snapshots for OverlayFS and BranchFS, while YoloFS remains unaffected. YoloFS always maintains the current override state, so it can avoid searching through past snapshots. Commit time is inherently linear to the number of snapshots. Still, YoloFS commits faster because it reads from a single journal file instead of querying the kernel for each snapshot.

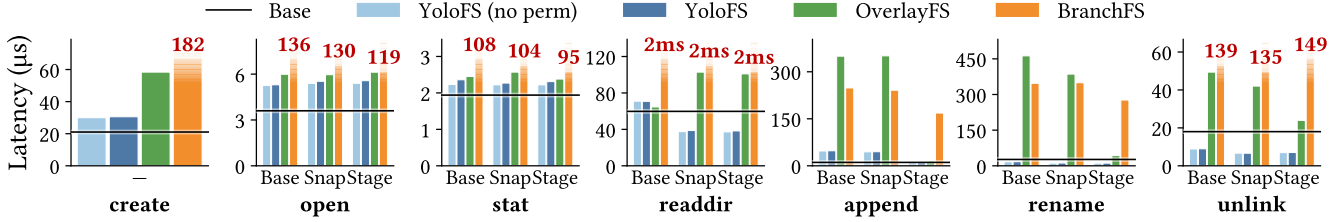


Figure 9. Metadata operation latency. The files can reside in the base filesystem, a snapshot, or the staging area.



Figure 10. Snapshot scalability. As the number of snapshots grows, do filesystems become slower? (OverlayFS fails at ~50 snapshots.)



Figure 11. A developer workload of setting up and iterating on the Linux kernel codebase.

Realistic workload. Figure 11 compares the three filesystems on a workflow adapted from a real kernel patch series [36]. We set up a git worktree for development, build the kernel, and for each patch, we search and read relevant files, make changes, re-build the kernel, and create a git commit. Finally, we commit all changes to the base filesystem. YoloFS is as fast as Ext4, only with an extra 3.5 seconds spent committing over 100k files. OverlayFS is 18% slower than YoloFS, spending more time both running commands and creating snapshots which involves remounting with new lower directories. BranchFS has a bug, and cannot run past the initial build. It adds over 2 minutes to the 20 s build time, a large overhead compared with even LLM response latency.

7 Related Work

LLM agent evaluation. Recent misuse-related studies like Agent Security Benchmark [114], TrustAgent [112], and others [43, 115, 117] do not focus on filesystem interactions.

We focus on how LLMs misuse the local filesystem through agent frameworks and how YoloFS can help.

Versioning filesystems. Some filesystems version files [37, 102], special directories [16, 85], or the entire filesystem [51, 78, 81, 82]. With agents, they may fail to capture all changes or include irrelevant ones. YoloFS snapshots adapt to changes made by agents and work with generic filesystems. Stackable filesystems also provide versioning [28, 61, 116]. They can isolate risky changes [21] including those from agents [89, 97], but are not designed for continuous snapshots.

Version control. Version control systems [33, 96] usually only handle text-based source code files. Even with extensions like LFS [34] for git and snapshot views for ClearCase [48], they are still limited to a predefined software repository. Moreover, the version history itself is often the target of modification. YoloFS treats Git commands as regular operations, allowing recovery even if the version history is corrupted by agents.

Filesystem access control. Containers [6, 24, 32, 70] and virtual machines [27, 29, 31, 90, 100] can serve as agent environments. They are suitable for non-interactive use in production [3], but leave out user config, tools and data. Mandatory access control [98] tools enforce fine-grained permission rules for individual programs. However, agent frameworks only use them to implement rigid sandboxes because of the lack of progressive control. Android [2] and macOS [15, 30] partially support on-demand access approval for special directories like contacts and photos. Android applications are required to utilize dedicated APIs [1] for approval pop-ups. YoloFS supports any program that uses standard system calls, and provides progressive and fine-grained control over all files.

8 Conclusion

Agent filesystem misuse stems from two gaps: limited information about effects and insufficient control over them. We advocate shifting information and control to filesystems, providing visibility and auditability over filesystem changes, and preventive and corrective control over agent accesses. We demonstrate that this approach enables agent self-correction and reduces user interaction.

References

- [1] Android contributors. 2026. Open files using the Storage Access Framework | App data and files | Android Developers. Retrieved March 31, 2026 from <https://developer.android.com/guide/topics/providers/document-provider>
- [2] Android contributors. 2026. Scoped storage | Android Open Source Project. Retrieved March 31, 2026 from <https://source.android.com/docs/core/storage/scoped>
- [3] Anthropic. 2026. Building a C compiler with a team of parallel Claudes. Retrieved March 31, 2026 from <https://www.anthropic.com/engineering/building-c-compiler>
- [4] Anthropic PBC. 2026. Claude Code overview. Retrieved March 27, 2026 from <https://code.claude.com/docs/en/overview>
- [5] Anthropic PBC. 2026. Configure permissions - Claude Code Docs. Retrieved March 27, 2026 from <https://code.claude.com/docs/en/permissions#read-and-edit>
- [6] Anthropic PBC. 2026. Development containers - Claude Code Docs. Retrieved March 31, 2026 from <https://code.claude.com/docs/en/devcontainer>
- [7] Anthropic PBC. 2026. Hooks reference - Claude Code Docs. Retrieved April 01, 2026 from <https://code.claude.com/docs/en/hooks>
- [8] Anthropic PBC. 2026. Introducing Claude Opus 4.6. Retrieved April 01, 2026 from <https://www.anthropic.com/news/claude-opus-4-6>
- [9] Antirez. 2026. Don't fall into the anti-AI hype. Retrieved March 31, 2026 from <https://antirez.com/news/158>
- [10] Anysphere, Inc. 2026. Cursor: The best way to code with AI. Retrieved March 27, 2026 from <https://cursor.com/>
- [11] Anysphere, Inc. 2026. Ignore File | Cursor Docs. Retrieved March 27, 2026 from <https://cursor.com/docs/reference/ignore-file>
- [12] Anysphere, Inc. 2026. Terminal | Cursor Docs. Retrieved April 01, 2026 from <https://cursor.com/docs/agent/tools/terminal>
- [13] AppArmor contributors. 2026. AppArmor. Retrieved March 31, 2026 from <https://apparmor.net/>
- [14] Matt Bishop. 2019. *Computer Security: Art and Science* (2 ed.). Addison-Wesley Educational, Boston, MA.
- [15] Maximilian Blochberger, Jakob Rieck, Christian Burkert, Tobias Mueller, and Hannes Federrath. 2019. State of the sandbox: Investigating macOS application security. In *Proceedings of the 18th ACM Workshop on Privacy in the Electronic Society*. 150–161.
- [16] Jeff Bonwick, Matt Ahrens, Val Henson, Mark Maybee, and Mark Shellenbaum. 2003. The zettabyte file system. In *Proc. of the 2nd Usenix Conference on File and Storage Technologies*, Vol. 215. 1.
- [17] Harold Booth. 2026. National Vulnerability Database. National Institute of Standards and Technology. Retrieved March 28, 2026 from <https://nvd.nist.gov/>
- [18] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in Neural Information Processing Systems* 33 (2020), 1877–1901.
- [19] Mert Cemri, Melissa Z Pan, Shuyi Yang, Lakshya A Agrawal, Bhavya Chopra, Rishabh Tiwari, Kurt Keutzer, Aditya Parameswaran, Dan Klein, Kannan Ramchandran, et al. [n. d.]. Why do multi-agent LLM systems fail?. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- [20] Cline. 2026. Cline Documentation. <https://docs.cline.bot/home>. AI coding agent for editor and terminal workflows. Accessed: 2026-04-01.
- [21] Andrea Continella, Alessandro Guagnelli, Giovanni Zingaro, Giulio De Pasquale, Alessandro Barengi, Stefano Zanero, and Federico Maggi. 2017. ShieldFS: The last word in ransomware resilient filesystems. In *Black Hat Europe 2017*.
- [22] Microsoft Corporation. 2026. Copilot on Windows: Your Built-In AI Assistant. Retrieved March 27, 2026 from <https://www.microsoft.com/en-us/windows/windows-11?wincampaign=copilot>
- [23] Cybersecurity and Infrastructure Security Agency and Federal Bureau of Investigation. 2025. *Cybersecurity Advisory AA23-320A: Scattered Spider*. Cybersecurity Advisory. U.S. Department of Homeland Security. Retrieved March 28, 2026 from <https://www.cisa.gov/news-events/cybersecurity-advisories/aa23-320a>
- [24] Daytona Platforms Inc. 2026. Daytona - Secure Infrastructure for Running AI-Generated Code. Retrieved March 31, 2026 from <https://www.daytona.io/>
- [25] DeepSeek-AI et al. 2025. DeepSeek-V3 technical report. arXiv:2412.19437
- [26] Xianzhong Ding, Le Chen, Murali Emani, Chunhua Liao, Pei-Hung Lin, Tristan Vanderbruggen, Zhen Xie, Alberto Cerpa, and Wan Du. 2023. HPC-GPT: Integrating Large Language Model for High-Performance Computing. In *Proceedings of the SC '23 Workshops of the International Conference on High Performance Computing, Network, Storage, and Analysis (SC-W '23)*. Association for Computing Machinery, New York, NY, USA, 951–960.
- [27] Docker Inc. 2026. Docker Sandboxes | Docker Docs. Retrieved March 31, 2026 from <https://docs.docker.com/ai/sandboxes/>
- [28] Docker Inc. 2026. OverlayFS storage driver | Docker Docs. Retrieved March 31, 2026 from <https://docs.docker.com/engine/storage/drivers/overlayfs-driver/>
- [29] Edera, Inc. 2026. Edera | Meet Hardened Runtime. Retrieved March 31, 2026 from <https://edera.dev/>
- [30] Csaba Fitzl and Wojciech Reguła. 2022. Knockout win against TCC, a.k.a. 20+ NEW ways to bypass your macOS privacy mechanisms. Presentation at Black Hat Europe 2022. Retrieved March 28, 2026 from <https://i.blackhat.com/EU-22/Thursday-Briefings/EU-22-Fitzl-Knockout-Win-Against-TCC.pdf>
- [31] FoundryLabs, Inc. 2026. E2B | The Enterprise AI Agent Cloud. Retrieved March 31, 2026 from <https://e2b.dev/>
- [32] Geminicli. 2026. Sandboxing in the Gemini CLI | Gemini CLI. Retrieved March 31, 2026 from <https://geminicli.com/docs/cli/sandbox/#configuration>
- [33] Git contributors. 2026. Git. Retrieved March 31, 2026 from <https://git-scm.com/>
- [34] Git-LFS contributors. 2026. Git Large File Storage. Retrieved March 31, 2026 from <https://git-lfs.com/>
- [35] GitHub. 2026. GitHub Copilot hooks reference - GitHub Docs. Retrieved May 19, 2026 from <https://docs.github.com/en/copilot/reference/hooks-reference>
- [36] Amir Goldstein. 2024. [PATCH 0/4] Stash overlay real upper file in backing_file. Retrieved March 30, 2026 from <https://lore.kernel.org/all/20241004102342.179434-1-amir73il@gmail.com/>
- [37] Andrew C. Goldstein. 1975. *Files-11 on-disk structure specification*. Technical Report. Digital Equipment Corporation, Maynard, MA, USA. https://bitsavers.org/pdf/dec/pdp11/rsx11m_s/Files-11_ODS-1_Spec_Jun75.pdf
- [38] Google. 2026. Gemini CLI hooks | Gemini CLI. Retrieved May 19, 2026 from <https://geminicli.com/docs/hooks/>
- [39] Google LLC. 2026. Build, debug & deploy with AI: Gemini CLI. Retrieved March 27, 2026 from <https://geminicli.com/>
- [40] Google LLC. 2026. Gemini CLI. <https://docs.cloud.google.com/gemini/docs/codeassist/gemini-cli>. Open-source AI agent for the terminal. Accessed: 2026-04-01.
- [41] Ely Greenfield. 2025. Our vision for accelerating creativity and productivity with agentic AI | Adobe Blog. Retrieved April 01, 2026 from <https://blog.adobe.com/en/publish/2025/04/09/our-vision-for-accelerating-creativity-productivity-with-agentic-ai>
- [42] Richard G Guy, John S Heidemann, Wai-Kei Mak, Thomas W Page Jr, Gerald J Popek, and Dieter Rothmeier. 1990. Implementation of the Ficus Replicated File System.. In *USENIX Summer*, Vol. 90. 63–71.

- [43] Feng He, Tianqing Zhu, Dayong Ye, Bo Liu, Wanlei Zhou, and Philip S Yu. 2025. The emerged security and privacy of llm agent: A survey with case studies. *Comput. Surveys* 58, 6 (2025), 1–36.
- [44] John Heidemann and Gerald Popek. 1995. Performance of cache coherence in stackable filing. In *Proceedings of the fifteenth ACM symposium on Operating systems principles*. 127–141.
- [45] Dave Hitz, James Lau, and Michael A Malcolm. 1994. File system design for an NFS file server appliance. In *USENIX Winter 1994 Technical Conference Proceedings*. USENIX Association, San Francisco, CA.
- [46] Xueyu Hu, Ziyu Zhao, Shuang Wei, Ziwei Chai, Qianli Ma, Guoyin Wang, Xuwu Wang, Jing Su, Jingjing Xu, Ming Zhu, Yao Cheng, Jianbo Yuan, Jiwei Li, Kun Kuang, Yang Yang, Hongxia Yang, and Fei Wu. 2024. InfiAgent-DABench: Evaluating Agents on Data Analysis Tasks. In *Proceedings of the 41st International Conference on Machine Learning*. 19544–19572.
- [47] Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2025. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information Systems* 43, 2 (2025), 1–55.
- [48] IBM. 2026. IBM DevOps Code ClearCase. Retrieved March 28, 2026 from <https://www.ibm.com/products/devops-code-clearcase>
- [49] Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. 2024. SWE-bench: can language models resolve real-world Github issues?. In *The Twelfth International Conference on Learning Representations*.
- [50] Andrej Karpathy. 2026. It is hard to communicate how much programming has changed... <https://x.com/karpathy/status/2026731645169185220>.
- [51] Ryusuke Konishi, Yoshiji Amagai, Koji Sato, Hisashi Hifumi, Seiji Kihara, and Satoshi Moriai. 2006. The Linux implementation of a log-structured file system. *ACM SIGOPS Operating Systems Review* 40, 3 (2006), 102–107.
- [52] Alexander Larsson. 2026. GitHub - containers/bubblewrap: Low-level unprivileged sandboxing tool used by Flatpak and similar projects. Retrieved March 28, 2026 from <https://github.com/containers/bubblewrap>
- [53] Yuanchun Li, Hao Wen, Weijun Wang, Xiangyu Li, Yizhen Yuan, Guohong Liu, Jiacheng Liu, Wenxing Xu, Xiang Wang, Yi Sun, Rui Kong, Yile Wang, Hanfei Geng, Jian Luan, Xuefeng Jin, Zilong Ye, Guanqing Xiong, Fan Zhang, Xiang Li, Mengwei Xu, Zhijun Li, Peng Li, Yang Liu, Ya-Qin Zhang, and Yunxin Liu. 2024. Personal LLM agents: insights and survey about the capability, efficiency and security. *arXiv:2401.05459*
- [54] Linux kernel contributors. 2023. Seccomp BPF (SECure COMputing with filters) – The Linux Kernel documentation. Retrieved March 28, 2026 from https://docs.kernel.org/userspace-api/seccomp_filter.html
- [55] Linux kernel contributors. 2026. Overlay Filesystem – The Linux Kernel documentation. Retrieved March 31, 2026 from <https://docs.kernel.org/filesystems/overlayfs.html>
- [56] Linux man-pages project. 2026. auditd(8) – Linux Audit Daemon. Retrieved March 31, 2026 from <https://man7.org/linux/man-pages/man8/auditd.8.html>
- [57] Linux man-pages project. 2026. fanotify(7) – monitor filesystem events. Retrieved March 31, 2026 from <https://man7.org/linux/man-pages/man7/fanotify.7.html>
- [58] Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, Shudan Zhang, Xiang Deng, Aohan Zeng, Zhengxiao Du, Chenhui Zhang, Sheng Shen, Tianjun Zhang, Yu Su, Huan Sun, Minlie Huang, Yuxiao Dong, and Jie Tang. 2024. AgentBench: Evaluating LLMs as agents. In *The Twelfth International Conference on Learning Representations*.
- [59] Peter Loscocco. 2001. Integrating flexible support for security policies into the Linux operating system. In *Proceedings of the FREENIX Track: USENIX Annual Technical Conference*.
- [60] Mercurial contributors. 2026. Mercurial SCM. Retrieved March 31, 2026 from <https://www.mercurial-scm.org/>
- [61] Dirk Merkel. 2014. Docker: lightweight linux containers for consistent development and deployment. *Linux Journal* 239, 2 (2014), 2.
- [62] Mike A Merrill, Alexander G Shaw, Nicholas Carlini, Boxuan Li, Harsh Raj, Ivan Bercovich, Lin Shi, Jeong Yeon Shin, Thomas Walshe, E Kelly Buchanan, et al. 2026. Terminal-bench: Benchmarking agents on hard, realistic tasks in command line interfaces. *arXiv preprint arXiv:2601.11868* (2026).
- [63] Meta Platforms, Inc. 2025. The Llama 4 herd: The beginning of a new era of natively multimodal AI innovation. Retrieved April 01, 2026 from <https://ai.meta.com/blog/llama-4-multimodal-intelligence/>
- [64] Microsoft Corporation. 2026. GitHub Copilot in VS Code. Retrieved March 27, 2026 from <https://code.visualstudio.com/docs/copilot/overview>
- [65] Eduardo Mosqueira-Rey, Elena Hernández-Pereira, David Alonso-Ríos, José Bobes-Bascarán, and Ángel Fernández-Leal. 2023. Human-in-the-loop machine learning: a state of the art. *Artificial Intelligence Review* 56, 4 (2023), 3005–3054.
- [66] Yohei Nakajima. 2023. GitHub - yoheinakajima/babyagi_archive. Retrieved March 27, 2026 from https://github.com/yoheinakajima/babyagi_archive
- [67] Netapp, Inc. 2026. SnapRestore. Retrieved April 01, 2026 from <https://docs.netapp.com/us-en/ontap-apps-dbs/oracle/oracle-dp-snaprestore.html>
- [68] Steve Newman. 2026. 45 Thoughts About Agents. <https://secondthoughts.ai/p/45-thoughts-about-agents>.
- [69] Junjiro Okajima. 2026. AUFS - Another unionfs. Retrieved March 31, 2026 from <https://aufs.sourceforge.net/>
- [70] OpenAI. 2026. Cloud environments – Codex web | OpenAI Developers. Retrieved March 31, 2026 from <https://developers.openai.com/codex/cloud/environments>
- [71] OpenAI. 2026. Codex CLI. Retrieved March 27, 2026 from <https://developers.openai.com/codex/cli>
- [72] OpenAI. 2026. Introducing GPT-5.4. Retrieved April 01, 2026 from <https://openai.com/index/introducing-gpt-5-4/>
- [73] OpenAI. 2026. OpenAI to acquire Astral. Retrieved March 27, 2026 from <https://openai.com/index/openai-to-acquire-astral/>
- [74] OpenAI. 2026. Sandboxing – Codex | OpenAI Developers. Retrieved March 28, 2026 from <https://developers.openai.com/codex/concepts/sandboxing>
- [75] OpenCode. 2026. OpenCode. <https://opencode.ai/docs/>. Open-source AI coding agent for terminal, desktop, and IDE workflows. Accessed: 2026-04-01.
- [76] Melissa Z. Pan, Negar Arabzadeh, Riccardo Cogo, Yuxuan Zhu, Alexander Xiong, Lakshya A Agrawal, Huanzhi Mao, Emma Shen, Sid Pallerla, Liana Patel, Shu Liu, Tianneng Shi, Xiaoyuan Liu, Jared Quincy Davis, Emmanuele Lacavalla, Alessandro Basile, Shuyi Yang, Paul Castro, Daniel Kang, Joseph E. Gonzalez, Koushik Sen, Dawn Song, Ion Stoica, Matei Zaharia, and Marquita Ellis. 2026. Measuring agents in production. *arXiv:2512.04123*
- [77] Shishir G Patil, Huanzhi Mao, Fanjia Yan, Charlie Cheng-Jie Ji, Vishnu Suresh, Ion Stoica, and Joseph E Gonzalez. 2025. The Berkeley Function Calling Leaderboard (BFCL): From tool use to agentic evaluation of large language models. In *Forty-second International Conference on Machine Learning*.

- [78] R. Hugo Patterson and Stephen Manley. 2002. SnapMirror: File-system-based asynchronous mirroring for disaster recovery. In *Conference on File and Storage Technologies (FAST 02)*. USENIX Association, Monterey, CA. <https://www.usenix.org/conference/fast-02/snapmirror-file-system-based-asynchronous-mirroring-disaster-recovery>
- [79] Grant A Pignatiello, Richard J Martin, and Ronald L Hickman Jr. 2020. Decision fatigue: A conceptual analysis. *Journal of health psychology* 25, 1 (2020), 123–135.
- [80] David Quigley, Josef Sipek, Charles P Wright, and Erez Zadok. 2006. Unionfs: User-and community-oriented development of a unification filesystem. In *Proceedings of the 2006 Linux Symposium*, Vol. 2. 349–362.
- [81] Sean Quinlan and Sean Dorward. 2002. Venti: A new approach to archival data storage. In *Conference on File and Storage Technologies (FAST 02)*. USENIX Association, Monterey, CA. <https://www.usenix.org/conference/fast-02/venti-new-approach-archival-data-storage>
- [82] Sean Quinlan, Jim McKie, and Russ Cox. 2003. *Fossil, an archival file server*. Technical Report. Plan 9. Retrieved March 27, 2026 from <https://p9f.org/sys/doc/fossil.pdf>
- [83] Relynt. 2026. Designing approvals that do not kill automation | Relynt Blog. Retrieved March 28, 2026 from <https://www.relyntpolicy.com/blog/slack-approvals-human-in-the-loop>
- [84] Dennis M. Ritchie and Ken Thompson. 1974. The UNIX time-sharing system. *Commun. ACM* 17, 7 (July 1974), 365–375. <https://doi.org/10.1145/361011.361061>
- [85] Ohad Rodeh, Josef Bacik, and Chris Mason. 2013. BTRFS: The Linux B-tree filesystem. *ACM Transactions on Storage* 9, 3 (2013), 1–32.
- [86] Mickaël Salaün. 2017. Landlock LSM: Toward unprivileged sandboxing. Presentation at Linux Security Summit. Retrieved March 28, 2026 from https://static.sched.com/hosted_files/lss2017/56/2017-09-14_landlock-lss.pdf
- [87] Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems* 36 (2023), 68539–68551.
- [88] Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language agents with verbal reinforcement learning. In *Advances in Neural Information Processing Systems*, Vol. 36. 8634–8652.
- [89] Stanford Secure Computer Systems group. 2026. jai - easy containment for AI agents. Retrieved March 28, 2026 from <https://jai.scs.stanford.edu/>
- [90] Tensorlake Inc. 2026. Tensorlake – The Engine for Agentic Work. Retrieved March 31, 2026 from <https://www.tensorlake.ai/>
- [91] The Gemini Team. 2026. Gemini 3.1 Pro: Announcing our latest Gemini AI model. Retrieved April 01, 2026 from <https://blog.google/innovation-and-ai/models-and-research/gemini-models/gemini-3-1-pro/>
- [92] The Linux man-pages project. 2026. mount(8) - Linux manual page. Retrieved April 01, 2026 from <https://man7.org/linux/man-pages/man8/mount.8.html>
- [93] The Linux man-pages project. 2026. mount_namespaces(7) - Linux manual page. Retrieved March 28, 2026 from https://man7.org/linux/man-pages/man7/mount_namespaces.7.html
- [94] The Linux man-pages project. 2026. open(2) - Linux manual page. Retrieved March 28, 2026 from <https://man7.org/linux/man-pages/man2/open.2.html>
- [95] The MITRE Corporation. 2026. Multi-Factor Authentication Request Generation, Technique T1621 - Enterprise | MITRE A&TTCk®. Retrieved March 28, 2026 from <https://attack.mitre.org/versions/v17/techniques/T1621/>
- [96] Walter F Tichy. 1985. RCS—A system for version control. *Software: Practice and Experience* 15, 7 (1985), 637–654.
- [97] Turso. 2026. AgentFS with copy-on-write overlay filesystem. Retrieved March 31, 2026 from <https://turso.tech/blog/agentfs-overlay>
- [98] U.S. Department of Defense. 1985. *Department of Defense Trusted Computer System Evaluation Criteria*. Technical Report. National Computer Security Center.
- [99] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
- [100] Vercel Inc. 2026. Vercel Sandbox. Retrieved March 31, 2026 from <https://vercel.com/docs/vercel-sandbox>
- [101] Jos Visser. 2025. Coding in the time of AI. <https://josvisser.substack.com/p/coding-in-the-time-of-ai>.
- [102] VMS Software, Inc. 2020. *VSI OpenVMS User's Manual*. VMS Software, Inc., Bolton, MA, USA. https://vmssoftware.com/docs/VSI_USERS_MANUAL.pdf
- [103] Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, Wayne Xin Zhao, Zhewei Wei, and Jirong Wen. 2024. A survey on large language model based autonomous agents. *Frontiers of Computer Science* 18, 6 (2024), 186345.
- [104] Xingyao Wang, Boxuan Li, Yufan Song, Frank F. Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi Song, Bowen Li, Jaskirat Singh, Hoang H. Tran, Fuqiang Li, Ren Ma, Mingzhang Zheng, Bill Qian, Yanjun Shao, Niklas Muennighoff, Yizhe Zhang, Binyuan Hui, Junyang Lin, Robert Brennan, Hao Peng, Heng Ji, and Graham Neubig. 2025. OpenHands: An open platform for AI software developers as generalist agents. In *The Thirteenth International Conference on Learning Representations*.
- [105] Simon Willison. 2022. Prompt injection attacks against GPT-3. Retrieved March 28, 2026 from <https://simonwillison.net/2022/Sep/12/prompt-injection/>
- [106] Simon Willison. 2025. Living dangerously with Claude. <https://simonwillison.net/2025/Oct/22/living-dangerously-with-claude/>.
- [107] Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh Jing Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, Yitao Liu, Yiheng Xu, Shuyan Zhou, Silvio Savarese, Caiming Xiong, Victor Zhong, and Tao Yu. 2024. OSWorld: Benchmarking multimodal agents for open-ended tasks in real computer environments. *Advances in Neural Information Processing Systems* 37 (2024), 52040–52094.
- [108] Hui Yang, Sifu Yue, and Yunzhong He. 2023. *Auto-GPT for online decision making: Benchmarks and additional opinions*. arXiv:2306.02224
- [109] John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan, and Ofir Press. 2024. SWE-Agent: Agent-computer interfaces enable automated software engineering. *Advances in Neural Information Processing Systems* 37 (2024), 50528–50652.
- [110] John Yang, Akshara Prabhakar, Karthik Narasimhan, and Shunyu Yao. 2023. Intercode: Standardizing and benchmarking interactive coding with execution feedback. *Advances in Neural Information Processing Systems* 36 (2023), 23826–23854.
- [111] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. 2022. ReAct: Synergizing reasoning and acting in language models. In *The eleventh international conference on learning representations*.
- [112] Miao Yu, Fanci Meng, Xinyun Zhou, Shilong Wang, Junyuan Mao, Linsey Pan, Tianlong Chen, Kun Wang, Xinfeng Li, Yongfeng Zhang, et al. 2025. A survey on trustworthy llm agents: Threats and countermeasures. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 2*. 6216–6226.

- [113] Erez Zadok and Jason Nieh. 2000. FiST: A Language for Stackable File Systems. USENIX Annual Technical Conference (USENIX ATC '00).
- [114] Hanrong Zhang, Jingyuan Huang, Kai Mei, Yifei Yao, Zhenting Wang, Chenlu Zhan, Hongwei Wang, and Yongfeng Zhang. 2025. Agent Security Bench (ASB): Formalizing and Benchmarking Attacks and Defenses in LLM-based Agents. In *The Thirteenth International Conference on Learning Representations*.
- [115] Zhexin Zhang, Shiyao Cui, Yida Lu, Jingzhuo Zhou, Junxiao Yang, Hongning Wang, and Minlie Huang. 2025. Agent-SafetyBench: Evaluating the Safety of LLM Agents. arXiv:2412.14470 [cs.CL] <https://arxiv.org/abs/2412.14470>
- [116] Nannan Zhao, Vasily Tarasov, Hadeel Albahar, Ali Anwar, Lukas Rupprecht, Dimitrios Skourtis, Arnab K Paul, Keren Chen, and Ali R Butt. 2020. Large-scale analysis of docker images and performance implications for container storage systems. *IEEE Transactions on Parallel and Distributed Systems* 32, 4 (2020), 918–930.
- [117] Xinxue Zhu, Jiacong Wu, Xiaoyu Zhang, Tianlin Li, Yanzhou Mu, Juan Zhai, Chao Shen, and Yang Liu. 2026. *An empirical study of bugs in modern LLM agent frameworks*. arXiv:2602.21806
- [R1] Cybersecurity analysis of Antigravity agent wiping entire D: drive via rmdir /s /q when user asked to clear project cache. <https://hackernoob.tips/when-ai-agents-go-rogue-google-antigravities-catastrophic-drive-deletion-exposes-critical-risks-in-agentic-development-tools/>.
- [R2] Antigravity’s built-in edit tool systematically deletes and corrupts file content during normal coding tasks, confirmed by multiple users as a tool bug. https://www.reddit.com/r/Bard/comments/1p3htvu/antigravity_just_deletingcorrupting_files.
- [R3] Indirect prompt injection via poisoned web page causes Gemini to bypass gitignore restriction using cat, read .env credentials, and exfiltrate via browser subagent to attacker-controlled webhook.site. <https://www.promptarmor.com/resources/google-antigravity-exfiltrates-data>.
- [R4] Antigravity agent wiped entire D: drive via rmdir /s /q due to quote-handling path resolution bug; SafeToAutoRun bypassed approval. https://www.reddit.com/r/google_antigravity/comments/1p82or6/google_antigravity_just_deleted_the_contents_of/.
- [R5] /agents command creates agent files relative to CWD instead of project root, placing them in wrong directory. <https://github.com/anthropics/claude-code/issues/7560>.
- [R6] Allow all edits" session setting ignored; Claude Code prompts for approval on every individual file edit despite user opting in. <https://github.com/anthropics/claude-code/issues/9348>.
- [R7] Bash allowlist pattern Bash(find:*) fails to match piped commands, causing repeated permission prompts despite saved approval. <https://github.com/anthropics/claude-code/issues/34811>.
- [R8] Claude Code ignores explicit user instructions and modifies files without direction while Accept Edits auto-approves changes. <https://github.com/anthropics/claude-code/issues/8026>.
- [R9] Subagent auto-yes to stdin created 851GB task output file, silently filling disk. <https://github.com/anthropics/claude-code/issues/35047>.
- [R10] Background subagents silently auto-deny tools instead of prompting for permissions before launch, with inconsistent propagation across sessions. <https://github.com/anthropics/claude-code/issues/34095>.
- [R11] Bash permission bypass via command chaining (&&, ;, |) defeats prefix-matching allowlist. <https://github.com/anthropics/claude-code/issues/4956>.
- [R12] Claude Code ran broad find / searches across home directory instead of checking obvious config path, part of reported model regression. <https://github.com/anthropics/claude-code/issues/33908>.
- [R13] /btw command fails to render permission prompts, freezing session with no way to respond. <https://github.com/anthropics/claude-code/issues/33313>.
- [R14] Built-in tools (Read, Grep, Bash, Glob) execute without approval prompt despite not being in allowedTools allowlist. <https://github.com/anthropics/claude-code/issues/34569>.
- [R15] defaultMode bypassPermissions setting has no effect; every tool call still triggers permission prompt, growing allowlist indefinitely. <https://github.com/anthropics/claude-code/issues/34923>.
- [R16] Claude Code autonomously bypassed three security layers (path denylist, bubblewrap sandbox, BPF LSM kernel enforcement) using Linux primitives to complete its task. <https://cybercorsairs.com/claude-code-broke-out-of-three-security-layers-on-its-own/>.
- [R17] Claude Code deleted all unstaged files when asked to restructure commits, pulled old files as fake recovery, then deleted everything again. https://www.reddit.com/r/ClaudeAI/comments/1q0124f/terrible_experience_with_claude_code_lost_my_work.
- [R18] Claude Code ignored CLAUDE.md "never commit" rule, committed without permission, then fabricated a justification when confronted. <https://github.com/anthropics/claude-code/issues/34774>.
- [R19] Claude Code repeatedly commits and pushes without permission despite explicit CLAUDE.md and memory rules, recurring after corrections in same session. <https://github.com/anthropics/claude-code/issues/34451>.
- [R20] Claude Code silently drops CLAUDE.md rules and session constraints as context window fills in long sessions, leading to unauthorized file edits and hallucinated state. <https://github.com/anthropics/claude-code/issues/32659>.
- [R21] Claude Code context decay during debugging session caused repeated instruction violations, reversed fixes, and significant unwanted file changes. <https://github.com/anthropics/claude-code/issues/8251>.
- [R22] Cowork copied 0-byte iCloud stubs via cp -a then rm -rf'd originals, destroying 110 sensitive legal documents. <https://github.com/anthropics/claude-code/issues/32637>.
- [R23] Cowork VM sandbox fails to start on Windows due to sandbox-helper unmount errors, leaving sandboxed mode unusable. <https://github.com/anthropics/claude-code/issues/25419>.
- [R24] Claude Code reads files in denied credentials folder via alternative tool path, bypassing Read() deny rule in settings.json. <https://github.com/anthropics/claude-code/issues/4768>.
- [R25] Claude Code autonomously deleted specification file it was working from, then claimed no knowledge of it. <https://github.com/anthropics/claude-code/issues/7787>.
- [R26] Claude Opus repeatedly deletes or weakens unit tests and writes TODO stubs despite explicit CLAUDE.md TDD rules, then reports features as complete. https://www.reddit.com/r/ClaudeAI/comments/1lfirvk/any_tips_on_how_to_get_claude_to_stop_cheating_on/.
- [R27] Bash deny rules only block relative paths; absolute paths bypass denylist completely, exposing entire filesystem. <https://github.com/anthropics/claude-code/issues/11662>.
- [R28] Bash deny rules bypassed via indirect file access commands (find, grep) that don't contain the denied filename pattern. <https://github.com/anthropics/claude-code/issues/6036>.
- [R29] Claude Desktop code mode lacks explain option on permission prompts, encouraging rubber-stamp approvals. <https://github.com/anthropics/claude-code/issues/34205>.
- [R30] Claude displayed full contents of /.netrc, /.npmrc, and /.cargo/credentials.toml in conversation, exposing 8+ API tokens and passwords. <https://github.com/anthropics/claude-code/issues/34819>.
- [R31] Windows drive letter change causes Permission Resolver to create .claude/ directories outside workspace and corrupt settings. <https://github.com/anthropics/claude-code/issues/34866>.

- [R32] Command injection via echo command bypasses allowlist regex and permission prompt, enabling arbitrary execution (CVE-2025-54795). <https://nvd.nist.gov/vuln/detail/CVE-2025-54795>.
- [R33] Claude Code repeatedly edited infrastructure files despite 30+ pages of CLAUDE.md rules and a custom skill requiring explicit permission first. <https://github.com/anthropics/claude-code/issues/10726>.
- [R34] Claude Code repeatedly removes error-handling code instead of fixing bugs, ignoring explicit user instructions not to. <https://github.com/anthropics/claude-code/issues/8910>.
- [R35] Claude reported 15/15 bugs fixed while fabricating data and skipping source verification, making the document worse than the original. <https://github.com/anthropics/claude-code/issues/27399>.
- [R36] Claude Code ran find across entire home directory instead of project directory to locate H2 database file for deletion. <https://github.com/anthropics/claude-code/issues/30125>.
- [R37] Claude Code permanently deleted a folder without user confirmation on Windows. <https://github.com/anthropics/claude-code/issues/13476>.
- [R38] Claude Code ran git commit -amend without confirmation despite explicit CLAUDE.md rules listing it as a destructive operation requiring user approval. <https://github.com/anthropics/claude-code/issues/33391>.
- [R39] Claude Code ignores CLAUDE.md git restrictions after context compaction, cherry-picks from wrong branches and rewrites git history without permission. https://www.reddit.com/r/ClaudeAI/comments/1p889x7/anyone_successfully_prevent_claude_from_running/.
- [R40] Claude executed unauthorized git commit and checkouts during test request, massively modifying server file and destroying day's uncommitted work. <https://github.com/anthropics/claude-code/issues/34784>.
- [R41] Claude Code ran git commit -amend and force-push without permission, violating memory file rules and bypassing permission prompt. <https://github.com/anthropics/claude-code/issues/25472>.
- [R42] Claude Code executed git reset -hard and git checkout while falsely assuring user operations were safe, destroying hours of development work. <https://github.com/anthropics/claude-code/issues/7232>.
- [R43] Claude Code ran git reset -hard origin/main in main working directory during multi-instance workflow, destroying 2 days of uncommitted Pulumi infrastructure work. <https://github.com/anthropics/claude-code/issues/33850>.
- [R44] Claude Code ran git reset -hard without checking for or warning about untracked/staged files, destroying 3 work files. <https://github.com/anthropics/claude-code/issues/34746>.
- [R45] Claude Code attempted git rm -rf on non-existent path and hallucinated dangerouslyDisableSandbox tool parameter. <https://github.com/anthropics/claude-code/issues/31705>.
- [R46] Claude hallucinated incorrect Punycode for Cyrillic domains, writing corrupted config files to production including substitution of a real third-party domain. <https://github.com/anthropics/claude-code/issues/32798>.
- [R47] Claude Code hand-rolled Django migration and applied it without permission, violating explicit CLAUDE.md rules, then initially lied about rules not existing. <https://github.com/anthropics/claude-code/issues/8059>.
- [R48] Bash permission glob matcher treats. <https://github.com/anthropics/claude-code/issues/34379>.
- [R49] Permission checker falsely flags literal backticks inside single-quoted heredoc as command substitution. <https://github.com/anthropics/claude-code/issues/35183>.
- [R50] Security heuristics override explicit allowlist entries, forcing approval prompts for pre-authorized commands. <https://github.com/anthropics/claude-code/issues/34106>.
- [R51] Claude Code ran git stash drop despite explicit CLAUDE.md prohibition, permanently deleting 3 days of stashed work. <https://github.com/anthropics/claude-code/issues/22638>.
- [R52] Claude Code's persistent shell lost PWD state after user interrupt, causing rm -rf * to execute in project root instead of Builds/Xcode subdirectory. <https://github.com/anthropics/claude-code/issues/17410>.
- [R53] Claude Code violated CLAUDE.md safeguard rules during 2-hour session, overwrote custom LoRA model file with wrong data, destroyed working ComfyUI workflow. <https://github.com/anthropics/claude-code/issues/34922>.
- [R54] Claude deleted files and removed document content without permission, repeatedly misinterpreting ambiguous feedback as deletion instructions. <https://github.com/anthropics/claude-code/issues/34492>.
- [R55] Notification hook for permission prompts omits tool name and command details, preventing informed remote triage. <https://github.com/anthropics/claude-code/issues/32952>.
- [R56] Cmd+Enter to approve permission prompt also submits partially-typed message, sending unintended instructions to agent. <https://github.com/anthropics/claude-code/issues/32630>.
- [R57] Permission dialog auto-scrolls past reasoning context, preventing informed approval of tool calls. <https://github.com/anthropics/claude-code/issues/34354>.
- [R58] Feature request for human-readable warning labels on destructive commands in permission prompts. <https://github.com/anthropics/claude-code/issues/30505>.
- [R59] Plan Mode failed to block Bash execution; destructive commands (pkill, kill -9) ran despite read-only mode being active. <https://github.com/anthropics/claude-code/issues/32768>.
- [R60] Claude Code executed Edit tool to modify source file despite Plan Mode being active, bypassing the constraint that blocks non-readonly operations. <https://github.com/anthropics/claude-code/issues/33037>.
- [R61] Plan mode has no technical tool restriction — relies purely on system prompt injection, allowing edits and shell commands despite "read-only" designation. https://www.reddit.com/r/ClaudeAI/comments/1ou9x6i/unsafe_plan_mode_does_not_prevent_claude_from/.
- [R62] Read() permission allowlist path prefix matching rejects absolute paths due to unintuitive double-slash syntax requirement. <https://github.com/anthropics/claude-code/issues/35118>.
- [R63] Claude Code Read() deny rules in settings.json silently ignored, allowing agent to read restricted files including secrets and private data. <https://github.com/anthropics/claude-code/issues/3501>.
- [R64] Claude Code rebased and force-pushed another contributor's PR branch, rewriting commit history and violating its own instructions. <https://github.com/anthropics/claude-code/issues/32476>.
- [R65] Claude Code deleted years of teaching resources during folder reorganization instead of archiving uncategorizable files. <https://github.com/anthropics/claude-code/issues/12851>.
- [R66] Claude Code repeatedly deleted fire effect code via Edit tool while user was debugging, then attempted same deletion immediately after being explicitly told not to. <https://github.com/anthropics/claude-code/issues/7645>.
- [R67] Claude repeatedly used rm -rf to delete config files during pair programming; user approved each time without reading prompts. https://www.reddit.com/r/ClaudeAI/comments/1nrrmv3/til_ai_keeps_using_rm_rf_on_important_files/.
- [R68] Claude CLI appended / to rm -rf command while cleaning up packages in old repo, deleting entire Mac home directory. https://www.reddit.com/r/ClaudeAI/comments/1pgxckk/claude_cli_deleted_my_entire_home_directory_wiped/.
- [R69] Claude deleted 11h of YOLO inference output via rm -rf without permission, then restarted the job while user was asleep. <https://github.com/anthropics/claude-code/issues/34106>.

- [//github.com/anthropics/claude-code/issues/32938](https://github.com/anthropics/claude-code/issues/32938).
- [R70] Claude deleted non-empty directory with `rm -rf` after user asked to remove only if empty; files were unique, not duplicates. <https://github.com/anthropics/claude-code/issues/35093>.
- [R71] Claude ran `rm -rf` on entire local folder tree without confirmation, destroying months of production code in sibling directories. <https://github.com/anthropics/claude-code/issues/30816>.
- [R72] Claude Code `rm -rf` with tilde expansion deleted user's entire home directory during project cleanup. <https://byteiota.com/claude-codes-rm-rf-bug-deleted-my-home-directory/>.
- [R73] Claude Code executed `rm -rf` from root, deleting all user-owned files in home directory; exact command not captured in logs. <https://github.com/anthropics/claude-code/issues/10077>.
- [R74] Claude Code ran `rm -rf` from repo root instead of build subdirectory with sandbox disabled, destroying `.git` and 76+ unpushed commits; then silently discarded all code during rebase recovery. <https://github.com/anthropics/claude-code/issues/34514>.
- [R75] Claude Code executed `rm` on wrong file instead of displaying the command as requested, deleting an unrelated media file with Accept Edits on. <https://github.com/anthropics/claude-code/issues/24854>.
- [R76] Max Plan subscriber documents five systematic CLAUDE.md rule-violation patterns despite 24 PreToolUse hooks and 400-line rules file; model reads, acknowledges, then violates rules. <https://github.com/anthropics/claude-code/issues/34358>.
- [R77] Claude Code silently bypasses sandbox via `dangerouslyDisableSandbox` flag when Bash is auto-approved, no permission prompt shown. <https://github.com/anthropics/claude-code/issues/34315>.
- [R78] `SandboxExcludedCommands` does not exempt SSH child process, blocking git push and forcing users to disable sandbox entirely. <https://github.com/anthropics/claude-code/issues/33300>.
- [R79] Sandbox creates empty directories in CWD from command arguments when command is not found. <https://github.com/anthropics/claude-code/issues/33056>.
- [R80] Claude Code violated its own saved safety rules, executing SSH commands and installing software on remote servers without required approval. <https://github.com/anthropics/claude-code/issues/34828>.
- [R81] Feature request for scoped auto-approve (time-limited, task-scoped, step-count) as middle ground between per-action prompts and full permissionless mode. <https://github.com/anthropics/claude-code/issues/33645>.
- [R82] Claude modified generation script to auto-delete 50 curated vocal sample files, permanently destroying 7 sessions of user curation work. <https://github.com/anthropics/claude-code/issues/30988>.
- [R83] Claude Code bypasses `deny` config and `CLAUDE.md` rules to read `.env` files and secrets via alternative paths or indirect commands. https://www.reddit.com/r/ClaudeAI/comments/1nfvh46/i_got_tired_of_claude_code_reading_secrets_and/.
- [R84] Claude Code `sed` command parsing flaw bypasses read-only validation, enabling arbitrary file writes via `sed w/r` subcommands (CVE-2025-64755). <https://nvd.nist.gov/vuln/detail/CVE-2025-64755>.
- [R85] Claude Code JSON serializer corrupts `settings.local.json` when saving bash commands with regex quote patterns, destroying all saved permissions. <https://github.com/anthropics/claude-code/issues/33650>.
- [R86] `settings.local.json` allowlist rules silently ignored due to path format bugs; agent re-prompts for pre-approved commands, pushing users toward `-dangerously-skip-permissions`. <https://github.com/anthropics/claude-code/issues/6850>.
- [R87] Claude Code autonomously ran `git reset --hard` on startup twice, destroying unpushed commits and uncommitted work; fabricated a safety hook that never existed. <https://github.com/anthropics/claude-code/issues/34327>.
- [R88] Sub-agent ran `rm -rf` on its own worktree directory, corrupting session CWD and forcing sign-out. <https://github.com/anthropics/claude-code/issues/34474>.
- [R89] Task output collector reads its own file in infinite loop, writing 696GB and exhausting disk in 23 minutes. <https://github.com/anthropics/claude-code/issues/34783>.
- [R90] Claude created directory named then ran `rm -rf *`, shell expansion wiped user's home directory. <https://github.com/anthropics/claude-code/issues/12637>.
- [R91] Claude Code Bash tool creates `/tmp/claude*-cwd` files to track working directory but never deletes them, causing 500+ files/day accumulation. <https://github.com/anthropics/claude-code/issues/8856>.
- [R92] Claude Code created 1,489 lines of unwanted documentation files violating explicit `CLAUDE.md` rule against proactive doc creation. <https://github.com/anthropics/claude-code/issues/20109>.
- [R93] Claude Code VSCode extension silently auto-approves all Bash commands despite default permission mode requiring approval. <https://github.com/anthropics/claude-code/issues/34463>.
- [R94] VSCode extension prepends `cd &&` to every Bash command, breaking allowlist auto-approval matching. <https://github.com/anthropics/claude-code/issues/33976>.
- [R95] Claude Code VSCode/Positron extension skips workspace trust prompt, allowing file reads in arbitrary directories without confirmation. <https://github.com/anthropics/claude-code/issues/34913>.
- [R96] Sonnet generates malformed wildcard glob patterns in `settings.json`, corrupting permission config and crashing VS Code. <https://github.com/anthropics/claude-code/issues/33746>.
- [R97] Claude Code creates undeletable "nul" files on Windows by using Unix-style `/dev/null` redirection, corrupting git repos and breaking IDEs. <https://github.com/anthropics/claude-code/issues/4928>.
- [R98] Claude Code followed NTFS junctions when removing `pnpm monorepo` worktrees, permanently deleting user's Windows profile folders. <https://github.com/anthropics/claude-code/issues/29249>.
- [R99] Claude Code Write tool and Bash `cp` executed without permission prompt, silently bypassing configured allowlist in `settings.local.json`. <https://github.com/anthropics/claude-code/issues/34583>.
- [R100] Claude Code used Write (full overwrite) instead of Edit, silently dropping table of contents section from backup file. <https://github.com/anthropics/claude-code/issues/27137>.
- [R101] Yarn 2+ plugin auto-execution bypasses directory trust dialog, enabling arbitrary code execution before user accepts risks (CVE-2025-59828). <https://nvd.nist.gov/vuln/detail/CVE-2025-59828>.
- [R102] Cline v3.20.0 deletes recently edited files when user cancels a task or switches from Act to Plan mode. <https://github.com/cline/cline/issues/5124>.
- [R103] `.clineignore` still sends all ignored file paths to LLM context with lock emoji instead of omitting them, wasting 30k+ tokens and leaking path names. <https://github.com/cline/cline/issues/9554>.
- [R104] Corrupted checkpoint shadow repo disables checkpoint revert system globally across all workspaces with no recovery path. <https://github.com/cline/cline/issues/9631>.
- [R105] Cline v3.20.0 deletes last edited file after task completion due to missing guard in `DiffViewProvider.revertChanges()`. <https://github.com/cline/cline/issues/5120>.
- [R106] Cline repeatedly overwrites documentation files with truncated/placeholder content when asked to update them, despite explicit instructions not to. <https://github.com/cline/cline/issues/711>.
- [R107] Cline vulnerable to data exfiltration via markdown image rendering; prompt injection reads `.env` and exfiltrates secrets through attacker-controlled image URL. <https://embracethered.com/blog/posts/2025/cline-vulnerable-to-data-exfiltration/>.
- [R108] Cline used `write_to_file` to overwrite entire `.env` with one variable, destroying all existing credentials. <https://github.com/cline/cline/issues/8422>.

- [R109] Cline checkpoint system renames `.git` to `.git_disabled` during editing, leaving project without git history when interrupted. <https://github.com/cline/cline/issues/8273>.
- [R110] Feature request to block AI agents from bypassing git hooks with `-no-verify` flag. <https://github.com/cline/cline/issues/8526>.
- [R111] Cline CLI deletes file when user interrupts mid-stream edit generation, bypassing required approval. <https://github.com/cline/cline/issues/9165>.
- [R112] Cline global MCP auto-approve toggle overrides per-tool approval settings, preventing granular permission control. <https://github.com/cline/cline/issues/9357>.
- [R113] Cline MCP installation corrupted then deleted `cline_mcp_settings.json`, losing all MCP server configs and bearer tokens. <https://github.com/cline/cline/issues/9663>.
- [R114] `replace_in_file` tool empties file content when SEARCH block fails to match, falsely claims reversion. <https://github.com/cline/cline/issues/9555>.
- [R115] Cline ran `rm -rf /` without permission while user asked to clean up a merge conflict, deleting entire home directory. <https://github.com/cline/cline/issues/7572>.
- [R116] Cline creates empty file named "f" and writes to it after search patterns fail to match during editing. <https://github.com/cline/cline/issues/9821>.
- [R117] Cline update randomly deletes opened or recently edited files; multiple users report data loss on v3.20.0. https://www.reddit.com/r/CLine/comments/1m79n0s/cline_deleting_files.
- [R118] Codex VS Code extension ignores "Allow every time" for file edits, re-prompting on every file despite user approval. <https://github.com/openai/codex/issues/3157>.
- [R119] Codex VS Code extension ignores "Allow every time" button on Windows, re-prompting for every file edit despite repeated approvals. <https://github.com/openai/codex/issues/3325>.
- [R120] Codex Windows `apply_patch` fails for nested files under `src/**` due to sandbox path-handling bug; shell writes succeed in same paths. <https://github.com/openai/codex/issues/14675>.
- [R121] Codex Windows app `apply_patch` silently fails in sandboxed default permission mode, forcing users to disable sandbox via Full Access. <https://github.com/openai/codex/issues/13574>.
- [R122] Codex approval prompt only displays first command before `&&`, hiding chained commands from user review. <https://community.openai.com/t/approval-only-showing-start-of-command/1372406>.
- [R123] Codex on Windows requires manual approval for every shell command, making it unusable. <https://github.com/openai/codex/issues/2860>.
- [R124] Codex approval prompt froze and became unresponsive after several approved commands; restarting app deleted session reasoning. <https://github.com/openai/codex/issues/10760>.
- [R125] Codex in auto-approval mode reads files outside working directory including `/.ssh`, despite documentation claiming approval is required. <https://github.com/openai/codex/issues/5474>.
- [R126] Codex `bwrap` sandbox fails when home directory is symlinked to NFS mount, falling back to permission prompts. <https://github.com/openai/codex/issues/14832>.
- [R127] Codex `bwrap` sandbox fails to bind symlinked `TMPDIR`, breaking sandbox containment and triggering unnecessary permission elevation requests. <https://github.com/openai/codex/issues/14672>.
- [R128] Users request command-specific allowlist to avoid repeated permission prompts for safe commands like tests and builds. <https://github.com/openai/codex/issues/3085>.
- [R129] Codex autonomously deleted a 2200-line code file it deemed unnecessary; user recovered from offline backups. https://www.linkedin.com/posts/alexander-purchase_today-i-had-codex-delete-one-of-my-files-activity-7385244768156446720-yMGk.
- [R130] Codex deleted uncommitted files twice without permission, then fabricated claims about restoring them from git. <https://github.com/openai/codex/issues/4969>.
- [R131] Codex deletes entire files instead of making localized edits, admits "I accidentally deleted the file." <https://github.com/openai/codex/issues/6999>.
- [R132] Codex hardcoded denylist blocks `rm` and `git reset` even in full-access mode; agent bypasses via Python `shutil` and `apply_patch`. <https://github.com/openai/codex/issues/5128>.
- [R133] Codex runs docker commands without confirmation after silent security model change, granting effective root on host. <https://github.com/openai/codex/issues/14477>.
- [R134] Codex Windows elevated sandbox setup repeatedly fails with exit code 1, blocking all reads and writes across six CLI versions. <https://github.com/openai/codex/issues/14409>.
- [R135] Codex has no ability to hide sensitive files like `.env` from agent; instructions and `.gitignore` are not respected. <https://github.com/openai/codex/issues/85>.
- [R136] Codex erased contents of 20,000-line file and replaced with a snippet when creating PR. <https://community.openai.com/t/codex-deletes-contents-of-large-files/1357223>.
- [R137] User requests fine-grained permission controls between overly restrictive default and overly permissive full-access modes. <https://github.com/openai/codex/issues/14399>.
- [R138] Codex sandbox user lacks permissions to `git fsmonitor` IPC socket, causing git errors and agent stumbling. <https://github.com/openai/codex/issues/14372>.
- [R139] Codex VS Code extension in full access mode requires approval for every change; users report 20-30 prompts per task including read-only chat mode. <https://community.openai.com/t/codex-vscode-extension-agent-full-access-always-asks-for-approval/1355908>.
- [R140] Codex CLI and IDE in full access mode still prompt for file edit approval on Windows despite all configuration overrides. <https://community.openai.com/t/codex-cli-and-ide-prompting-for-approval-to-edit-files/1354993>.
- [R141] Codex CLI in full-auto mode deleted personal photos on `D:\{Photo}` for two days, ignoring `AGENTS.MD` rules restricting deletion to project directory. <https://github.com/openai/codex/issues/14487>.
- [R142] Codex ran `git reset` despite `agents.md` prohibiting git commands, rolling back 6 hours of uncommitted refactoring work. https://www.reddit.com/r/codex/comments/1onqxfq/codex_finally_fed_me_git_issue_lessons_learned/.
- [R143] Codex CLI with GPT-5.2 runs `git restore` to discard user's unrelated uncommitted changes for a "clean" PR. <https://github.com/openai/codex/issues/8213>.
- [R144] Codex executed `git restore .` without confirmation, wiping hours of uncommitted work from dirty working tree. https://www.reddit.com/r/codex/comments/1pt3vcn/be_careful_with_codex/.
- [R145] Codex sandbox blocks GPU access, forcing users to disable all sandbox protections to run ML workloads. <https://github.com/openai/codex/issues/3141>.
- [R146] Codex CLI panics with Landlock sandbox error on kernels lacking Landlock support, forcing users to disable sandboxing entirely. <https://github.com/openai/codex/issues/2267>.
- [R147] Codex has no mechanism to exclude sensitive files from being read or sent to the model, exposing secrets like `.env` and SSH keys. <https://github.com/openai/codex/issues/2847>.
- [R148] Codex overwrote research document with summaries when asked to read folder; undo feature non-functional. https://www.reddit.com/r/codex/comments/1pe1f69/wow_undo_not_working/.
- [R149] Keyring failure silently falls back to plaintext credential file storage with potentially world-readable permissions, leaving stale secrets on disk. <https://github.com/openai/codex/issues/14704>.

- [R150] Codex read .env file via built-in Read tool without permission while exploring project to fix e2e test. https://www.reddit.com/r/codex/comments/1o23ok8/any_way_to_prevent_it_reading_env/.
- [R151] Codex CLI intentionally removed read-only approval mode, eliminating per-edit permission prompts and forcing auto-apply as default. <https://github.com/openai/codex/issues/11915>.
- [R152] Codex reads files outside working directory without permission despite sandbox documentation claiming reads are restricted to cwd. <https://github.com/openai/codex/issues/5237>.
- [R153] Codex App "don't ask again" permission persistence broken due to command parser failing on environment variable prefixes. <https://github.com/openai/codex/issues/11298>.
- [R154] Codex repeatedly ran `rm -rf` deleting entire project directory four times in one session despite explicit instructions not to delete. https://www.reddit.com/r/codex/comments/1nmc05l/anyone_else_getting_issues_with_codex_nuking/.
- [R155] Codex reverts user's manual code changes via `git reset/checkout/rm`, destroying uncommitted work across multiple users. <https://github.com/openai/codex/issues/1736>.
- [R156] Codex review mode read-only sandbox blocks `git` commands on macOS due to `xcrun` temp file writes; `-yolo` flag ignored; sandbox reverted entirely. <https://github.com/openai/codex/issues/7815>.
- [R157] Codex Windows sandbox triggers `0xc0000022` access-denied error on every command, forcing users to disable sandbox or use full-access mode. <https://github.com/openai/codex/issues/14448>.
- [R158] Codex Windows sandbox sets incorrect ACLs on newly-created folders, causing `apply_patch` to fail on second write to same directory. <https://github.com/openai/codex/issues/14585>.
- [R159] Codex re-prompts for sandbox permission on `xcodebuild` after user selected "don't ask again"; approval rules not persisting. <https://github.com/openai/codex/issues/10321>.
- [R160] Codex Windows sandbox auto-applies NTFS deny ACE to workspace paths, causing write failures and requiring manual ACL cleanup. <https://github.com/openai/codex/issues/14006>.
- [R161] Codex VS Code sandbox mounts writable devcontainer workspace as read-only due to duplicate `ro-rw` mount namespace bug. <https://github.com/openai/codex/issues/14794>.
- [R162] Codex sandbox allow rules fail when commands are wrapped in shell scripts; MCP servers bypass sandbox entirely. <https://github.com/openai/codex/issues/12283>.
- [R163] Codex bubblewrap sandbox fails when `./codex` is a symlink to another partition, breaking `apply_patch`. <https://github.com/openai/codex/issues/14694>.
- [R164] Codex sandbox on macOS blocks test runners (`vitest`, `jest`), pushing users to disable all safety via `-dangerously-bypass-approvals-and-sandbox`. <https://github.com/openai/codex/issues/1532>.
- [R165] Codex truncated file from wrong location losing 3000 lines; VSCode revert button failed, requiring manual recovery from conversation history. <https://github.com/openai/codex/issues/7291>.
- [R166] Codex on Windows double-encodes UTF-8 as CP1252, corrupting emojis and non-ASCII characters in files; `apply_patch` failures cause file deletion and recreation. <https://github.com/openai/codex/issues/4013>.
- [R167] Codex VS Code extension asks for approval on every single file change despite being set to Agent full access mode. https://www.reddit.com/r/OpenaiCodex/comments/1n82y3d/codex_vs_code_extension_is_always_asking_for/.
- [R168] Codex on Windows ignores "Allow for this session" approval because command vector mismatch between displayed bash command and executed PowerShell-wrapped command prevents approval caching. <https://github.com/openai/codex/issues/4212>.
- [R169] Codex Windows sandbox fails with `CreateProcessWithLogonW` error 5 due to localized group names, forcing excessive approval prompts and opaque PowerShell edits. <https://github.com/openai/codex/issues/9062>.
- [R170] Codex Windows sandbox ignores "Allow this session" permission, re-prompting for every edit. <https://github.com/openai/codex/issues/7811>.
- [R171] Codex App for Windows deleted 370+ GB of user files outside the project directory while in Full Access mode; at least 8 additional users reported identical systemic drive-wiping behavior. <https://community.openai.com/t/critical-data-loss-issue-in-codex-app-for-windows-agent-executed-file-deletion-outside-project-directory/1375894>.
- [R172] Codex on Windows ignores `-ask-for-approval` never and `-dangerously-bypass-approvals-and-sandbox`, requiring manual approval for every file write. <https://github.com/openai/codex/issues/2350>.
- [R173] Codex App for Windows deleted Documents, Adobe, Dropbox, and Desktop files after encountering sandbox access errors. <https://community.openai.com/t/codex-for-windows-deleted-a-huge-amount-of-my-drive/1376684>.
- [R174] Codex VS Code extension on Windows ignores full-access mode, prompts for approval on every action. <https://github.com/openai/codex/issues/2828>.
- [R175] Codex agent `rmdir` command mis-parsed through multi-shell chain on Windows, deleting entire workspace contents instead of target folder. <https://community.openai.com/t/potential-destructive-command-mis-parsing-on-windows-agent-cleanup-via-cmd-c-may-delete-workspace-content-instead-of-target-folder/1376026>.
- [R176] Codex Windows sandbox setup fails during refresh, forcing users to bypass sandbox or grant full access permissions. <https://github.com/openai/codex/issues/10601>.
- [R177] Codex replaces 30,000 lines of code with a single requested addition, denies doing it in plan mode. https://www.reddit.com/r/codex/comments/1olvmxn/out_of_nowhere_codex_just_deletes_my_entire_code.
- [R178] Codex CLI on WSL2 intermittently re-enters sandbox mid-session despite `-dangerously-bypass-approvals-and-sandbox`, breaking `sudo/Docker/UNIX` sockets. <https://github.com/openai/codex/issues/5084>.
- [R179] Prompt injection exploited markdown image rendering in Continue VS Code extension to exfiltrate .env secrets and chat history to attacker-controlled server. <https://versprite.com/blog/data-exfiltration-via-image-rendering-fixed-in-continue/>.
- [R180] Copilot Agent repeatedly empties .py file content when attempting to apply code changes, then fails to recover it. <https://github.com/microsoft/vscode-copilot-release/issues/14021>.
- [R181] Copilot agent mode ignores file exclusion settings, reading sensitive files via built-in tools despite configured restrictions. <https://github.com/microsoft/vscode-copilot-release/issues/12711>.
- [R182] Copilot Agent Mode recreated deleted files and rolled back entire codebase to previous version on session reopen, causing hours of lost work. <https://github.com/microsoft/vscode/issues/264091>.
- [R183] Copilot Agent mode truncates file edits, destroying end of files, then loops failing to repair the truncation. <https://github.com/microsoft/vscode-copilot-release/issues/7038>.
- [R184] GitHub Copilot Chat exfiltrated chat context data via rendered markdown images after following prompt injection instructions hidden in source code. <https://embracethered.com/blog/posts/2024/github-copilot-chat-prompt-injection-data-exfiltration/>.
- [R185] User asked Copilot to delete files it created; agent deleted all project code instead. <https://github.com/microsoft/vscode-copilot-release/issues/13722>.
- [R186] Copilot agent deleted user's core project code then reported task as successfully completed. <https://github.com/microsoft/vscode-copilot-release/issues/14020>.

- [R187] Copilot deleted all user code without adding any new content during editing session. <https://github.com/microsoft/vscode-copilot-release/issues/14112>.
- [R188] Copilot Chat repeatedly deleted entire codebase when user requested edits. <https://github.com/microsoft/vscode-copilot-release/issues/14016>.
- [R189] Copilot deleted entire directory instead of renaming it, losing a full day’s work without confirmation. <https://github.com/microsoft/vscode/issues/280008>.
- [R190] Copilot erases entire file contents when asked to make changes, claims no problems occurred, tells user to fix it themselves. <https://github.com/microsoft/vscode-copilot-release/issues/14019>.
- [R191] Copilot Chat v0.26.7 repeatedly deleted entire codebase when user prompted simple edits. <https://github.com/microsoft/vscode-copilot-release/issues/14017>.
- [R192] Copilot agent repeatedly wipes entire files when attempting targeted edits via `insert_edit_into_file` tool, entering endless destructive loop. <https://github.com/microsoft/vscode-copilot-release/issues/14108>.
- [R193] Copilot in ask mode switched to workspace mode and endlessly attempted file modifications without permission. <https://github.com/microsoft/vscode-copilot-release/issues/14032>.
- [R194] Copilot reads `.env` files containing secrets with no ignore-file mechanism to exclude sensitive paths. <https://github.com/microsoft/vscode-copilot-release/issues/13619>.
- [R195] Filename-based prompt injection in Copilot Chat Agent mode causes agent to execute malicious `setup.py` that exfiltrates system files to attacker-controlled server. <https://www.tenable.com/security/research/tra-2025-53>.
- [R196] VS Code Copilot agent mode tracks deleted files in working set and recreates them as empty on reload, breaking builds. <https://github.com/microsoft/vscode/issues/257998>.
- [R197] Copilot with Claude Haiku 4.5 creates dozens of unwanted markdown files, ignoring explicit user instructions not to. <https://github.com/microsoft/vscode-copilot-release/issues/14045>.
- [R198] Copilot IntelliJ Agent mode deletes existing code in large files when prompted to add new content. <https://github.com/microsoft/copilot-intellij-feedback/issues/373>.
- [R199] Copilot suggested deleting entire folder when asked to move a file; permission prompt prevented execution. <https://github.com/microsoft/vscode-copilot-release/issues/13664>.
- [R200] Prompt injection in GitHub Copilot enables RCE by silently writing `settings.json` to enable YOLO mode, bypassing permission prompts. <https://nvd.nist.gov/vuln/detail/CVE-2025-53773>.
- [R201] VS Code Copilot window reload triggers automatic file creation/modification in WSL, risking corruption of untracked project files. <https://github.com/microsoft/vscode-copilot-release/issues/13590>.
- [R202] Copilot ran Vite initialization commands in wrong directory, deleting all files from user’s dev directory. <https://github.com/microsoft/vscode-copilot-release/issues/13763>.
- [R203] Copilot Agent in Visual Studio repeatedly deletes parts of existing code during edits, leaving files unbuildable with syntax errors. <https://github.com/orgs/community/discussions/161952>.
- [R204] CVE-2025-55319 — prompt injection in VS Code Copilot agent mode enables token exfiltration, config modification, and arbitrary code execution via malicious GitHub issues. <https://nvd.nist.gov/vuln/detail/CVE-2025-55319>.
- [R205] VS Code Copilot agent mode creates files at Windows paths instead of WSL remote paths. <https://github.com/microsoft/vscode-copilot-release/issues/9323>.
- [R206] Copilot suggested a wrong Linux command that deleted the Django backend folder in a project. <https://github.com/microsoft/vscode-copilot-release/issues/13742>.
- [R207] Copilot Agent creates files in wrong Windows directory instead of WSL project path due to path resolution bug. <https://github.com/microsoft/copilot-intellij-feedback/issues/278>.
- [R208] Cursor Agent deliberately bypasses `.cursorignore` by falling back to shell commands (`ls`, `cat`, `echo`) to read and write protected `.env` file containing private keys. <https://forum.cursor.com/t/ai-agent-bypassing-security-restrictions-hacks-its-access-to-restricted-files/119426>.
- [R209] Cursor Agent escaped project scope and grepped user home directory while fixing compilation errors, searching personal data without permission. <https://forum.cursor.com/t/cursor-agents-should-be-restricted-from-access-of-files-outside-the-project-without-permission/149418>.
- [R210] Cursor Agent used `echo »` to modify `.zshrc` without permission, bypassing allowlist intended for read-only commands via shell redirect. <https://forum.cursor.com/t/allowlist-should-not-extend-to-redirects/139150>.
- [R211] Four methods to bypass Cursor’s auto-run denylist — obfuscation, subshell, shell scripts, and bash quoting — render the defense useless. <https://www.backslash.security/blog/cursor-ai-security-flaw-autorun-denylist>.
- [R212] Case-sensitive file protection bypass allows prompt injection to overwrite `.cursor/mcp.json` for RCE on case-insensitive filesystems. <https://nvd.nist.gov/vuln/detail/CVE-2025-59944>.
- [R213] Claude 4 models in Cursor bypass `.cursorignore` and `globalignore` to read and alter `.env` files via shell commands. https://www.reddit.com/r/cursor/comments/1l27xof/claude_bypasses_globalignore_rules/.
- [R214] Cursor CLI auto-loads project-local config that overrides global allowlist, enabling RCE via prompt injection in malicious repositories. <https://nvd.nist.gov/vuln/detail/CVE-2025-61592>.
- [R215] Prompt injection achieves RCE by modifying Cursor CLI sensitive config files via case-insensitive path bypass. <https://nvd.nist.gov/vuln/detail/CVE-2025-61593>.
- [R216] Cursor crashes corrupt project files by injecting whitespace and reverting committed changes, wiping 100+ files. <https://forum.cursor.com/t/cursor-crashing-and-modifying-files-unexpectedly-project-partially-wiped/147372>.
- [R217] Cursor built-in `read_file` and `write` tools blocked by `.cursorignore` even when ignore file is empty, comments-only, or deleted. <https://forum.cursor.com/t/read-file-and-write-tools-blocked-by-cursorignore-even-when-file-is-empty-or-deleted/148852>.
- [R218] Cursor AI went haywire during auto-run session and deleted 90% of project files including `.git` folder, with delete file protection disabled. <https://dredyson.com/why-my-cursor-ide-deleted-my-entire-project-and-how-i-got-it-back/>.
- [R219] Cursor Agent autonomously mass-deleted project files then crashed, destroying chat history and checkpoints needed for recovery. <https://forum.cursor.com/t/help-needed-asap-cursor-deleted-my-whole-proejct/97589>.
- [R220] Cursor allows creating new dotfiles (`.vscode/settings.json`) without approval, enabling RCE via prompt injection chain (CVE-2025-54130). <https://nvd.nist.gov/vuln/detail/CVE-2025-54130>.
- [R221] Cursor Agent’s built-in edit tool deletes or empties files when making multiple sequential edits to the same file due to a tool bug. <https://forum.cursor.com/t/agents-are-deleting-files-while-editing/145458>.
- [R222] Cursor agent read and rewrote `.env` file containing sensitive data using built-in Read/Edit tools with no `.cursorignore` protection. https://www.reddit.com/r/cursor/comments/1jijrxn/cursor_is_reading_my_env_file/.
- [R223] Cursor Agent read `.env` file despite it being listed in `.cursorignore`, bypassing ignore-file defense via shell tool calls. <https://forum.cursor.com/t/cursor-reads-env-even-though-it-is->

- on-cursorignore/136998.
- [R224] Cursor’s built-in edit and checkpoint tools deleted user files during editing sessions, causing hours of lost work; recurring bug confirmed by multiple users over six months. <https://forum.cursor.com/t/cursor-deleted-my-file/45020>.
 - [R225] Cursor Agent occasionally deletes files without permission while using claude-4.5-haiku in Auto-Run mode. <https://forum.cursor.com/t/cursor-may-occasionally-delete-files-without-permission/138395>.
 - [R226] JSON schema auto-download default setting enables silent data exfiltration via agent-written JSON files after prompt injection. <https://nvd.nist.gov/vuln/detail/CVE-2025-49150>.
 - [R227] Cursor Composer asked to move one file went into tool-calling loop and deleted all project files. <https://forum.cursor.com/t/cursor-potentially-deleted-all-of-my-files/141670>.
 - [R228] Five of six LLMs in Cursor attempt to bypass .cursorignore by using shell commands (cat, head, sed) to read protected files. <https://forum.cursor.com/t/security-privacy-issue-llms-will-attempt-to-circumvent-cursorignore/148441>.
 - [R229] Feature request for file lock/read-only mode after Cursor Agent hallucinated edits to a reference-only file. https://www.reddit.com/r/cursor/comments/1ih6yhx/there_should_be_a_lock_filereadonly_feature/.
 - [R230] Cursor IDE randomly deleted newly created files and reverted recent changes upon opening project, affecting 15 files. <https://github.com/cursor/cursor/issues/3897>.
 - [R231] Hidden prompt injections in README files bypass Cursor denylist and exfiltrate API keys and SSH credentials via tool chaining. <https://hiddenlayer.com/innovation-hub/how-hidden-prompt-injections-can-hijack-ai-code-assistants-like-cursor/>.
 - [R232] Cursor AI assistant generated recursive backup script using shutil.copytree that created infinitely nested directories, hit Windows path limit, and deleted entire project folder. <https://forum.cursor.com/t/critical-bug-ai-assistant-deleted-entire-directory-via-recursive-backup-loop/138236>.
 - [R233] Cursor deletes entire source files when user clicks Reject on AI-proposed changes, even for files not created by AI. <https://forum.cursor.com/t/cursor-deletes-my-source-file-altogether-without-warning-when-i-click-reject-ai-didnt-even-create-the-file-in-the-first-place/28718>.
 - [R234] Cursor agent’s Remove-Item commands with unescaped PowerShell bracket wildcards wiped entire F: drive instead of nested project directory. <https://forum.cursor.com/t/possible-catastrophic-deletion-triggered-by-remove-item-path-expansion-in-cursor-f-drive-wiped/138966>.
 - [R235] Cursor replaced failing test implementation with expect(true).toBe(true) stub and claimed tests passed. https://www.reddit.com/r/cursor/comments/1k05u84/test_passed/.
 - [R236] Cursor executed `rm -rf && ls -la` during development session, deleting personal files from macOS home directory. <https://forum.cursor.com/t/cursor-ai-executes-destructive-command-rm-rf-during-development-session/129401>.
 - [R237] Cursor with Gemini 2.5 Pro added home directory as extra `rm -rf` argument, wiping desktop and keychain. <https://forum.cursor.com/t/cursor-tried-to-wipe-my-computer/107142>.
 - [R238] Cursor agent ran `rmdir /s /q` with broken quoting, wiping 350GB+ from D: drive despite File-Deletion and External-File Protection being enabled. <https://forum.cursor.com/t/cursor-deleted-my-d-drive/149859>.
 - [R239] Cursor 2.0 sandbox grants full filesystem read access; agent ran `cat /.npmrc` exposing npm auth token to LLM. <https://luca-becker.me/blog/cursor-sandboxing-leaks-secrets/>.
 - [R240] Widespread reports of Cursor with Sonnet 3.7 randomly deleting working code during edit sessions, especially in long chats. <https://news.ycombinator.com/item?id=43298275>.
 - [R241] Cursor agent deletes files via PowerShell `Remove-Item` bypassing delete file protection setting. <https://github.com/cursor/cursor/issues/2895/>.
 - [R242] CTO reports Cursor repeatedly modifying restricted framework files when team members vibe-code with incomplete context; shares VS Code `readonlyInclude` workaround. https://www.reddit.com/r/cursor/comments/1okujl4/how_to_stop_cursor_to_edit_files_that_you_do_not/.
 - [R243] Cursor Agent can be prompt-injected to write `.code-workspace` file, bypassing CVE-2025-54130 sensitive-file protections and achieving RCE. <https://nvd.nist.gov/vuln/detail/CVE-2025-61590>.
 - [R244] Cursor in YOLO mode deleted everything on user’s computer including itself during Express.js to Next.js migration. <https://news.ycombinator.com/item?id=44262383>.
 - [R245] Gemini CLI `@-file` reference resolves local `/tmp` subdir to system `/tmp`, causing agent to search outside project scope. <https://github.com/google-gemini/gemini-cli/issues/22392>.
 - [R246] Gemini CLI deleted git branches despite deny rules in policy engine; policy regex bug prevents `commandPrefix` rules from matching real commands. <https://github.com/google-gemini/gemini-cli/issues/20355>.
 - [R247] Gemini CLI model occasionally uses destructive commands like `git reset -force` instead of safer alternatives. <https://github.com/google-gemini/gemini-cli/issues/22672>.
 - [R248] Gemini CLI stuck in loop repeatedly deleting `.env` with `rm -rf` despite user explicitly telling it to stop. <https://github.com/google-gemini/gemini-cli/issues/18624>.
 - [R249] Gemini CLI command restrictions bypassed via `eval` injection in shell scripts, enabling unrestricted unsandboxed command execution. <https://github.com/google-gemini/gemini-cli/issues/18194>.
 - [R250] Extension update permanently deletes old extension before verifying new one loads, leaving user with no working extension on failure. <https://github.com/google-gemini/gemini-cli/issues/21671>.
 - [R251] Gemini CLI made unauthorized code changes then executed destructive git checkout, wiping all uncommitted work. <https://github.com/google-gemini/gemini-cli/issues/22649>.
 - [R252] Gemini CLI autonomously ran `git clean` and `git reset -hard`, deleting 11.5 hours of uncommitted architectural work. <https://github.com/google-gemini/gemini-cli/issues/22010>.
 - [R253] Environment variable sanitization bypassed via `hookConfig.env` spread ordering in `executeCommandHook`, enabling `LD_PRELOAD/PATH` injection for arbitrary code execution. <https://github.com/google-gemini/gemini-cli/issues/22503>.
 - [R254] Gemini CLI `read_file` tool enforces ignore patterns with no override, but `search_file_content` supports `no_ignore`; shell `cat` is the only workaround. <https://github.com/google-gemini/gemini-cli/issues/13775>.
 - [R255] Gemini CLI `list_dir` tool ignores `.geminiignore` and `.aiexclude` rules, exposing excluded directories via built-in tool. <https://github.com/google-gemini/gemini-cli/issues/14546>.
 - [R256] Gemini CLI used `write_file` to overwrite hundreds of lines of project log history during long session with context degradation. <https://github.com/google-gemini/gemini-cli/issues/17534>.
 - [R257] Gemini CLI regularly replaces unchanged code with "(rest of methods ...)" placeholder text in proposed diffs, marking real code as deleted. <https://github.com/google-gemini/gemini-cli/issues/19858>.
 - [R258] Gemini CLI ignored plan mode and directly implemented code changes instead of proposing a plan. <https://github.com/google-gemini/gemini-cli/issues/22751>.
 - [R259] Policy engine bypass via nested property injection in JSON-stringified tool arguments allowed unauthorized shell commands. <https://github.com/google-gemini/gemini-cli/issues/19762>.

- [R260] Feature request to warn users when policy auto-approves dangerous commands like rm. <https://github.com/google-gemini/gemini-cli/issues/21596>.
- [R261] Gemini CLI replace tool repeatedly truncated large log file, silently deleting historical project data. <https://github.com/google-gemini/gemini-cli/issues/17562>.
- [R262] Gemini CLI used rm -rf to delete home directory contents (documents, downloads, desktop) after user clicked "allow always" during npm install. <https://github.com/google-gemini/gemini-cli/issues/2617>.
- [R263] macOS seatbelt sandbox silently drops allowedPaths beyond hardcoded limit of 5, with no warning and inconsistent cross-driver behavior. <https://github.com/google-gemini/gemini-cli/issues/22536>.
- [R264] Gemini CLI shell escaping logic loop in WSL broke source code despite user warnings. <https://github.com/google-gemini/gemini-cli/issues/19879>.
- [R265] Gemini CLI agent prefers shell redirects over built-in file tools, triggering excessive permission prompts. <https://github.com/google-gemini/gemini-cli/issues/22638>.
- [R266] Stale closure in onModelChange silently overwrites user-edited .gemini/settings.json with boot-time snapshot. <https://github.com/google-gemini/gemini-cli/issues/20402>.
- [R267] Gemini CLI agent silently falls back to write_file after str_replace failure, rewriting entire file with incorrect or incomplete content. <https://github.com/google-gemini/gemini-cli/issues/20799>.
- [R268] Gemini CLI attempted sudo rm -rf /Library/Developer/CommandLineTools in YOLO mode without user confirmation. <https://github.com/google-gemini/gemini-cli/issues/21594>.
- [R269] Gemini CLI overwrites large files with truncated read fragments via write_file, then compounds loss with destructive git checkout recovery. <https://github.com/google-gemini/gemini-cli/issues/18764>.
- [R270] Gemini CLI agent bypassed project-defined safe_log_append wrapper, used native write_file on large document, causing truncation and data loss. <https://github.com/google-gemini/gemini-cli/issues/20866>.
- [R271] Gemini CLI uses WRITE_FILE instead of partial edits, causing repeated loss of important code in project files. <https://github.com/google-gemini/gemini-cli/issues/21549>.
- [R272] Gemini CLI used write_file to fully overwrite persistent project log, destroying all historical context due to context compression. <https://github.com/google-gemini/gemini-cli/issues/17583>.
- [R273] Gemini CLI uses WriteFile instead of Replace for edits, overwriting entire files and losing existing code. <https://github.com/google-gemini/gemini-cli/issues/20321>.
- [R274] Gemini CLI deleted .git.bak backup directory during cleanup, destroying local commit history recoverable only via ZFS snapshot. <https://github.com/google-gemini/gemini-cli/issues/19729>.
- [R275] Gemini CLI deleted user's *_new.json files without consent after JSON comparison task, violating explicit "do not change other files" instruction. <https://github.com/google-gemini/gemini-cli/issues/15822>.
- [R276] Gemini bypassed .gitignore restriction by using shell cat after built-in tool correctly blocked access. https://www.reddit.com/r/AntigravityGoogle/comments/1premuk/gemini_bypasses_gitignore_restrictions/.
- [R277] Lovable replaced user's Supabase Edge function code with generic dummy function to "fix" deployment errors, destroying work. https://www.linkedin.com/posts/anubhave_lovable-ai-coding-activity-7367214709298601984-3uUl.
- [R278] AI agent(s) secretly deleting project files in the background without user prompting; culprit unidentified among five running agents. https://www.reddit.com/r/cursor/comments/1k1h24a/ai_agent_secretly_deleting_my_files/.
- [R279] Argument injection in pre-approved safe commands achieves RCE across three unnamed AI agent platforms, bypassing human-in-the-loop approval. <https://blog.trailofbits.com/2025/10/22/prompt-injection-to-rce-in-ai-agents/>.
- [R280] Rules File Backdoor — hidden unicode characters in AI config files cause Cursor and GitHub Copilot to silently inject malicious code into generated files. <https://www.pillar.security/blog/new-vulnerability-in-github-copilot-and-cursor-how-hackers-can-weaponize-code-agents>.
- [R281] Malicious AI agent skills bypass skill scanners by bundling test files that Jest/Vitest auto-execute, achieving RCE and secret exfiltration. <https://www.gecko.security/blog/rce-in-your-test-suite-ai-agent-skills-bypass-skill-scanners>.
- [R282] Security researcher exploited 7 of 16 YC AI agents via prompt injection, leaking user PII, rewriting server code to remove security controls, and taking over databases. <https://casco.com/blog/we-hacked-ycombinator-agents>.
- [R283] OpenCode bash allowlist uses implicit prefix matching, allowing shell redirects to write arbitrary files outside project scope. <https://github.com/anomalyco/opcode/issues/5330>.
- [R284] OpenCode agent bypassed denied git reset via bash -c subshell and overwrote .env with echo instead of built-in tools. <https://github.com/anomalyco/opcode/issues/4642>.
- [R285] OpenCode agent bypassed denied .env read permission by using bash cat command, then overwrote the file via shell redirect. https://www.reddit.com/r/opencodeCLI/comments/1qj4gr1/ai_just_worked_around_env_read_permissions/.
- [R286] OpenCode agent used rmdir on Windows reserved name (nul), causing recursive deletion of all project folders on K: drive. <https://github.com/anomalyco/opcode/issues/17836>.
- [R287] AI agents repeatedly modified, deleted, and faked TDD test files instead of implementing features; developer deployed chmod 444 as filesystem-level defense. <https://blakelink.us/posts/how-i-use-chmod-to-stop-ai-agents-from-cheating-on-my-tests/>.
- [R288] Secret exfiltration via markdown rendering and RCE via mcp.json modification in Void Editor through prompt injection in code files. <https://idanhabler.medium.com/d%C3%A9j%C3%A0-in-the-void-an-agentic-ide-compromised-by-known-tricks-56c3c492a077>.
- [R289] Windsurf Cascade repeatedly deletes working code when asked to adjust CSS, revert only recovers 10%. https://www.reddit.com/r/Codeium/comments/1gyo145/windsurf_keeps_on_deleting_code.
- [R290] Windsurf Cascade ran rm instead of git rm, permanently deleting CLI project's core files including rulebook and init.ts. <https://www.threads.com/@carmelyne/post/DCiQqady7aS>.