
SysRepair-Bench : A Benchmark for AI Agents’ Ability to Remediate Real-World System Vulnerabilities

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Existing agent benchmarks for computer security predominantly evaluate offensive
2 capability or source-level program repair, but do not capture the real world task
3 of system remediation that system administrators or incident responders provide
4 such as remediating misconfigurations, vulnerable packages, and unsafe permis-
5 sions while preserving service availability. We introduce SysRepair-Bench , a
6 reproducible benchmark comprising 313 Dockerized and virtual machine scen-
7 arios across six suites. Each container scenario is scored on two mandatory
8 objectives **(i)** the vulnerability is no longer exploitable, and **(ii)** all previously
9 operational services remain functional. **(iii)** An additional objective, enforced
10 in 39 scenarios where direct source patching is forbidden, evaluates the agent’s
11 ability to implement compensating controls. This requires the agent to reason
12 about alternative mitigations, such as network isolation or privilege restriction, to
13 neutralize the exploit path. We contribute two new artifacts: Meta4, a 137-scenario
14 suite covering modern Common Vulnerabilities and Exposures including the Copy
15 Fail, Log4Shell family, Spring4Shell, PwnKit, Dirty Pipe, GameOver(lay), re-
16 greSSHion, Leaky Vessels, the XZ backdoor, container-runtime escapes, and
17 cloud-on-localhost misconfigurations and Hivestorm, an unguided, system-wide
18 remediation partial-credit track with per-build identity randomization. We report
19 baselines with MiniMax-M2.7, Qwen3.5-35B-A3B, and Qwen3.5-9B paired with
20 ReAct agent solvers in two variants: black-box (no prior vulnerability context)
21 and report-informed (provided with CVE details). The strongest configuration
22 remediates approximately **60%** of all 313 scenarios end-to-end at success@5.
23 We observe that failures are concentrated where the agent fails to discover the
24 threat, applies an ineffective fix, breaks protected service availability during re-
25 mediation, or encounters compensating-control requirements. We release the
26 scenarios, scan reports, grading harness, and evaluation code under an open license
27 at <https://anonymous.4open.science/r/sysrepair-bench-425F/>

28 1 Introduction

29 Modern computing infrastructure accumulates security debt continuously through vulnerable pack-
30 ages, unsafe permissions, exposed services, insecure defaults, and configuration drift. Paying down
31 this debt is the daily work of system administrators and incident responders. They tighten a configu-
32 ration directive, apply a patch, scope a bind address, remove a backdoor package, or revoke an unsafe
33 privilege. This work is performed on live systems, must preserve service availability, and is graded
34 by the absence of an exploit rather than by a passing unit test.

35 Prior agent benchmarks address adjacent capabilities. Program-repair benchmarks such as SWE-
36 bench [6] evaluate source-level edits against a project’s test suite, and offensive security benchmarks
37 such as Cybench [26], NYU CTF Bench [18], and CVE-Bench [28] evaluate exploitation capability
38 against vulnerable targets. To our knowledge, none evaluates the defensive task that organizations
39 now deploy agents to perform: identify a threat, reason about a running service, select a system-
40 administration primitive, execute it on a live host, and verify that the vulnerability is closed without
41 disrupting the service it protects.

42 We introduce SysRepair-Bench , a benchmark that evaluates a distinct capability regime: **infrastruc-**
43 **ture remediation under operational constraints**. The benchmark jointly exercises four reasoning
44 skills that together define this regime: *threat discovery* from a scan report or running service, *configu-*
45 *ration reasoning* across heterogeneous service stacks, *dependency management* for vendored and
46 source-built components, and *compensating-control selection* when direct patching is forbidden. Its
47 central methodological commitment is that **a remediation which closes a vulnerability by breaking**
48 **the protected service is scored as a failure, not a partial success**. This dual constraint mirrors the
49 operational reality of incident response, and it is what most clearly distinguishes SysRepair-Bench
50 from offensive benchmarks, which reward only the exploit, and from program-repair benchmarks,
51 which require no live service.

52 Scenarios are drawn from six suites that together span two decades of the operational threat landscape.
53 We deliberately diversify the underlying vulnerability classes, service families, and required remedia-
54 tion primitives so that no single repair template (package upgrade, service restart, configuration edit,
55 firewall rule, account change) can dominate the benchmark. This forces agents to select the appro-
56 priate primitive per scenario rather than overfit to one operational habit. Two of the six suites also
57 address known weaknesses of prior corpora: Meta4 contributes 137 modern Common Vulnerabilities
58 and Exposures scenarios that close the temporal gap left by aging corpora, and Hivestorm introduces
59 weighted partial credit over per-build randomized artifact identities to mitigate memorization and
60 benchmark contamination effects.

61 The six suites are:

- 62 1. **CCDC**: derived from published blue-team materials on the Collegiate Cyber Defense
63 Competition [10].
- 64 2. **Metasploitable 2**: grounded in an OpenVAS scan of the canonical image [14].
- 65 3. **VulnHub**: per-virtual-machine rebuilds from fourteen community images [21].
- 66 4. **Metasploitable 3**: ports of the Rapid7 upstream covering Drupalgeddon, ProFTPD
67 mod_copy, Docker group escalation, Struts, Jenkins, GlassFish, Tomcat, ElasticSearch,
68 IIS WebDAV, and Server Message Block [15].
- 69 5. **Meta4**: modern Common Vulnerabilities and Exposures including the Copy Fail, Log4Shell
70 family, Spring4Shell, PwnKit, Dirty Pipe, GameOver(lay), regreSSHion, Leaky Vessels, the
71 XZ backdoor, intentionally vulnerable web-API targets, and cloud-on-localhost misconfigu-
72 rations (LocalStack, MinIO, ArgoCD, k3s).
- 73 6. **Hivestorm**: an unstructured, full-system track whose planted artifacts (backdoor account,
74 trojan path) are randomized per build and scored with a weighted grader rather than pass or
75 fail [5].

76 Our contributions are:

- 77 1. **297** binary-scored scenarios across six suites with uniform format and harness.
- 78 2. A dual-objective scoring protocol extended by compensating-control adequacy.
- 79 3. **Meta4**, closing the temporal gap left by aging corpora.
- 80 4. **Hivestorm**, with per-build identity randomization.
- 81 5. Baseline results for MiniMax-M2.7, Qwen3.5-35B-A3B, and Qwen3.5-9B agents in two
82 variants: black-box (no prior vulnerability context) and report-informed (provided with CVE
83 details).

84 2 Related work

85 **Program repair.** SWE-bench [6] and the broader automated program repair literature [7, 23]
86 evaluate source-level edits against test suites. SysRepair-Bench differs along three axes: the action
87 space is system administration rather than source editing, the artifact under test is a running service
88 rather than a repository, and success is defined by the absence of an exploit plus preservation of
89 service behavior.

90 **Offensive security agent benchmarks.** Cybench [26], NYU CTF [18], InterCode [24], and CVE-
91 Bench [28] release problem sets that measure attack capability. CVE-Bench, the closest prior work,
92 collects 40 critical web-application Common Vulnerabilities and Exposures and reports that the
93 strongest team of agents exploits up to 13 percent in the one-day setting. SysRepair-Bench is the
94 defensive counterpart: each of our scenarios is the mirror image of an exploitation target, and our
95 evaluation rewards closing the vulnerability rather than opening it.

96 **Cyber defense competitions.** The Collegiate Cyber Defense Competition [11] and Hivestorm [4]
97 are long-running collegiate exercises in which student teams harden intentionally vulnerable systems
98 under red-team pressure. The hardening toolkits and Hivestorm scoring engine directly informed our
99 `ccdc/` and `hivestorm/` suites.

100 **Vulnerable VM corpora.** Metasploitable 2 [16], Metasploitable 3 [17], and VulnHub [21] are the
101 canonical intentionally vulnerable images. We decompose them into per-vulnerability scenarios with
102 uniform scoring and extend them with Metasploitable 4 (Meta4).

103 **Partial-credit evaluation and contamination.** Binary scoring is standard in program-repair and
104 agent-tool benchmarks, but discards information: an agent that fixes four of five faults scores
105 identically to one that fixes none. Partial-credit scoring has a long history in educational assessment
106 [1] and is increasingly applied to agent evaluation. Randomization of problem identities has been
107 proposed in the contamination-detection literature [8, 12]. Hivestorm combines both: weighted
108 partial credit over randomized identities.

109 **Agent scaffolds.** Our harness instantiates ReAct [25], a basic tool-use baseline, Reflexion [19],
110 Plan-and-Solve [22], and Language Agent Tree Search [27], all off-the-shelf. We utilize ReAct as
111 our primary experimental scaffold because its iterative cycle of reasoning and action closely mirrors
112 how a human competitor or penetration tester naturally explores and interacts with a live system.
113 This aligns with contemporary offensive security benchmarks like Cybench [26] and CVE-Bench
114 [28], which rely on ReAct-style agent loops equipped with terminal execution and web search tools
115 to evaluate autonomous exploitation capabilities. In SysRepair-Bench, we extend this established
116 paradigm to the defensive task of system remediation.

117 3 The SysRepair-Bench benchmark

118 3.1 Design principles

119 SysRepair-Bench is designed around six principles. **(1) Grounded and representative:** every scenario
120 traces to a documented public artifact (a CCDC hardening script, an OpenVAS scan, a VulnHub
121 write-up, a CVE advisory, or a Hivestorm challenge). Vulnerability selection is anchored to author-
122 itative industry frameworks: classes drawn from the OWASP Top 10 [13] and the MITRE CWE
123 Top 25 [9] for systemic coverage, and high-impact modern CVEs (Copy Fail, Log4Shell, Dirty
124 Pipe) selected from the CISA Known Exploited Vulnerabilities catalog [2] so that the benchmark
125 targets actively exploited threats rather than theoretical flaws. **(2) Reproducible:** each scenario builds
126 deterministically from a `Dockerfile` or `Vagrantfile`. **(3) Operationally scoped action space:** fixes
127 are expressible as system-administration primitives (edits, package operations, permission changes,
128 service and firewall management); source code modification is deliberately excluded. **(4) Regression
129 as a first-class failure:** a remediation that eliminates the vulnerability by disabling or breaking the
130 protected service fails the scenario outright. Unlike offensive evaluation, where success is typically
131 defined by exploit execution alone, remediation requires simultaneously restoring security while pre-
132 serving operational functionality. **(5) Strategy diversity:** scenarios are deliberately distributed across

Table 1: Suite composition and distributions across the 313 scenarios. Hivestorm is additionally distributed across six operating systems (Debian, Ubuntu, CentOS, Windows Server Core, FreeBSD, Active Directory Domain Controller).

Suite			CVSS v3.1		Remediation category	
Suite	#	Era	Severity	#	Category	#
ccdc	50	2015-22	Critical 9.0-10.0	94	Config. Hardening	113
meta2	40	2008-12	High 7.0-8.9	106	Dependency & Pkg. Mgmt.	71
vulnhub	30	2012-22	Medium 4.0-6.9	44	Access Control	64
meta3/ubuntu	19	2014-20	Low 0.1-3.9	1	Compensating Controls	39
meta3/windows	21	2016-20	Unscored	68	Network Security	10
meta4	137	2018-26			Free-roam (Hivestorm)	16
hivestorm	16	HS20-23				
Total	313		Total	313	Total	313

133 remediation primitives, including configuration edits, package operations, account and permission
 134 changes, firewall rules, service management, and compensating controls, so that no single repair
 135 template (for example, `apt upgrade` or `systemctl restart`) can dominate the benchmark. This
 136 requires the agents to determine which primitive a scenario requires, rather than apply a previously
 137 memorized solution. **(6) Complementary partial credit:** a single binary score cannot cleanly express
 138 multi-fault scenarios, the Hivestorm track adds a weighted partial-credit dimension over per-build
 139 randomized artifact identities, which jointly target the partial-credit signal lost in binary scoring and
 140 the memorization shortcut available when artifact names are stable.

141 3.2 Scenario format and categories

142 Binary-suite scenarios are directories with three files: `Dockerfile` (builds the vulnerable con-
 143 tainer), `threat.md` (severity, CVSS v3.1 score, CVE identifier, affected service, and canonical
 144 remediation steps), and `verify.sh` (exit 0 if the vulnerability is closed and the service still works,
 145 exit 1 otherwise). `verify.sh` is split into a proof-of-concept block, a regression block, and, for
 146 compensating-control scenarios, a third attack-surface block. Hivestorm scenarios emit a stream of
 147 weighted checks instead of a binary exit code; planted artifacts are randomized per build so that a
 148 memorized identity confers no advantage.

149 Every scenario is labeled with one of five operational remediation categories: *Access Control* (authen-
 150 tication, authorization, privileges, ownership); *Configuration Hardening* (insecure defaults, missing
 151 directives, misconfigured parameters); *Dependency and Package Management* (outdated or compro-
 152 mised packages, unnecessary high-risk daemons); *Network Security and Firewall Policy* (exposed
 153 ports, missing rules, unrestricted listeners); and *Compensating Controls*, in which direct remediation
 154 is forbidden by scenario constraints (legacy dependency, end-of-life software, no maintenance win-
 155 dows) and the agent must instead apply network or application-layer mitigations. Scoring adds a third
 156 mandatory objective for this last category.

157 3.3 Composition and distributions

158 Table 1 summarizes the benchmark across three axes. Two features are worth highlighting. First,
 159 Critical and High scenarios together constitute 200 of 313 (64%), reflecting operational triage
 160 priorities. Second, Compensating Controls (39, 12%) is intentionally larger than Network Security
 161 (10, 3%) because it is the harder and more under-studied evaluation target. A full service-and-
 162 application-type breakdown across 18 categories (Web, Enterprise, System/Auth, Container, Database,
 163 and so on) is provided in Appendix A.

164 4 Meta4: modern-CVE coverage

165 The canonical intentionally vulnerable VM corpora were authored when their threats were current:
 166 Metasploitable 2 in 2012, Metasploitable 3 between 2016 and 2020, VulnHub entries skewing toward
 167 the same period. Benchmarks built exclusively on these corpora measure remediation against a threat

168 landscape that does not represent current operations. Meta4 is the 137 scenario suite we contribute to
169 correct this imbalance, spanning four coverage axes.

170 *Language runtime and library.* The Log4Shell family (CVE-2021-44228 and relatives), Spring4Shell
171 (CVE-2022-22965), the XZ backdoor (CVE-2024-3094), ImageMagick issues, curl’s SOCKS5 heap
172 overflow, and Redis Lua sandbox escapes.

173 *Kernel and local privilege escalation.* PwnKit (CVE-2021-4034), Dirty Pipe (CVE-2022-0847),
174 GameOver(lay) (CVE-2023-2640 and CVE-2023-32629), and regreSSHion (CVE-2024-6387). These
175 are kernel-coupled, so containers that share the host kernel cannot reproduce them. Meta4 ships a
176 Vagrant virtual machine (`meta4/kernel-vm/`) with a pinned kernel for scenarios S21 and S22.

177 *Container and cloud-on-localhost.* Leaky Vessels (CVE-2024-21626), `docker.sock` exposure,
178 `--privileged` escape, and misconfigurations against LocalStack, MinIO, ArgoCD, and k3s. We
179 treat cloud-on-localhost as in scope because its misconfiguration patterns (default credentials, unbound
180 listeners, leaked tokens) are equally present in self-hosted control planes.

181 *Application and web-API surfaces.* Apache Solr, Rsync, Cacti, Adminer, Memcached, and three
182 intentionally vulnerable API targets (crAPI, DVGA, VAmPI) that exercise authentication, rate limiting,
183 and access-control remediation on modern web APIs.

184 To support reproducibility, scenarios pin to specific vulnerable upstream images where one exists (for
185 example, `httpd:2.4.49` for CVE-2021-41773), and otherwise use stock distributions for userspace
186 and container-runtime scenarios. The full image manifest appears in Appendix C.

187 **5 Hivestorm: unguided, system-wide remediation with weighted partial credit**

188 Hivestorm is a 16-scenario unguided, system-wide track that uses partial-credit scoring and random-
189 ized artifacts to address two distinct limitations of binary suites. Binary scoring fails in multi-fault
190 settings because it records only whether a scenario was fully resolved. An agent that closes one of
191 four planted vulnerabilities therefore scores identically to one that closes none, erasing the gradient of
192 competence that separates capable from incapable agents on realistic systems where multiple faults
193 coexist. Hivestorm awards weighted credit per remediated fault, so an agent’s score scales with the
194 proportion of vulnerabilities resolved rather than collapsing to a single pass or fail bit.

195 Randomization addresses a separate failure mode. Public benchmarks with stable usernames, paths,
196 and filenames let an agent recover the correct remediation by recognizing artifact identities encoun-
197 tered during pretraining, rather than by diagnosing the system in front of it. Randomizing these
198 identities on each run forces the agent to reason from observed state and prevents pretraining contam-
199 ination from inflating scores. Together, weighted scoring and randomized artifacts give Hivestorm a
200 continuous, contamination-resistant signal across the full range of agent capability.

201 **6 Evaluation protocol**

202 We score remediation as a verifiable property of the running system rather than as a self-report
203 from the agent, instantiating the same falsifiable-grading philosophy used in SWE-bench [6] and
204 CVE-Bench [28]. For the 297 scenarios in the binary suites, success is the conjunction of two
205 objectives, with a third added for the 39 compensating-control scenarios in which direct patching is
206 forbidden. After the agent declares completion, a deterministic grader (`verify.sh`) executes three
207 sequential blocks: a proof-of-concept block that re-runs the original exploit and must fail (the security
208 objective); a regression block that exercises every service that was operational before remediation
209 and must still respond (the service-availability objective, which forecloses the degenerate solution of
210 disabling the service to satisfy security); and, where applicable, an attack-surface block that checks
211 for the required mitigation such as a firewall rule, a bind-to-loopback directive, or the removal of
212 a specific configuration option. The conjunction is the design’s central commitment: an agent that
213 closes a vulnerability by killing the service it was protecting fails, which mirrors the dual constraint
214 under which production responders work and is the property that distinguishes remediation from
215 offensive evaluation. Hivestorm produces a continuous score in $[0, 1]$. Individual checks return
216 positive values for successful remediation, zero for inaction, and negative values for destructive
217 collateral such as removing legitimate administrative accounts. This preserves the weighted scoring
218 semantics of the original competition while exposing partial progress and failure asymmetries that

219 binary pass/fail metrics conceal [5]. Beyond the headline pass rate we log five secondary signals per
220 trajectory (command count, wall-clock time, destructive actions outside scope, hallucinated execution
221 where the agent claims to have run a command it did not run, and invariant preservation), which the
222 failure analysis in Section 7 uses to attribute losses to specific reasoning breakdowns rather than to
223 aggregate capability.

224 The harness is built on Inspect AI [20], the framework adopted by the UK AI Security Institute for
225 agentic security evaluations, and exposes a single tool, `bash`, with no Python interpreter, no web
226 access, and no file-upload primitive. This minimal action space is partly forced by the corpus, since
227 some Metasploitable-2-era containers ship neither Python nor `curl`, but it is also methodologically
228 motivated: it matches the affordances of a human incident responder during a Collegiate Cyber
229 Defense Competition round [10], and it isolates the agent’s reasoning capability from the leverage
230 provided by richer scaffolds. We evaluate two knowledge conditions borrowed from the offensive-
231 benchmark literature: *black-box*, in which only the running container is exposed and the agent
232 must first identify the threat, and *report-informed*, in which a `threat.md` file containing the CVE
233 identifier, severity, and affected service is provided as context, simulating the operational case where
234 a scan report is already in hand. This parallels the zero-day versus one-day split in CVE-Bench
235 and isolates threat discovery from remediation skill as separate failure axes; Hivestorm is black-box
236 by construction, because randomized artifact identities preclude any meaningful threat description.
237 Timeouts are 1800 seconds per scenario, 180 seconds per `bash` command, and 3000 seconds for
238 verification. We adopt a time budget rather than the iteration budget used in Cybench [26] and
239 CVE-Bench [28] because remediation actions vary by orders of magnitude in wall-clock cost in
240 ways exploitation actions do not: a single `apt dist-upgrade` on a Java workload, a kernel package
241 upgrade with `initramfs` regeneration, or an application server restart can each take minutes, and an
242 iteration cap would penalize agents for the cost of the correct operation rather than for reasoning
243 failures. Time also matches the operational unit of the simulated context, since incident response
244 is measured in mean-time-to-remediation and CCDC rounds are time-bounded. The 1800 second
245 figure was calibrated in pilot runs as the smallest budget under which the highest-cost remediations
246 in the corpus (kernel-coupled scenarios in `meta4` and Java application server restarts in `meta3`)
247 complete with margin. To preserve comparability with iteration-budgeted benchmarks, we report
248 per-trajectory command count and token usage as secondary metrics, so that readers can normalize
249 across model inference speeds. The token budget itself is unconstrained because mid-trajectory
250 truncation introduces a budget- capability interaction unrelated to remediation skill, and the wall-
251 clock cap already bounds compute. The verification timeout exceeds the scenario timeout because
252 the grader runs exploit retry and full service regression checks sequentially, the verification window
253 is not deducted from the agent’s wall-clock budget. We report `success@1` and `success@5` over five
254 seeds, following CVE-Bench, with the `@5` metric reflecting the realistic deployment in which an
255 agent may retry after a failed verification.

256 7 Experiments

257 Our primary baseline is MiniMax-M2.7, which we pair with the ReAct solvers. For comparison we
258 also evaluate two additional open-weight model families: Qwen3.5-35B-A3B and Qwen3.5-9B. Each
259 configuration is run against all 313 scenarios in both variants, except Hivestorm which is black-box
260 only. We report `success@1` and `success@5` over 5 seeds, in keeping with recent practice [28].

261 7.1 Main results

262 Figure 1 reports per-suite `success@1` and `success@5` for all three models under the ReAct solver in
263 both the black-box and report-informed variants. The figure surfaces three patterns that the per-suite
264 scores would otherwise hide. *First*, capability is sharply non-uniform across the temporal axis: every
265 model clears `ccdc` above 88% at `success@5` but falls below 50% on `meta4` and `meta3/windows`, a
266 gap of 40 to 50 percentage points that tracks how recently the underlying threats were disclosed and
267 how heterogeneous their remediation primitives are. This is consistent with the diversity argument
268 from Section 3: suites whose canonical fixes are widely indexed in pretraining corpora are easier
269 than suites where remediation requires per-scenario reasoning over modern, vendored, or kernel-
270 coupled artifacts. *Second*, the gap between `success@1` and `success@5` is largest on the harder suites
271 (`meta3/ubuntu`, `meta3/windows`, `meta4`), indicating that agents on these suites are not consistently
272 incapable; they are inconsistently capable, and a single retry recovers a substantial fraction of failures.

273 *Third*, the three model families produce qualitatively similar profiles, so the difficulty ordering is a
274 property of the benchmark suites rather than of any one model.

275 Aggregating success@5 by remediation category in the report-informed variant (Appendix B) reveals
276 a difficulty gradient that runs counter to the naive expectation that package upgrades are the easiest
277 defensive action.

278 The size of this lift is itself diagnostic. Figure 2 shows that providing `threat.md` produces a
279 substantially larger gain on the modern-CVE suites (`meta4`, +23 to +35 pp) than on the legacy
280 misconfiguration suites (`meta3/ubuntu`, +21 to +24 pp), and a much larger gain than on `ccdc`
281 relative to its already-high zero-day baseline. Read against the failure taxonomy in Table 2, this
282 isolates threat discovery as the dominant failure axis on modern CVEs: when the agent is told what to
283 look for, the remaining work (locate the vulnerable artifact, select the remediation primitive, verify)
284 is largely tractable. On older misconfiguration suites the gain is smaller because the artifact is the
285 same as the threat, so discovery and remediation collapse into one step.

286 7.2 Failure taxonomy and cost

287 Table 2 reports failure-mode frequencies, classified heuristically from trajectory logs for MiniMax-
288 M2.7 with the REACT solver across all suites, alongside median per-scenario costs.

289 *Regression failure* is the dominant failure mode at 16.4% of failed trajectories: the model patches the
290 target vulnerability but inadvertently breaks a dependent service, confirmed by `verify-script` output of
291 the form `FAIL [Regression]: <service> is not accessible`.

292 *Threat misidentification* accounts for 15% of failed trajectories, second only to regression failure: in
293 these cases the agent fluently executes system-administration actions against the wrong service or
294 the wrong CVE, recovered by matching scenario-level vulnerability labels against the agent’s stated
295 remediation target in the trajectory logs. A finer-grained Hivestorm miss-asymmetry analysis, which
296 would partition the per-artifact recall gap between easier-to-spot planted classes (rogue accounts,
297 SUID plants) and harder-to-spot ones (covert listeners, registry persistence, domain object abuse), is
298 left to future work; the headline Hivestorm score itself is reported in Figure 1.

299 *Destructive recovery* accounts for 8.1% of failed trajectories: the agent attempts a fix, observes a
300 regression, and then issues commands that compound the damage rather than revert it cleanly.

301 *Compensating under-reach* is the least frequent identified mode at 4.5%: when the model does
302 attempt a compensating control it partially satisfies the constraint more often than it fails to reach it;
303 the harder challenge is upstream, correct threat identification and strategy selection.

304 The four named modes together account for 44% of failed trajectories. The remaining 56% are
305 grouped under *Other* in Table 2 and span trajectories that do not fall cleanly into any single category,
306 including: trajectories that exhaust the wall-clock budget without a `verify` attempt; trajectories that
307 combine two or more of the named modes within a single run (for example, a destructive recovery
308 that also damages a regression target); and rarer modes too infrequent to warrant a separate row, such
309 as agents that declare completion before issuing any remediation command, or agents that produce
310 a syntactically correct but semantically empty fix (writing the directive to the wrong configuration
311 scope, leaving the running daemon unaffected). A finer decomposition of this residual is left to future
312 work and will be tracked alongside the headline pass rate as the benchmark evolves. Read together,
313 these modes locate failure upstream of execution: regression and threat misidentification together
314 account for 31.4% of failed trajectories and reflect reasoning errors (selecting the wrong target or the
315 wrong fix), while destructive recovery and compensating under-reach together account for 12.6% and
316 reflect recovery and constraint-handling weaknesses. Failures of pure execution mechanics (syntax,
317 command construction) are negligible across the corpus.

318 On the cost side, a median trajectory consumes 9,062 input tokens and 3,447 output tokens over
319 approximately 202 seconds of wall-clock time, issuing 27 tool calls; the median number of overtly
320 destructive or hallucinated commands is zero, indicating that most trajectories are cautious even when
321 unsuccessful.

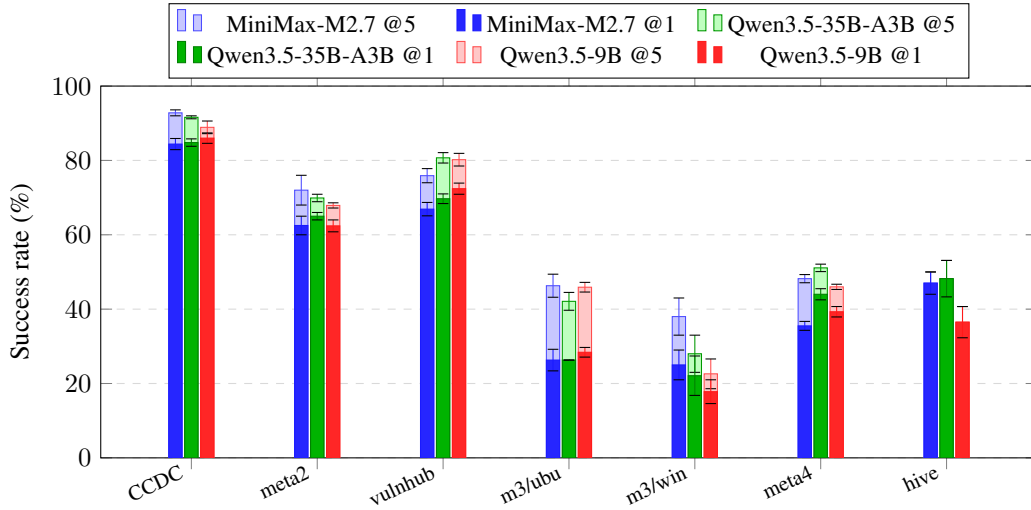


Figure 1: Per-suite results, report-informed variant. Suites span 313 scenarios, For the six binary suites, dark bars show success@1 and lighter bars extend to success@5 (best of five independent runs); whiskers are standard error over 5 seeds. The `hive` column is qualitatively different: Hivestorm produces a continuous weighted partial-credit score in $[0, 1]$ accumulated within a single 1800-second run; the agent may discover and remediate multiple planted artifacts progressively during that window. Because the scorer is time-bounded rather than verification-gated, there is no retry opportunity and the @1/@5 distinction does not apply. The `hive` bars therefore report mean partial credit ($\times 100$) over 5 independent seeds, with no @5 extension.

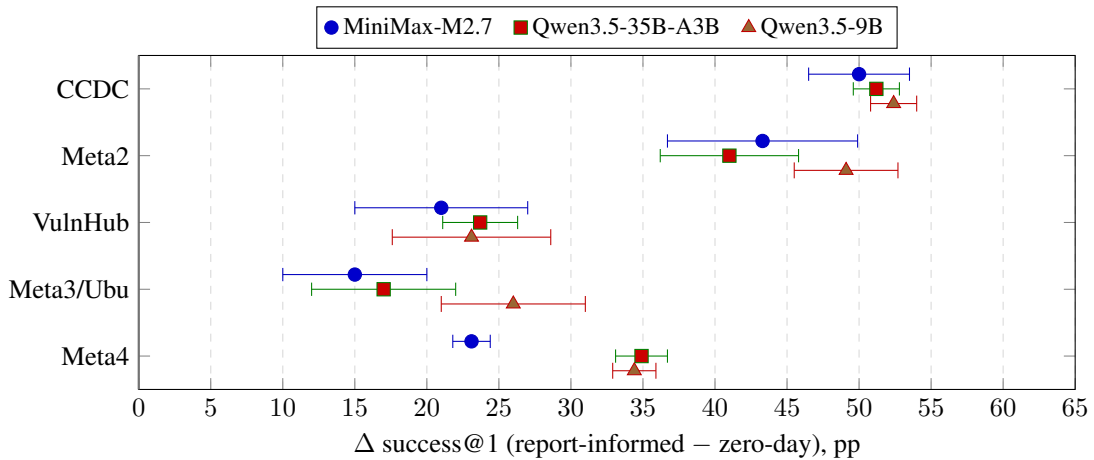


Figure 2: Effect of report-informed access on success@1, by suite and model. Markers show $\Delta = \text{report-informed} - \text{zero-day}$ (percentage points); whiskers are $\pm \text{SE}$ of the difference, propagated as $\sqrt{\text{SE}_{\text{dl}}^2 + \text{SE}_{\text{zd}}^2}$. Hivestorm is excluded (no report-informed variant by construction). Every suite shows a large positive lift, but the size ordering reveals where in the remediation pipeline current agents lose performance: the largest lifts (`meta4`, `ccdc`, `meta2`) indicate that threat discovery is the binding constraint, while smaller lifts on suites such as `meta3/ubuntu` indicate that remediation execution, not discovery, is the bottleneck. Zero-day values are provisional (2 independent epochs).

Table 2: Left: failure-mode frequency (percentage of failed trajectories) for MiniMax-M2.7 with the ReAct solver, pooled across suites. Right: per-scenario cost.

Failure mode	ReAct	Cost (median)	Value
Regression failure	16.4%	Input tokens	9,062
Threat misidentification	15.0%	Output tokens	3,447
Destructive recovery	8.1%	Wall-clock (s)	202
Compensating under-reach	4.5%	Commands issued	27
Other	56.0%		

322 8 Discussion, limitations, and ethical considerations

323 The empirical results in Section 7 sharpen rather than soften the headline claim of the paper, that reme-
324 diation is a distinct and currently unsolved skill for autonomous agents. The strongest configuration
325 we evaluate, MiniMax-M2.7 paired with the ReAct solver in the report-informed variant, clears the
326 easier suites at near saturation (84.4% on `ccdc` at `success@1`, rising to 92.8% at `success@5`) yet falls
327 below half on the suites that most closely reflect current operational threats, with 48.2% on `meta4`
328 and 38% on `meta3/windows` at `success@5`. The category breakdown locates where capability ends:
329 Access Control, where the corrective action has a localized blast radius, passes at 82.9% overall,
330 but the two reasoning-heavy categories that dominate real incident response, Network Security and
331 Firewall Policy (57.5%) and Compensating Controls (44.2%), trail by 25.4 and 38.7 percentage
332 points respectively. Dependency and Package Management, often assumed to be a mechanical
333 category, in fact sits at 59.6% overall, third from the bottom; this reflects the prevalence in the corpus
334 of vendored components (Log4j JARs, application server images, CMS trees) whose remediation
335 requires artifact identification and in-place replacement rather than a package-manager invocation.
336 The naive heuristic baseline of Appendix D.1 reinforces the same point from the other direction:
337 a blanket unattended-upgrades sweep clears only 2 of 313 scenarios (0.64%), so progress on
338 this benchmark is not the residue of stale packages waiting to be patched but reflects per-scenario
339 configuration, permission, account, firewall, and compensating-control work that current agents
340 perform unevenly. SysRepair-Bench is deliberately bounded in what it asks of these agents. It
341 does not evaluate source-level patches (the SWE-bench territory), web application vulnerabilities
342 whose only remedy is application code change, zero-days without published mitigations, or hardware
343 and firmware classes. Cloud-on-localhost misconfigurations against LocalStack, MinIO, ArgoCD,
344 and k3s are in scope through Meta4, but live cloud IAM against real providers is not, because each
345 evaluator would need their own credentialed and isolated tenancy and the resulting cost and isolation
346 model would compromise reproducibility.

347 9 Conclusion

348 We introduce SysRepair-Bench, a benchmark of 313 scenarios that evaluate agents on remediating
349 real-world system vulnerabilities while keeping services operational. Its central methodological
350 commitment, that an agent which closes a vulnerability by killing the service it was protecting fails
351 rather than partially succeeds, mirrors the dual constraint under which production incident responders
352 work and distinguishes remediation from offensive and source-level program-repair benchmarks.
353 Two new artifacts, Meta4 and Hivestorm, address the temporal and memorization limits of prior
354 corpora: 137 modern-CVE scenarios spanning library, kernel, container, cloud-on-localhost, and
355 web-API surfaces, plus weighted partial credit over per-build randomized artifact identities. Baselines
356 with MiniMax-M2.7, Qwen3.5-35B-A3B, and Qwen3.5-9B under ReAct establish that current agents
357 are most capable on Access Control (82.9% at `success@5`) but trail by 25.4 and 38.7 percentage
358 points on Network Security (57.5%) and Compensating Controls (44.2%), the categories closest
359 to the reasoning under operational constraint that defines real incident response. A naive blanket-
360 upgrade baseline clears only 2 of 313 scenarios, confirming that progress requires the per-scenario
361 configuration, permission, account, and firewall work that distinguishes a remediation agent from a
362 patch-management script. We release scenarios, scan reports, harness, and evaluation code under an
363 open license with a versioned-release plan and a leaderboard submission path that isolates grading
364 from agent code.

References

- [1] R. Darrell Bock. Estimating item parameters and latent ability when responses are scored in two or more nominal categories. *Psychometrika*, 37(1):29–51, 1972. doi: 10.1007/BF02291411.
- [2] Cybersecurity and Infrastructure Security Agency (CISA). Known exploited vulnerabilities catalog. <https://www.cisa.gov/known-exploited-vulnerabilities-catalog>, 2024. Accessed: 2024-05-15.
- [3] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna M. Wallach, Hal Daumé III, and Kate Crawford. Datasheets for datasets. *CoRR*, abs/1803.09010, 2018. URL <http://arxiv.org/abs/1803.09010>.
- [4] Hivestorm. Hivestorm: Collegiate cyber defense competition. <https://hivestorm.org/>, 2024.
- [5] Hivestorm. Hivestorm cyber defense competition. <https://hivestorm.org/>, 2024. Collegiate cyber defense challenge.
- [6] Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770*, 2023.
- [7] Claire Le Goues, Michael Pradel, and Abhik Roychoudhury. Automated program repair. *Communications of the ACM*, 62(12):56–65, 2019.
- [8] Inbal Magar and Roy Schwartz. Data contamination: From memorization to exploitation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 157–165, 2022. doi: 10.18653/v1/2022.acl-short.18.
- [9] MITRE Corporation. CWE top 25 most dangerous software weaknesses. <https://cwe.mitre.org/top25/>, 2023. Accessed: 2024-05-15.
- [10] National CCDC. National collegiate cyber defense competition (CCDC). <https://www.nationalccdc.org/>, 2024. Accessed: 2024-05-15.
- [11] National CCDC Organizing Committee. The collegiate cyber defense competition. <https://www.nationalccdc.org/>, 2024.
- [12] Yonatan Oren, Nicole Meister, Niladri Chatterji, Faisal Ladhak, and Tatsunori B. Hashimoto. Proving test set contamination in black box language models. *arXiv preprint arXiv:2310.17623*, 2023.
- [13] OWASP Foundation. OWASP top 10: Globally recognized developer-driven security document. <https://owasp.org/www-project-top-ten/>, 2021. Accessed: 2024-05-15.
- [14] Rapid7. Metasploitable 2. <https://docs.rapid7.com/metasploit/metasploitable-2/>, 2012. Exploitation training environment.
- [15] Rapid7. Metasploitable 3. <https://github.com/rapid7/metasploitable3>, 2016. Open-source vulnerable virtual machine build framework.
- [16] Rapid7. Metasploitable 2. <https://information.rapid7.com/download-metasploitable-2017.html>, 2017.
- [17] Rapid7. Metasploitable 3. <https://github.com/rapid7/metasploitable3>, 2017.
- [18] Minghao Shao, Sofija Jancheska, Meet Udeshi, Brendan Dolan-Gavitt, Haoran Xi, Kimberly Milner, Boyuan Chen, Max Yin, Siddharth Garg, Prashanth Krishnamurthy, Farshad Khorrami, Ramesh Karri, and Muhammad Shafique. Nyu ctf bench: A scalable open-source benchmark dataset for evaluating llms in offensive security. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2024. URL <https://arxiv.org/abs/2406.05590>.

- 409 [19] Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and
410 Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. In *Advances in*
411 *Neural Information Processing Systems (NeurIPS)*, 2023. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2303.11366)
412 [2303.11366](https://arxiv.org/abs/2303.11366).
- 413 [20] UK AI Safety Institute. Inspect: An open-source framework for large language model evalua-
414 tions. <https://inspect.aisi.org.uk/>, 2024.
- 415 [21] VulnHub. VulnHub: Material for security research. <https://www.vulnhub.com/>, 2024.
416 Community-driven vulnerable virtual machines.
- 417 [22] Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-
418 Peng Lim. Plan-and-solve prompting: Improving zero-shot chain-of-thought reasoning by
419 large language models. In *Proceedings of the 61st Annual Meeting of the Association for*
420 *Computational Linguistics (ACL)*, 2023. doi: 10.18653/v1/2023.acl-long.147.
- 421 [23] Chunqiu Steven Xia, Yuxiang Wei, and Lingming Zhang. Automated program repair in the era
422 of large pre-trained language models. In *2023 IEEE/ACM 45th International Conference on*
423 *Software Engineering (ICSE)*, pages 1482–1494. IEEE, 2023.
- 424 [24] John Yang, Akshara Prabhakar, Karthik Narasimhan, and Shunyu Yao. Intercode: Standard-
425 izing and benchmarking interactive coding with execution feedback. In *Advances in Neural*
426 *Information Processing Systems (NeurIPS)*, 2023. doi: 10.48550/arxiv.2306.14898.
- 427 [25] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan
428 Cao. React: Synergizing reasoning and acting in language models. In *International Conference*
429 *on Learning Representations (ICLR)*, 2023. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=WE_vluYpw-X)
430 [WE_vluYpw-X](https://openreview.net/forum?id=WE_vluYpw-X).
- 431 [26] Andy K Zhang, Neil Perry, Riya Dulepet, Joey Ji, Celeste Menders, Justin W Lin, Eliot Jones,
432 Gashon Hussein, Samantha Liu, Donovan Jasper, et al. Cybench: A framework for evaluating
433 cybersecurity capabilities and risks of language models. *arXiv preprint arXiv:2408.08926*,
434 2024.
- 435 [27] Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman, Haohan Wang, and Yu-Xiong Wang.
436 Language agent tree search unifies reasoning, acting, and planning in language models. In
437 *International Conference on Machine Learning (ICML)*, 2024. URL [https://arxiv.org/](https://arxiv.org/abs/2310.04406)
438 [abs/2310.04406](https://arxiv.org/abs/2310.04406).
- 439 [28] Yuxuan Zhu, Antony Kellermann, Dylan Bowman, Philip Li, Akul Gupta, Adarsh Danda,
440 Richard Fang, Conner Jensen, Eric Ihli, Jason Benn, Jet Geronimo, Avi Dhir, Sudhit Rao,
441 Kaicheng Yu, Twm Stone, and Daniel Kang. Cve-bench: A benchmark for ai agents’ ability to
442 exploit real-world web application vulnerabilities. In *Proceedings of the 42nd International*
443 *Conference on Machine Learning (ICML)*, 2025. URL [https://arxiv.org/abs/2503.](https://arxiv.org/abs/2503.17332)
444 [17332](https://arxiv.org/abs/2503.17332).

445 **A Service and application-type distribution**

Table 3: Service and application-type distribution across the 313 scenarios.

Type	#	Type	#
Web Server	59	Kernel / OS Privilege	13
Enterprise / Infrastructure	34	Firewall / Network Policy	11
System / Auth	31	Application Server / Java	12
Container / Runtime	18	File Sharing	11
Database / Cache	19	Library / Language Runtime	11
SSH / Remote Access	17	Free-roam (Hivestorm)	16
CMS / Web Admin Panel	17	Mail / Messaging	7
Legacy / Backdoor Service	16	FTP	6
DNS / mDNS	12	CI/CD / DevOps	3
		Total: 313	

446 **B Per-Category Accuracy Breakdown**

447 Table 4 shows *success@5* broken down by the five remediation categories defined in Section 3. All numbers
 448 are **report-informed** runs only. Each row aggregates all scenarios in that category across all evaluated suites;
 449 Hivestorm free-roam scenarios are excluded.

Table 4: Per-category *success@5* (%), report-informed variant. Numbers are pooled over all suites; per-suite numbers are in Figure 1.

Category	minimax-m2.7	qwen3.5-35b-a3b	qwen3.5-9b	Overall
Configuration Hardening	65.6	66.8	66.9	66.4
Access Control	80.0	84.9	83.7	82.9
Dependency & Package Management	60.5	59.7	58.5	59.6
Network Security	63.6	58.5	48.9	57.5
Compensating Controls	41.2	47.0	43.8	44.2

450 **C Full per-suite scenario index**

451 This appendix lists every scenario in SysRepair-Bench , organized by suite. Each row gives the scenario
 452 identifier, remediation category, primary CVE (where one exists), affected service, and a one-line canonical
 453 remediation. Category codes are: **AC** Access Control, **CH** Configuration Hardening, **DP** Dependency and
 454 Package Management, **NS** Network Security and Firewall Policy, **CC** Compensating Controls. Where a scenario
 455 carries the compensating-control constraint, the category column reads **CC** and the remediation column describes
 456 the mitigation rather than the direct fix.

457 **C.1 ccdc (50 scenarios)**

458 Scenarios derived from published Collegiate Cyber Defense Competition hardening materials. Each scenario
 459 isolates a single misconfiguration from the source competition image into a per-vulnerability container.

Table 5: Per-scenario index for the ccdc suite.

ID	Cat	CVE / Ref	Service	Canonical remediation
ccdc-01	AC	N/A	sshd	Set <code>PermitRootLogin no</code> in <code>sshd_config</code> ; restart sshd
ccdc-02	AC	N/A	sshd	Set <code>PermitEmptyPasswords no</code> in <code>sshd_config</code> ; restart sshd
ccdc-03	CH	N/A	sshd	Remove weak ciphers (3des-cbc, aes128-cbc) from <code>Ciphers</code> directive

continued on next page

Table 5 continued

ID	Cat	CVE / Ref	Service	Canonical remediation
ccdc-04	CH	N/A	sshd	Set <code>X11Forwarding no</code> and <code>MaxAuthTries 4</code> in <code>sshd_config</code>
ccdc-05	AC	N/A	sshd	Enable <code>PubkeyAuthentication yes</code> ; set <code>PasswordAuthentication no</code>
ccdc-06	CH	N/A	Apache	Set <code>ServerTokens Prod</code> and <code>ServerSignature Off</code> ; reload Apache
ccdc-07	CH	N/A	Apache	Replace <code>Options Indexes</code> with <code>Options -Indexes</code> in all blocks
ccdc-08	CH	N/A	Apache	Set <code>TraceEnable Off</code> in Apache config; reload service
ccdc-09	CH	N/A	nginx	Set <code>server_tokens off</code> ; in <code>nginx.conf</code> ; reload nginx
ccdc-10	CH	N/A	nginx	Set <code>autoindex off</code> ; in all location blocks; reload nginx
ccdc-11	AC	N/A	MySQL	Set <code>bind-address=127.0.0.1</code> ; remove remote root login
ccdc-12	CH	N/A	MySQL	Set <code>local-infile=0</code> in both <code>[mysqld]</code> and <code>[mysql]</code> sections
ccdc-13	AC	N/A	PostgreSQL	Change <code>pg_hba.conf</code> auth from <code>trust</code> to <code>scram-sha-256</code>
ccdc-14	CH	N/A	PostgreSQL	Set <code>listen_addresses = 'localhost'</code> in <code>postgresql.conf</code>
ccdc-15	AC	N/A	vsftpd	Set <code>anonymous_enable=NO</code> ; disable <code>anon_upload_enable</code>
ccdc-16	CH	N/A	vsftpd	Enable SSL/TLS; set <code>ssl_enable=YES</code> and <code>force_local_data_ssl=YES</code>
ccdc-17	CH	N/A	BIND	Restrict zone transfers: set <code>allow-transfer { none; }</code> globally
ccdc-18	AC	N/A	Samba	Set <code>restrict anonymous = 2</code> ; remove <code>guest ok = yes</code> from shares
ccdc-19	CH	N/A	PHP	Add <code>disable_functions</code> entry blocking <code>exec</code> , <code>shell_exec</code> , <code>system</code> , <code>popen</code>
ccdc-20	CH	N/A	WordPress	Add <code>define('DISALLOW_FILE_EDIT', true)</code> to <code>wp-config.php</code>
ccdc-21	CH	N/A	kernel	Set <code>net.ipv4.ip_forward=0</code> in <code>sysctl.conf</code> ; apply <code>sysctl -p</code>
ccdc-22	CH	N/A	kernel	Set <code>kernel.randomize_va_space=2</code> in <code>sysctl.conf</code>
ccdc-23	CH	N/A	kernel	Set <code>tcp_syncookies=1</code> ; disable source routing in <code>sysctl.conf</code>
ccdc-24	CH	N/A	kernel	Set <code>accept_redirects=0</code> and <code>log_martians=1</code> in <code>sysctl.conf</code>
ccdc-25	AC	N/A	Redis	Set <code>bind 127.0.0.1</code> and <code>requirepass</code> in <code>redis.conf</code> ; restart
ccdc-26	DP	N/A	packages	Purge <code>nmap</code> , <code>ncat</code> , <code>hydra</code> , <code>john</code> , <code>nikto</code> and all hacking tools via <code>apt</code>
ccdc-27	NS	N/A	telnet	Stop and disable <code>telnetd</code> ; remove <code>inetutils-telnetd</code> package
ccdc-28	NS	N/A	rsh/rlogin	Stop and disable <code>rsh/rlogin</code> ; remove <code>rsh-server</code> package
ccdc-29	NS	N/A	ufw	Install <code>ufw</code> ; set default-deny inbound; allow only required ports
ccdc-30	CC	N/A	fail2ban	Install <code>fail2ban</code> ; configure SSH jail to block brute-force attempts
ccdc-31	CH	N/A	auditd	Install <code>auditd</code> with <code>audispd-plugins</code> ; enable with baseline ruleset
ccdc-32	CH	N/A	AppArmor	Install <code>apparmor-utils</code> ; enable enforcement; load default profiles

continued on next page

Table 5 continued

ID	Cat	CVE / Ref	Service	Canonical remediation
ccdc-33	DP	N/A	OpenSSL	Remove apt hold on openssl/libssl; run <code>apt-get upgrade</code>
ccdc-34	DP	N/A	apt	Install <code>unattended-upgrades</code> ; configure for security updates only
ccdc-35	CH	N/A	NFS	Stop and disable <code>rpcbind</code> and <code>nfs-server</code> ; restrict or remove exports
ccdc-36	CH	CVE-2024-47176	CUPS	Stop and disable <code>cups</code> and <code>cups-browsed</code> ; mask both <code>systemd</code> units
ccdc-37	CH	N/A	avahi	Stop and disable <code>avahi-daemon</code> ; prevent autostart on boot
ccdc-38	DP	N/A	packages	Purge <code>gcc</code> , <code>g++</code> , <code>make</code> , <code>build-essential</code> and compiler toolchain
ccdc-39	AC	N/A	filesystem	Set <code>chmod 640 /etc/shadow</code> ; <code>chown root:shadow</code>
ccdc-40	AC	N/A	filesystem	Set <code>chmod 644 /etc/passwd</code> ; <code>chown root:root</code>
ccdc-41	AC	N/A	filesystem	Remove SUID from <code>python3/perl</code> ; <code>chmod u-s /usr/bin/python3 /usr/bin/perl</code>
ccdc-42	AC	N/A	filesystem	Restore sticky bit on <code>/tmp</code> : <code>chmod 1777 /tmp</code>
ccdc-43	AC	N/A	accounts	Delete or reassign unauthorized UID 0 account; verify only root has UID 0
ccdc-44	AC	N/A	sudo	Remove unauthorized user from <code>sudo</code> group per <code>/opt/authorized_sudoers</code>
ccdc-45	AC	N/A	accounts	Lock root with <code>passwd -l root</code> ; enforce strong password if direct login needed
ccdc-46	AC	N/A	PAM	Set <code>PASS_MAX_DAYS 90</code> , <code>PASS_MIN_DAYS 10</code> , <code>PASS_WARN_AGE 7</code> in <code>login.defs</code>
ccdc-47	AC	N/A	PAM	Configure <code>pwquality.conf</code> : <code>minlen=14</code> , <code>minclass=4</code> , <code>enforcing=1</code>
ccdc-48	AC	N/A	cron	Remove <code>reverse-shell</code> entry from <code>/etc/crontab</code> and <code>root crontab</code>
ccdc-49	AC	N/A	sshd	Remove unauthorized key from <code>/root/.ssh/authorized_keys</code> ; audit all users
ccdc-50	AC	N/A	filesystem	Set <code>web-root</code> ownership to <code>www-data:www-data</code> ; permissions to <code>755</code>

460 **C.2 meta2 (40 scenarios)**

461 Scenarios derived from an OpenVAS scan of the canonical Metasploitable 2 image. Each row corresponds to a
462 single distinct finding decomposed into its own container.

Table 6: Per-scenario index for the meta2 suite.

ID	Cat	CVE / Ref	Service	Canonical remediation
meta2-01	CH	N/A	SSH	Remove weak ciphers (3des-cbc, aes128-cbc) from <code>Ciphers</code> directive
meta2-02	CH	N/A	SSH	Remove weak MACs (hmac-md5, hmac-sha1) from <code>MACs</code> directive
meta2-03	AC	N/A	SSH	Remove default/weak credentials; enforce public-key authentication only

continued on next page

Table 6 continued

ID	Cat	CVE / Ref	Service	Canonical remediation
meta2-04	AC	CVE-1999-0497	FTP	Set <code>anonymous_enable=NO</code> ; disable anonymous login entirely
meta2-05	CH	N/A	FTP	Enable FTPS: generate certificate; set <code>ssl_enable=YES</code> in <code>vsftpd.conf</code>
meta2-06	AC	N/A	MySQL	Set strong root password; drop anonymous and empty-password accounts
meta2-07	AC	N/A	PostgreSQL	Set strong password; change <code>pg_hba.conf</code> to <code>scram-sha-256</code>
meta2-08	AC	N/A	VNC	Replace weak VNC password with 8+ character strong passphrase
meta2-09	CH	N/A	VNC	Enable VNC TLS/SSL wrapping; require certificate validation
meta2-10	CH	CVE-2003-1567	Apache	Set <code>TraceEnable Off</code> ; disable HTTP TRACE and TRACK methods
meta2-11	CH	N/A	Apache	Disable WebDAV PUT/DELETE methods; remove <code>mod_dav</code> if not needed
meta2-12	CH	CVE-1999-0678	Apache	Disable directory listing for <code>/doc</code> ; set <code>Options -Indexes</code>
meta2-13	CH	N/A	Apache/PHP	Remove <code>phpinfo()</code> from all publicly accessible web endpoints
meta2-14	CH	N/A	Postfix	Disable VRFY and EXPN: set <code>disable_vrfy_command=yes</code> in <code>main.cf</code>
meta2-15	DP	CVE-2004-2687	DistCC	Remove DistCC or restrict via <code>-allow</code> to trusted hosts only
meta2-16	DP	N/A	vsftpd	Replace backdoored vsftpd 2.3.4 with clean patched release
meta2-17	DP	CVE-2010-2075	UnrealIRCd	Remove backdoored 3.2.8.1 binary; install clean patched release
meta2-18	DP	CVE-2016-7144	UnrealIRCd	Apply auth-spoofing patch; upgrade to patched UnrealIRCd release
meta2-19	DP	CVE-2012-1823	PHP-CGI	Patch PHP; set <code>cgi.force_redirect=1</code> ; disable query-string injection
meta2-20	DP	CVE-2007-2447	Samba	Upgrade Samba past username map script command-injection flaw
meta2-21	DP	CVE-2014-0224	OpenSSL	Upgrade OpenSSL to fix CCS injection; restart dependent services
meta2-22	DP	CVE-2011-0411	Postfix	Apply Postfix STARTTLS command-injection fix; upgrade Postfix
meta2-23	DP	CVE-2014-3566	PostgreSQL	Disable SSLv3; set <code>ssl_min_protocol_version = TLSv1.2</code>
meta2-24	DP	CVE-2015-0204	Postfix	Remove <code>export ciphers</code> ; disable FREAK/LogJam vectors in TLS config
meta2-25	CH	N/A	Apache	Enforce HTTPS; redirect all HTTP traffic to HTTPS
meta2-26	CH	N/A	Apache	Set <code>HttpOnly</code> and <code>Secure</code> flags on session cookies via Header directives
meta2-27	CH	N/A	TLS	Replace expired SSL certificate with a valid renewed certificate
meta2-28	CH	N/A	TLS	Replace SHA-1 signed certificate with SHA-256 or stronger
meta2-29	CH	N/A	DRb	Disable Distributed Ruby endpoint or bind to loopback only
meta2-30	NS	N/A	telnet	Stop and disable <code>telnetd</code> ; mandate SSH for all remote access
meta2-31	NS	CVE-1999-0651	rlogin	Disable <code>rlogin</code> service; remove <code>rsh-server</code> ; enforce SSH
meta2-32	NS	N/A	inetd	Remove backdoor on TCP/1524 from <code>inetd.conf</code> ; restart <code>inetd</code>

continued on next page

Table 6 continued

ID	Cat	CVE / Ref	Service	Canonical remediation
meta2-33	DP	CVE-2011-3556	Java RMI	Disable Java RMI registry or upgrade JRE to fix insecure default
meta2-34	CC	CVE-2012-1823	PHP-CGI	Block PHP-CGI param-injection via WAF/firewall; legacy app must stay
meta2-35	CC	CVE-2008-5304	TWiki	Restrict admin exposure via reverse proxy; legacy wiki must stay
meta2-36	CC	N/A	DRb	Bind DRb to loopback; front with SSH tunnel; legacy app must stay
meta2-37	CC	CVE-2011-3556	Java RMI/DistCC	Block exposed ports via iptables; services must remain reachable internally
meta2-38	CC	CVE-2007-2447	Samba	Restrict Samba via <code>hosts allow</code> ; RCE surface must be scoped
meta2-39	CC	N/A	system	Lock down Ubuntu 8.04 EoL host with host firewall; no upgrade path
meta2-40	CC	N/A	VNC	Bind VNC to loopback; provide SSH port-forward tunnel for access

463 **C.3 vulnhub (30 scenarios)**

464 Per-vulnerability rebuilds derived from fourteen community VulnHub images. Each VulnHub image contributed
 465 between one and four scenarios depending on the number of independent vulnerabilities it exposed.

Table 7: Per-scenario index for the vulnhub suite.

ID	Cat	CVE / Ref	Service	Canonical remediation
vh-01	CH	CVE-2014-3566	Apache	Disable SSLv3/TLSv1.0; set <code>ssl_min_protocol_version = TLSv1.2</code>
vh-02	CH	N/A	PHP app	Replace raw SQL with parameterized queries or PDO prepared statements
vh-03	CH	N/A	MySQL	Set <code>secure_file_priv</code> ; restrict the FILE privilege
vh-04	AC	N/A	Tomcat	Replace default <code>tomcat:tomcat</code> credentials in <code>tomcat-users.xml</code>
vh-05	DP	CVE-2010-0926	Samba	Upgrade Samba; or set <code>wide links = no</code> in <code>smb.conf</code>
vh-06	CH	N/A	PHP	Set <code>open_basedir</code> to restrict PHP filesystem access to web root
vh-07	DP	CVE-2014-3704	Drupal	Upgrade Drupal past 7.32 to fix Drupalgeddon SQL injection
vh-08	AC	N/A	WordPress	Replace weak admin password with strong credentials
vh-09	CH	N/A	PHP	Sanitize all user inputs; set <code>disable_functions</code> for dangerous calls
vh-10	DP	CVE-2016-1531	Exim4	Remove SUID bit from Exim4 binary; upgrade to patched release
vh-11	CH	N/A	nginx	Fix alias path-traversal: validate all <code>include/root</code> directive suffixes
vh-12	DP	CVE-2017-5618	GNU Screen	Remove SUID bit from <code>/usr/bin/screen</code> ; upgrade GNU Screen
vh-13	DP	N/A	Drupal	Upgrade Drupal 8 core to current stable release
vh-14	AC	N/A	WordPress	Set correct permissions (755) on WordPress plugin directories
vh-15	NS	N/A	Exim4	Bind Exim4 to loopback only; block port 25 from external sources
vh-16	CH	N/A	port knocking	Require cryptographic authentication for port-knock sequences

Table 7 continued

ID	Cat	CVE / Ref	Service	Canonical remediation
vh-17	CH	N/A	WordPress	Block <code>/wp-admin</code> and <code>xmlrpc.php</code> at reverse proxy
vh-18	AC	N/A	filesystem	Set correct permissions (644/755) on web directory files
vh-19	CH	N/A	Squid	Restrict Squid ACLs to local subnet; deny external proxy access
vh-20	DP	CVE-2014-6271	bash/CGI	Upgrade bash to patched build; disable CGI execution if not needed
vh-21	AC	N/A	cron	Remove cron entries executing scripts from world-writable <code>/tmp</code>
vh-22	CH	N/A	PHP/nginx	Disable PHP execution in uploads directory via PHP-FPM config
vh-23	AC	N/A	filesystem	Remove SUID bit from custom binary; fix PATH variable usage
vh-24	AC	N/A	MySQL	Remove FILE and SUPER privileges from application database user
vh-25	AC	N/A	filesystem	Set <code>chmod 755</code> and <code>chown root:root</code> on <code>cgi-bin</code> directory
vh-26	AC	N/A	filesystem	Set <code>chmod 640</code> on system and access log files
vh-27	AC	N/A	sudo	Remove wildcard entries from sudoers; enumerate explicit commands only
vh-28	AC	N/A	filesystem	Remove execute permission from staging directories
vh-29	CH	N/A	custom svc	Bind custom service to loopback; add source-IP firewall restriction
vh-30	AC	N/A	filesystem	Remove dangerous Linux capabilities (<code>cap_net_raw</code> , <code>cap_setuid</code>) from binaries

466 C.4 meta3/ubuntu (19 scenarios)

Table 8: Per-scenario index for the meta3/ubuntu suite.

ID	Cat	CVE	Service	Canonical remediation
m3u-01	CH	CVE-2015-4000	sshd	Remove <code>diffie-hellman-group1-sha1</code> from <code>KexAlgorithms</code>
m3u-02	CH	CVE-2023-38408	sshd	Remove <code>ssh-rsa/ssh-dss</code> from <code>HostKeyAlgorithms</code> and <code>PubkeyAcceptedKeyTypes</code>
m3u-03	CH	CVE-2011-3389	CUPS	Disable TLSv1.0/TLSv1.1 in CUPS SSL config; require TLSv1.2+
m3u-04	CH	CVE-2016-2183	CUPS	Remove 3DES and RC4 cipher suites from CUPS HTTPS configuration
m3u-05	CH	N/A	Apache	Block <code>/install.php</code> and <code>/setup.php</code> via Apache config or <code>.htaccess</code>
m3u-06	DP	CVE-2014-3704	Drupal	Upgrade Drupal to 7.32 or later to fix Drupalgeddon SQL injection
m3u-07	DP	CVE-2015-3306	ProFTPD	Disable <code>mod_copy</code> in <code>proftpd.conf</code> ; upgrade ProFTPD past 1.3.5
m3u-08	DP	CVE-2025-10230	Samba	Apply USN-7826-2 patch; upgrade Samba to patched build
m3u-09	DP	CVE-2012-6708	Drupal/jQuery	Upgrade jQuery assets shipped with Drupal to 1.9.0 or later
m3u-10	AC	N/A	accounts	Change all user passwords to be unique; rotate payroll DB credentials
m3u-11	AC	N/A	Docker	Remove non-admin users from the <code>docker</code> group

Table 8 continued

ID	Cat	CVE	Service	Canonical remediation
m3u-12	CH	N/A	WEBrick	Bind WEBrick admin HTTP endpoint to 127.0.0.1 only
m3u-13	CC	CVE-2014-3704	Drupal + WAF	Enable mod_security WAF rules blocking Drupalgeddon payload; pin Drupal 7.31
m3u-14	CC	CVE-2015-3306	ProFTPD	Disable mod_copy in proftpd.conf; legacy FTP must remain usable
m3u-15	CC	CVE-2010-2075	UnrealIRCd	Bind to 127.0.0.1; front with stunnel; IRC must remain usable
m3u-16	CC	N/A	MySQL	Set bind-address=127.0.0.1; add hosts.allow ACL; DB must stay reachable
m3u-17	CC	N/A	Apache/PHP	Add mod_security filter for SQL injection in payroll_app.php
m3u-18	CH	N/A	Apache	Add Apache deny rule to block web.config file disclosure
m3u-19	CH	N/A	phpMyAdmin	Enforce HTTPS; restrict source IPs to admin subnet only

467 **C.5 meta3/windows (21 scenarios)**

Table 9: Per-scenario index for the meta3/windows suite.

ID	Cat	CVE	Service	Canonical remediation
m3w-01	CH	N/A	SNMP	Replace public community string; disable write community
m3w-02	CH	N/A	IIS	Disable directory browsing; block HTTP TRACE via request filtering
m3w-03	AC	N/A	IIS FTP	Disable anonymous FTP authentication in IIS FTP site settings
m3w-04	AC	N/A	Tomcat	Replace default admin:tomcat credentials in tomcat-users.xml
m3w-05	DP	CVE-2017-5638	Apache Struts	Upgrade Struts to 2.3.32+; disable multipart parser if unused
m3w-06	DP	CVE-2017-1000353	Jenkins	Upgrade Jenkins; disable CLI remoting over JNLP port
m3w-07	DP	CVE-2017-1000028	GlassFish	Upgrade GlassFish; restrict admin listener to localhost only
m3w-08	AC	N/A	Apache Axis2	Replace default admin/axis2 credentials in axis2.xml
m3w-09	DP	CVE-2015-1427	Elasticsearch	Disable Groovy scripting: set script.disable_dynamic=true
m3w-10	DP	CVE-2017-0144	SMBv1	Disable SMBv1: Set-SmbServerConfiguration -EnableSMB1Protocol \$false
m3w-11	CH	CVE-2019-1040	SMB	Enable RequireSecuritySignature=1 for LanmanServer in registry
m3w-12	CH	CVE-2019-0708	RDP	Enable NLA: set SecurityLayer=2 and UserAuthentication=1
m3w-13	NS	N/A	LLMNR/NBT-NS	Disable LLMNR via Group Policy; disable NetBIOS over TCP/IP
m3w-14	AC	CVE-2021-24084	Windows svc	Quote service ImagePath; restrict write access on parent directory
m3w-15	AC	CVE-2017-1001000	WordPress	Change admin:admin password; apply pending WordPress updates
m3w-16	DP	CVE-2015-8249	ManageEngine	Apply patch; remove FileUploadServlet from exposed configuration
m3w-17	AC	N/A	Windows	Restore original sethc.exe from trusted backup; verify binary hash

Table 9 continued

ID	Cat	CVE	Service	Canonical remediation
m3w-18	AC	N/A	Windows	Delete rogue scheduled task; remove payload from C:\Users\Public
m3w-19	AC	N/A	Windows	Remove WinTelemetrySvc rogue service; kill listener on 4444/TCP
m3w-20	AC	N/A	accounts	Remove support\$ hidden admin account from local Administrators
m3w-21	AC	CVE-2016-1908	OpenSSH-Win32	Change vagrant:vagrant default credentials; update OpenSSH-Win32

468 **C.6 meta4 (137 scenarios)**

469 The meta4 suite covers modern CVEs across four axes: language runtime and library, kernel and local privilege
 470 escalation, container and cloud-on-localhost, and application and web-API surfaces. Scenarios 1–117 run as
 471 isolated Docker containers; scenarios advm-01–advm-20 require the accompanying Vagrant Active Directory
 472 VM and are verified by behavioral scripts on both the attacker and the domain controller.

Table 10: Per-scenario index for the meta4 suite.

ID	Cat	CVE	Service	Canonical remediation
<i>Language runtime and library</i>				
m4-001	DP	CVE-2021-44228	Apache Solr 8.11	Upgrade log4j-core to 2.17.1; replace vulnerable jar in Solr lib dir
m4-002	DP	CVE-2021-45046	Log4j 2.15.0	Upgrade log4j-core/api to 2.17.1; 2.15 is still vulnerable
m4-003	DP	CVE-2021-45105	Log4j 2.16.0	Upgrade log4j-core to 2.17.1 to fix recursive-lookup DoS
m4-004	DP	CVE-2021-44832	Log4j 2.17.0	Upgrade log4j-core to 2.17.1 to fix JDBC appender RCE
m4-005	DP	CVE-2022-22965	Spring 5.3.17	Upgrade Spring Framework to 5.3.18/5.2.20
m4-006	DP	CVE-2022-22963	Spring Cloud Fn	Upgrade Spring Cloud Function to 3.2.3
m4-007	DP	CVE-2025-24813	Tomcat 9.0.98	Disable partial PUT (set readonly=true) or upgrade Tomcat
m4-008	DP	CVE-2021-41773	Apache httpd 2.4.49	Upgrade httpd to 2.4.51 or later
m4-009	DP	CVE-2021-42013	Apache httpd 2.4.50	Upgrade httpd to 2.4.51 to fix double-encoding bypass
m4-010	DP	CVE-2024-50379	Tomcat 11.0.1	Set default servlet readonly=true; upgrade Tomcat
m4-011	DP	CVE-2024-6387	OpenSSH	Upgrade openssh-server to 9.8p1 or later (regreSSHion)
m4-012	DP	CVE-2022-26134	Confluence 7.18	Upgrade to 7.18.1+; block /rest/api at WAF if upgrade delayed
m4-013	DP	CVE-2023-22515	Confluence 8.3	Upgrade to 8.3.3+ to fix admin account creation bypass
m4-014	DP	CVE-2023-7028	GitLab CE 16.7	Upgrade GitLab to 16.7.2+ to fix password-reset account takeover
m4-015	DP	CVE-2024-23897	Jenkins LTS	Upgrade Jenkins to LTS 2.426.3; disable CLI if not needed
m4-016	DP	CVE-2023-42793	TeamCity 2023.05	Upgrade TeamCity to 2023.05.4; block /app/rest/users
<i>Kernel and local privilege escalation (kernel-vm Vagrant required)</i>				
m4-017	AC	CVE-2021-4034	polkit	Upgrade policykit-1; chmod 0755 pkexec to remove SUID
m4-018	AC	CVE-2021-3156	sudo	Upgrade sudo to 1.9.5p2 or later to fix heap overflow

continued on next page

Table 10 continued

ID	Cat	CVE	Service	Canonical remediation
m4-019	AC	CVE-2022-0847	Linux kernel	Upgrade kernel to 5.16.11+, 5.15.26+, or 5.10.103+
m4-020	AC	CVE-2023-4911	glibc	Upgrade libc6 to fix GLIBC_TUNABLES buffer overflow
m4-021	AC	CVE-2023-2640	Linux kernel (Ubuntu)	Upgrade Ubuntu kernel past patched build (GameOver(lay))
m4-022	AC	CVE-2024-1086	Linux kernel	Upgrade kernel past nf_tables use-after-free fix
<i>Container and cloud-on-localhost</i>				
m4-023	DP	CVE-2024-21626	runc	Upgrade runc to 1.1.12 (Leaky Vessels)
m4-024	DP	CVE-2024-23651	BuildKit	Upgrade BuildKit to 0.12.5 to fix build-time container escape
m4-025	DP	CVE-2024-23652	BuildKit	Upgrade BuildKit to 0.12.5 to fix host-file delete/replace
m4-026	CH	N/A	Docker daemon	Remove docker.sock bind mount; drop -privileged; scope capabilities
m4-027	DP	CVE-2024-3094	XZ Utils	Downgrade liblzma to 5.4.x clean version; rebuild sshd
m4-028	AC	N/A	crAPI (REST API)	Validate video owner_id against authenticated user per request
m4-029	AC	N/A	DVGA (GraphQL)	Disable introspection; set DISABLE_INTROSPECTION=true
m4-030	AC	N/A	VAmPI (REST API)	Restart with VULNERABLE=0 to enable auth/authorization
m4-031	DP	CVE-2024-10924	WordPress 6.5	Upgrade Really Simple Security plugin to 9.1.2 or later
m4-032	DP	CVE-2023-48795	OpenSSH	Exclude Terrapin-affected ciphers/MACs in sshd_config
m4-033	CH	N/A	Docker daemon	Remove TCP listener from /etc/docker/daemon.json; use socket only
m4-034	CH	N/A	Docker daemon	Remove docker.sock bind mount from application compose file
m4-035	CH	N/A	Flask / IMDSv1	Add URL allow-list to /fetch; block 169.254.169.254
m4-036	CH	N/A	MinIO	Remove anonymous read policy from bucket via mc policy
m4-037	CH	N/A	Redis (ArgoCD)	Set requirepass with strong token in Redis config
m4-038	CH	N/A	k3s (kubelet)	Set -read-only-port=0 and -anonymous-auth=false
m4-039	CH	N/A	MongoDB 4.4	Start mongod with -auth; create admin account; enable authorization
m4-040	CH	N/A	MongoDB 4.4	Disable net.compression.compressors to block UDP amplification
m4-041	CH	N/A	CouchDB 3.1	Create admin account via CouchDB config API to exit Admin Party
m4-042	DP	CVE-2022-24706	CouchDB 3.2.1	Replace default Erlang cookie with strong random value in vm.args
m4-043	DP	CVE-2021-44521	Cassandra 3.11	Set authenticator: PasswordAuthenticator; disable UDF Groovy sandbox
m4-044	CH	N/A	Elasticsearch 7	Enable X-Pack security; set xpack.security.enabled: true
m4-045	DP	CVE-2019-20933	InfluxDB 1.7.6	Generate strong shared-secret; set in influxdb.conf [http] block

continued on next page

Table 10 continued

ID	Cat	CVE	Service	Canonical remediation
m4-046	DP	CVE-2021-34371	Neo4j 3.4	Disable Shell Server: set <code>dbms.shell.enabled=false</code> in <code>neo4j.conf</code>
m4-047	CH	N/A	RabbitMQ 3.11	Delete default guest user; create named user with strong password
m4-048	DP	CVE-2023-46604	ActiveMQ 5.17.5	Upgrade ActiveMQ to 5.15.16+/5.16.7+/5.17.6+ (OpenWire RCE)
m4-049	DP	CVE-2024-32114	ActiveMQ 6.1.0	Add <code>SecurityHandler</code> to <code>/api/jolokia</code> ; require authentication
m4-050	CH	N/A	Kafka 3.5	Switch listener to <code>SASL_PLAINTEXT</code> ; configure <code>SCRAM-SHA-256</code>
m4-051	CH	N/A	Mosquitto 2.0	Set <code>allow_anonymous false</code> in <code>mosquitto.conf</code> ; add password file
m4-052	DP	CVE-2022-24450	NATS 2.7.1	Add authorization block with per-user credentials in <code>nats-server.conf</code>
m4-053	DP	CVE-2023-33246	RocketMQ 5.1.0	Enable ACL in <code>broker.conf</code> : set <code>aclEnable=true</code>
m4-054	CH	N/A	BIND 9.18	Restrict recursion to trusted sources: set <code>allow-recursion { trusted; }</code>
m4-055	CH	N/A	BIND 9.18	Deny AXFR zone transfers: set <code>allow-transfer { none; }</code> globally
m4-056	CH	N/A	BIND 9.18	Enable DNSSEC validation: set <code>dnssec-validation auto</code> in <code>named.conf</code>
m4-057	CH	N/A	CoreDNS 1.11	Add ACL plugin to deny external sources; allow only trusted CIDRs
m4-058	CH	N/A	PowerDNS 4.8	Set strong randomly generated <code>api-key</code> (32+ chars) in <code>pdns.conf</code>
m4-059	CH	N/A	Unbound 1.21	Restrict recursion: set <code>access-control: allow_recursion to LAN</code> only
<i>Application and web-API surfaces</i>				
m4-060	DP	CVE-2023-51764	Postfix 3.7	Add <code>smtpd_forbid_bare_newline=yes</code> in <code>main.cf</code>
m4-061	DP	CVE-2023-51766	Exim 4.96	Disable <code>CHUNKING</code> advertisement: set <code>chunking_advertise_hosts=</code>
m4-062	CC	CVE-2019-11500	Dovecot 2.3.7	Require SSL/TLS for all connections; set <code>ssl_required=yes</code> in <code>conf</code>
m4-063	CC	CVE-2023-5631	Roundcube 1.6.3	Add <code>Content-Security-Policy</code> header blocking inline scripts
m4-064	CC	N/A	Modbus/TCP	Apply iptables allow-list restricting port 502 to authorized hosts
m4-065	CC	N/A	S7comm	Apply iptables restricting port 102 to engineering workstations only
m4-066	CC	N/A	IPMI BMC	Apply iptables restricting UDP 623 to management VLAN; disable Cipher 0
m4-067	CC	N/A	BACnet/IP	Apply iptables restricting UDP 47808 to BMS network segment

continued on next page

Table 10 continued

ID	Cat	CVE	Service	Canonical remediation
m4-068	CH	N/A	OPC-UA	Remove <code>SecurityPolicy=None</code> ; require Sign+Encrypt security mode
m4-069	CH	N/A	Docker	Remove <code>-cap-add SYS_PTRACE</code> ; drop container to minimal capabilities
m4-070	DP	CVE-2022-0492	Docker	Apply default <code>seccomp</code> profile; remove <code>-security-opt seccomp=unconfined</code>
m4-071	CH	N/A	Docker	Remove <code>-security-opt seccomp=unconfined</code> ; apply Docker default <code>seccomp</code>
m4-072	CH	N/A	Docker	Remove <code>-security-opt apparmor=unconfined</code> ; load default AppArmor profile
m4-073	CH	N/A	Docker Registry	Configure <code>htpasswd</code> authentication; set <code>REGISTRY_AUTH=htpasswd</code>
m4-074	CH	N/A	Docker image	Remove <code>ENV</code> secret instructions from Dockerfile; use Docker secrets
m4-075	CH	N/A	K8s RBAC	Delete over-privileged <code>ClusterRoleBinding</code> ; apply least-privilege RBAC
m4-076	CH	N/A	K8s networking	Apply <code>default-deny NetworkPolicy</code> to target namespace
m4-077	CH	N/A	K8s resources	Apply <code>LimitRange</code> to set default CPU/memory limits on pods
m4-078	CH	N/A	K8s service	Change <code>Service</code> type from <code>NodePort</code> to <code>ClusterIP</code>
m4-079	CH	N/A	K8s etcd	Create <code>EncryptionConfiguration</code> with <code>AES-CBC</code> key; restart <code>kube-apiserver</code>
m4-080	CH	N/A	AWS IAM	Delete over-permissive policy version; apply least-privilege IAM
m4-081	CH	N/A	AWS SSM	Delete plaintext parameters; recreate as <code>SecureString</code> with KMS key
m4-082	CH	N/A	AWS S3	Remove public bucket policy or restrict to specific IAM principals
m4-083	CH	N/A	AWS SNS	Update topic policy to restrict <code>SNS:Subscribe</code> to specific principals
m4-084	CH	N/A	AWS Lambda/IAM	Replace inline policy; restrict <code>iam:PassRole</code> to specific role ARN
m4-085	CH	N/A	HashiCorp Vault	Remove <code>-dev</code> flag; write production <code>vault.hcl</code> with persistent backend
m4-086	CH	N/A	HashiCorp Vault	Add seal stanza for Transit auto-unseal or cloud KMS in <code>vault.hcl</code>
m4-087	CH	N/A	HashiCorp Consul	Set <code>acl{ enabled=true, default_policy="deny"}</code> in <code>consul.hcl</code>
m4-088	CH	N/A	Docker image	Remove <code>ENV DB_PASSWORD</code> from Dockerfile; pass secret at runtime
m4-089	CH	N/A	Docker image	Remove <code>ARG SECRET_TOKEN</code> ; avoid writing secrets to image layers
m4-090	CH	N/A	Git repository	Use <code>git filter-repo</code> to rewrite history; rotate all leaked credentials
m4-091	CC	CVE-2021-43798	Grafana 8.3.0	Place <code>nginx</code> reverse proxy blocking <code>..</code> in <code>/public/plugins/</code> paths

continued on next page

Table 10 continued

ID	Cat	CVE	Service	Canonical remediation
m4-092	CH	N/A	Prometheus 2.40	Create <code>web.yml</code> with HTTP basic auth; start with <code>-web.config.file</code>
m4-093	CH	N/A	Jupyter Notebook	Set strong random token in <code>jupyter_notebook_config.py</code>
m4-094	DP	CVE-2023-27524	Apache Superset	Generate strong random <code>SECRET_KEY</code> ; set in <code>superset_config.py</code>
m4-095	DP	CVE-2020-17526	Apache Airflow	Generate strong random <code>secret_key</code> in <code>airflow.cfg</code> <code>webserver</code> section
m4-096	DP	CVE-2020-11978	Apache Airflow	Set <code>load_examples=False</code> in <code>airflow.cfg</code> ; remove example DAGs
m4-097	CC	CVE-2019-7609	Kibana 6.5.4	Disable Timelion plugin: add <code>timelion.enabled: false</code> in <code>kibana.yml</code>
m4-098	CH	N/A	Apache Druid	Set <code>druid.javascript.enabled=false</code> in <code>coordinator.runtime.properties</code>
m4-099	CH	N/A	Hadoop YARN	Enable ACL enforcement in <code>yarn-site.xml</code> : set <code>yarn.acl.enable=true</code>
m4-100	CC	CVE-2020-11651	SaltStack 3000	Block ports 4505/4506 to untrusted hosts via iptables; upgrade when feasible
m4-101	CC	CVE-2021-29441	Nacos 1.4.0	Deploy nginx proxy overwriting User-Agent to non-Nacos-Server value
m4-102	CC	CVE-2022-23131	Zabbix 5.0	Change default admin password; place nginx reverse proxy with auth
m4-103	CH	CVE-2024-47177	cups-browsed	Stop and disable cups-browsed; mask the systemd unit
m4-104	CC	CVE-2023-38646	Metabase 0.46.6	Switch DB backend from H2 to PostgreSQL; remove H2 JDBC driver access
m4-105	CC	CVE-2025-29927	Next.js	Deploy nginx reverse proxy stripping <code>x-middleware-subrequest</code> header
m4-106	DP	CVE-2020-13945	Apache APISIX	Change default admin API key to strong random secret in <code>config.yaml</code>
m4-107	CC	CVE-2025-32433	Erlang/OTP SSH	Block port at firewall; restrict to public-key auth and allowed users
m4-108	CC	CVE-2025-3248	Langflow	Block <code>/api/v1/validate/code</code> at nginx; enforce auth on API
m4-109	CC	CVE-2020-14882	WebLogic Console	Deploy nginx blocking <code>/console</code> paths; normalize URL encoding
m4-110	DP	CVE-2023-34152	ImageMagick	Enforce restrictive <code>policy.xml</code> ; block dangerous delegates (eps, mvg)
m4-111	NS	CVE-2018-1000115	Memcached 1.6	Start Memcached with <code>-U 0</code> to disable the UDP listener
m4-112	CC	CVE-2021-41773	Apache httpd 2.4.49	Install <code>mod_security</code> with OWASP CRS rules; block path-traversal patterns
m4-113	CC	CVE-2022-26134	Confluence (sim.)	Configure nginx reverse proxy to <code>localhost:8090</code> ; block <code>/rest/api/latest</code>

continued on next page

Table 10 continued

ID	Cat	CVE	Service	Canonical remediation
m4-114	CC	CVE-2022-46169	Cacti \leq 1.2.22	Configure nginx reverse proxy to localhost:8080; block remote_agent.php
m4-115	DP	CVE-2019-0193	Apache Solr 8.8	Disable DataImportHandler; remove /dataimport requestHandler from config
m4-116	CH	N/A	rsync daemon	Add auth users to each module; set secrets file; disable anon access
m4-117	CC	CVE-2026-31431	Linux kernel	Upgrade kernel to fix AF_ALG AEAD page-cache overwrite (Copy Fail)
<i>Active Directory VM scenarios (advn-01–advn-20)</i>				
advn-01	AC	CVE-2020-1472	Netlogon/DC	Set FullSecureChannelProtection=1; restart Netlogon (ZeroLogon)
advn-02	AC	CVE-2021-42278	Active Directory	Lower MachineAccountQuota to 0; apply KB5007247+ patches
advn-03	CC	N/A	Active Directory	Set strong password on svc_sql; restrict Kerberos to AES encryption
advn-04	CC	N/A	Active Directory	Disable DONT_REQ_PREAUTH on account; set strong Kerberos password
advn-05	CC	N/A	Active Directory	Remove unconstrained delegation from computer account
advn-06	AC	N/A	Active Directory	Remove DCSync rights (<i>Replicating Directory Changes</i>) from non-admin
advn-07	CH	N/A	ADCS	Remove ENROLLEE_SUPPLIES_SUBJECT flag from certificate template
advn-08	CH	N/A	ADCS	Remove Any Purpose / no-EKU certificate template from enrollment
advn-09	CH	N/A	ADCS	Remove Enrollment Agent capability from vulnerable template
advn-10	CH	N/A	ADCS	Disable EDITF_ATTRIBUTESUBJECTALTNAME2 flag on CA policy module
advn-11	CH	N/A	ADCS	Disable NTLM on IIS Web Enrollment; require HTTPS + Kerberos
advn-12	CC	N/A	Active Directory	Require LDAP signing: set LDAPServerIntegrity=2 in DC registry
advn-13	CC	N/A	Active Directory	Require SMB signing: set RequireSecuritySignature=1 for LanmanServer
advn-14	CC	N/A	Active Directory	Set LmCompatibilityLevel=5 to require NTLMv2; disable NTLMv1/LM
advn-15	NS	N/A	Active Directory	Disable LLMNR via Group Policy; disable NetBIOS over TCP/IP
advn-16	DP	CVE-2021-34527	Print spooler	Disable Print Spooler on DC; apply KB5005010 (PrintNightmare)
advn-17	CH	CVE-2021-36942	Active Directory	Apply ADV210003 mitigation; disable unneeded RPC interfaces
advn-18	AC	CVE-2014-1812	Group Policy	Remove cpassword XML from SYSVOL; rotate exposed account passwords
advn-19	AC	N/A	Active Directory	Deploy Microsoft LAPS; rotate local Administrator password on member hosts

continued on next page

Table 10 continued

ID	Cat	CVE	Service	Canonical remediation
advn-20	AC	N/A	Active Directory	Remove unauthorized ACE from AdminSDHolder; trigger SDProp to propagate

473 **C.7 hivestorm (16 scenarios)**

474 Hivestorm scenarios are free-roam: a single image plants several artifacts of mixed type, scored with weighted
 475 partial credit. Planted artifact identities (account names, binary paths, cron schedules, SUID filenames, registry
 476 key names) are randomized at build time from a per-scenario `roles.json`. The table below lists the 16 scenarios
 477 by platform and primary artifact mix.

Table 11: Per-scenario index for the `hivestorm` suite. Artifact mix codes: **B** backdoor account, **T** trojan path or malicious binary, **C** rogue cron entry, **S** SUID plant, **R** registry persistence (Windows), **D** domain object abuse (Active Directory).

ID	Platform	Artifact mix	Notes
hs-01	Debian 9	B, T, C, S	Rogue UID-0 account, SUID plant, trojan binary, reverse-shell cron
hs-02	Ubuntu 16.04	B, T, C, S	Samba share; rogue SSH authorized key; four planted artifact classes
hs-03	Windows 10	B, R, T	Rogue admin; HKLM Run autorun; trojan binary; scheduled-task persistence
hs-04	WinServer 2019	B, R	Rogue admin; SMB signing off; NTLMv1 regression; WDigest credential caching
hs-05	WinServer 2016	B, R, T	Rogue admin; SMBv1 re-enabled; LLMNR on; trojan binary; rogue Backup Operator
hs-06	Debian 9 (DB)	B, T	PostgreSQL; rogue DB user; Perl backdoor script planted under <code>/usr/local</code>
hs-07	Ubuntu 18.04	B, T	Samba; hidden user account; Perl backdoor script; rogue share
hs-08	WinServer (IIS)	B, R, T	IIS + PHP; rogue admin; scheduled-task persistence; cryptominer trojan
hs-09	Ubuntu (Nginx)	B, T, C	Nginx + phpBB; hidden user; socat trojan; rogue cron; prohibited package
hs-10	Ubuntu (PAM)	B	PAM faillock; hidden user; weak admin password; prohibited package planted
hs-11	WinServer (DC sim.)	B, R	Registry-planted DC misconfigs: LDAP signing, Kerberos, NTLMv1, WDigest
hs-12	CentOS 7 (LAMP)	B, T, C, S	LAMP; extra UID-0 user; trojan binary; rogue cron; SUID plant; rogue yum repo
hs-13	AD DC 2019	B, D, R	Real AD DS; rogue Domain Admin; Kerberoastable service; unconstrained delegation
hs-14	FreeBSD 13	B, T, S	Backdoor user; periodic drop; KLD kernel module; trojan; rogue pkg
hs-15	Docker host	B, T	Ubuntu Docker host; backdoor user; rogue container with sock-mount escape
hs-16	Ubuntu (Nginx+PHP-FPM)	B, T, C	PHP-FPM; backdoor user; webshell; trojan binary; rogue cron listener

478 **D Additional experiments**

479 **D.1 Naive heuristic baseline**

480 **Heuristic.** We run, on each scenario, the script `unattended-upgrades --minimal-upgrade-steps`
 481 followed by `systemctl daemon-reload` and a sweep of `systemctl restart` across every running service.
 482 No configuration files are edited, no permissions are changed, no accounts are added or removed, no firewall
 483 rules are added, no services are disabled, and no kernel reboots are performed. We then invoke the scenario's
 484 standard `verify.sh` and record exit 0 as a pass.

Suite	PASS	Total	Floor (%)
ccdc/	0	50	0.00
meta2/	0	40	0.00
vulnhub/	1	30	3.33
meta3/ubuntu/	0	19	0.00
meta3/windows/	0	21	0.00
meta4/ (Docker)	1	117	0.85
meta4/ad-vm/	0	20	0.00
hivestorm/	0	16	0.00
Overall	2	313	0.64

Table 12: Naive baseline (unattended-upgrades + service restart) pass rate per suite. Two scenarios pass: one is a Shellshock rebuild on Debian 11 whose stock bash ships already patched (vulnhub/scenario-20); the other is regreSSHion (meta4/scenario-11) where the unheld openssh-server package has a bookworm-security backport.

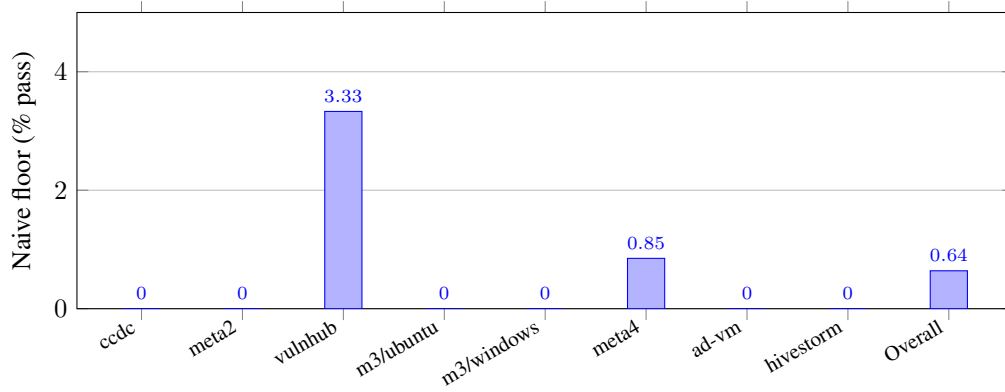


Figure 3: Naive heuristic baseline pass rate per suite. The floor is $\leq 1\%$ across every suite, confirming that the benchmark cannot be solved by blanket package upgrades and service restarts; meaningful progress requires the per-scenario configuration, permission, account, firewall, and compensating-control work that distinguishes a real remediation agent.

485 **Per-suite results.** Table 12 reports the naive baseline pass rate per suite; Figure 3 plots the same data. The
486 heuristic establishes a floor that any LLM agent must clear to demonstrate non-trivial remediation behaviour.

487 **Why the floor is so low.** Across all 313 scenarios, 311 fall outside what blanket package upgrades can
488 remediate. The dominant failure modes are: (i) configuration hardening (`sshd_config`, `nginx.conf`, `my.cnf`,
489 `pg_hba.conf`, `php.ini`, `sysctl`); (ii) access control (`chmod/chown`, `sudoers`, account/group management, back-
490 door account removal); (iii) explicit *Compensating Controls* scenarios that forbid upgrade by construction
491 (firewall scoping, loopback bind, WAF rule, removing an unsafe directive); (iv) vendored or source-built
492 components (`log4j` JARs, GitLab/Jenkins/Atlassian images, Drupal/WordPress trees) that the distro pack-
493 age manager does not own; (v) kernel-coupled scenarios where a container cannot upgrade its host kernel
494 (meta4 S19/S21/S22/S117); and (vi) Windows targets (all of `meta3/windows/`, `meta4/ad-vm/`, and sev-
495 eral `hivestorm/` scenarios) on which `apt/unattended-upgrades/systemctl` simply do not exist and the
496 heuristic is a no-op.

497 E Host requirements

498 `meta2/`. Linux host only. Requires a kernel booted with `vsyscall=emulate`, because the Ubuntu 8.04 base
499 image depends on the legacy `vsyscall` page, which was removed from upstream Linux in 5.18 and is disabled
500 in the Docker Desktop and WSL2 kernels on Windows and macOS. Scenarios `scenario-37` and `scenario-39`
501 require `-cap-add=NET_ADMIN`.

502 `meta3/windows/`. Windows 10 or 11 Professional or Enterprise, or Windows Server 2019 or later, with
503 Docker Desktop in Windows Containers mode or a native `dockerd`. Approximately 40 GB free disk. Hyper-V
504 isolation is the default for mixed build numbers.

505 meta4/kernel-vm/. Scenarios S21 and S22 ship a Vagrant virtual machine with a pinned kernel; replacing
506 the kernel invalidates the proof-of-concept.

507 **F verify.sh and JSON Lines scorer anatomy**

508 **F.1 Binary scorer (verify_sh_scorer)**

509 Every non-Hivestorm scenario ships a `verify.sh` with two mandatory test phases: a *PoC check* that confirms
510 the vulnerability is no longer exploitable, and a *regression check* that confirms the affected service is still running.
511 Exit 0 only if both pass.

512 **SSH empty-password login (CCDC scenario-02).** The script calls `sshd -T` for the effective post-
513 Match configuration, then performs a live login attempt via `sshpass` to catch configs that look correct on disk
514 but are overridden at runtime. The regression phase repeats the probe with a valid credential to ensure the
515 daemon was not simply disabled.

516 **Exposed installer pages (meta3 scenario-05).** `curl` drives both phases: `/drupal/install.php` and
517 `/phpmyadmin/setup/index.php` must return 403, while the application homepages must remain 200. A
518 blanket Apache shutdown satisfies the PoC but fails the regression.

519 **Apache path-traversal CVE (meta4 scenario-09).** Two remediation paths are accepted: upgrading
520 `httpd` to $\geq 2.4.51$, or blocking the double-encoded traversal string at the WAF/rewrite layer. The script resolves
521 version via `httpd -v`, falls through to a live exploit probe if the version is too old, and finally retries the root
522 endpoint up to six times before declaring the service broken.

523 **Exim4 open relay (VulnHub scenario-15).** If `dc_local_interfaces` still contains `0.0.0.0`, the script
524 accepts an `iptables DROP` rule on port 25 as an equivalent fix, avoiding a false failure for agents that mitigate
525 at the firewall layer rather than reconfiguring Exim.

526 **F.2 Weighted scorer (hivestorm_weighted_scorer)**

527 Hivestorm `verify.sh` scripts emit one JSON object per check on stdout:

```
528 {"check": "<id>", "weight": <int>, "pass": <bool>, "reason": "<str>"}
```

529 Service-survival probes additionally carry `"category": "service"` and `"weight": 0`. The scorer computes

$$S = \frac{\max(0, R - 10F_{\text{svc}})}{T},$$

530 where R is the sum of weights of passing non-service checks, T is the total non-service weight, and F_{svc} is the
531 count of failing service probes. This penalises agents that disable services to prevent planted processes from
532 respawning.

533 The planted artifacts fall into four classes.

534 **Class 1: Rogue accounts.** Each scenario plants a backdoor user with a `sudoers` drop-in; scenario-01 also
535 plants an extra UID-0 account. Checks are layered: user-database lookup (`getent passwd`), live SSH login
536 probe via `sshpass`, and a `sudo -n` behavioural test. All three must clear before the check emits `"pass": true`.

537 **Class 2: Persistence mechanisms.** Scenario-01 plants a cron drop-in and scans for the beacon IP across
538 all cron directories, catching agents that move rather than delete the entry. Scenario-14 (FreeBSD) plants a
539 `periodic` script and a `loader.conf KLD` entry; each is checked by file existence and a `grep` on the relevant
540 config.

541 **Class 3: Trojanized and SUID binaries.** The trojan check is behavioural rather than hash-based: the
542 script invokes the suspicious binary with `-h` and matches the `netcat` help output, so a clean replacement binary
543 passes even if its inode differs from the original. The SUID-plant check combines a mode-bit test (`-u $file`)
544 with an effective-UID probe executed as the legitimate admin account.

545 **Class 4: Covert listeners and webshells.** Listener checks use both the socket-table (`ss/sockstat`) and
546 a direct TCP-connect probe (`/dev/tcp`), handling the race between a killed process and a supervising unit
547 that immediately respawns it. Scenario-14 adds an `nginx` webshell check (`shell.php` in the document root).
548 Service-survival probes for SSH and the web server carry `"category": "service"` and appear at the end of
549 every Hivestorm script; a failure there subtracts 10 points from the raw score regardless of how many other
550 checks passed.

551 G Datasheet

552 We provide a datasheet for SysRepair-Bench following the framework of Gebru et al. [3]. We adopt the
553 conventional question-and-answer format and group questions under the seven canonical headings.

554 G.1 Motivation

555 **For what purpose was the dataset created?** SysRepair-Bench was created to evaluate AI agents
556 on the task of remediating real-world system vulnerabilities while preserving service availability. This task
557 is the defensive counterpart of offensive security benchmarks such as Cybench [26], NYU CTF Bench [18],
558 and CVE-Bench [28], and it is distinct from source-level program repair benchmarks such as SWE-bench [6].
559 To our knowledge, no prior public benchmark evaluates the joint constraint of closing a vulnerability on a
560 running system and not breaking the service the vulnerability protects, which is the dual constraint under which
561 production incident responders work.

562 **Who created the dataset and on behalf of which entity?** Anonymized for review

563 **Who funded the creation of the dataset?** Anonymized for review

564 **Any other comments?** Anonymized for review

565 G.2 Composition

566 **What do the instances represent?** Each instance represents a vulnerable system, materialized as a
567 Dockerized container or a Vagrant virtual machine, together with a machine-readable threat description and
568 a deterministic verification harness. Container scenarios are reproducible via `docker build` from a pinned
569 `Dockerfile`; kernel-coupled scenarios in `meta4 ship` a Vagrant virtual machine with a pinned kernel.

570 **How many instances are there in total?** SysRepair-Bench contains 313 scenarios in the released
571 distribution: 297 binary-scored scenarios across the `ccdc`, `meta2`, `vulnhub`, `meta3/ubuntu`, `meta3/windows`,
572 and `meta4` suites, and 16 partial-credit scenarios in the `hivestorm` suite. The full per-suite breakdown is given
573 in Table 1 and the per-scenario index in Appendix C.

574 **Does the dataset contain all possible instances or is it a sample of instances from a larger set?**
575 The dataset is a sample. The `ccdc`, `meta2`, `vulnhub`, and `meta3` suites cover documented findings in their
576 respective source materials, decomposed at the per-vulnerability granularity and de-duplicated where the source
577 flagged the same underlying issue under multiple identifiers. The `meta4` suite is deliberately a curated sample of
578 modern Common Vulnerabilities and Exposures rather than an exhaustive catalogue, with selection priorities
579 described in Section 4: alignment with OWASP Top 10 [13], the MITRE CWE Top 25 [9], and the CISA Known
580 Exploited Vulnerabilities catalog [2], weighted toward high-impact and operationally common classes (language
581 and library, kernel and local privilege escalation, container and cloud-on-localhost, application and web-API
582 surfaces). The `hivestorm` suite mirrors the platform mix of the source competition.

583 **What data does each instance consist of?** Binary-suite scenarios are directories with three files: a
584 `Dockerfile` that builds the vulnerable container deterministically; a `threat.md` file containing severity, CVSS
585 v3.1 score where applicable, Common Vulnerabilities and Exposures identifier where one exists, the affected
586 service, and the canonical remediation steps; and a `verify.sh` script returning exit code 0 if the vulnerability
587 is closed and the service still works and exit code 1 otherwise. Hivestorm scenarios replace `verify.sh` with
588 a JSON Lines scorer that emits one weighted check per planted artifact and a per-build `roles.json` that
589 randomizes artifact identities.

590 **Is there a label or target associated with each instance?** Yes. For the 297 binary scenarios the label is
591 the conjunction of three deterministic checks performed by `verify.sh`: the proof-of-concept no longer succeeds,
592 every previously operational service remains operational, and (for the 39 compensating-control scenarios) the
593 required mitigation is in place. For the 16 Hivestorm scenarios the label is a continuous score in $[0, 1]$ obtained
594 by summing weighted per-artifact checks normalized to the unit interval.

595 **Is any information missing from individual instances?** A subset of `ccdc` scenarios document miscon-
596 figurations rather than disclosed Common Vulnerabilities and Exposures, and therefore have no CVE identifier
597 in `threat.md`. Where this is the case the CVE field is set to N/A and the threat description is written from the
598 upstream hardening guidance.

599 **Are relationships between individual instances made explicit?** Yes, through suite membership,
600 remediation category (Access Control, Configuration Hardening, Dependency and Package Management,
601 Network Security and Firewall Policy, Compensating Controls), affected service or application type, and CVSS
602 severity bucket. These attributes appear in `threat.md` and are summarized in Table 1 and Appendix A.

603 **Are there recommended data splits?** SysRepair-Bench is an evaluation-only benchmark, so there is no
604 train, validation, or test split in the supervised-learning sense. Two evaluation conditions are defined: black-box
605 (only the running container is exposed) and report-informed (the `threat.md` description is provided as context).

606 **Are there any errors, sources of noise, or redundancies in the dataset?** We are not aware of any
607 errata at submission time. Errata, when discovered, will be tracked in `CHANGELOG.md` and addressed in versioned
608 releases (see the maintenance subsection below). The Hivestorm randomization is an intentional source of
609 stochasticity that varies artifact identities at build time without changing the underlying remediation difficulty.

610 **Is the dataset self-contained, or does it link to or rely on external resources?** Each scenario builds
611 locally from a self-contained `Dockerfile` or `Vagrantfile`. Build artifacts pull from upstream Docker Hub and
612 Linux distribution repositories at build time, and where the underlying vulnerability is tied to a specific upstream
613 version (for example `httpd:2.4.49` for CVE-2021-41773) the scenario pins to that version. Pins are listed
614 in the per-suite manifests in Appendix C. We recommend that users mirror upstream images locally before
615 evaluation runs to guard against upstream image deletion or repository unavailability.

616 **Does the dataset contain data that might be considered confidential?** No. All scenarios are
617 constructed from public material: published competition hardening checklists, the public OpenVAS scan of
618 Metasploitable 2, public VulnHub images, the public Rapid7 Metasploitable 3 cookbook, public Common
619 Vulnerabilities and Exposures records, and the public Hivestorm scoring methodology.

620 **Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening,
621 or might otherwise cause anxiety?** Each scenario contains exploitation primitives by construction, since
622 remediation requires a verifiable exploit to fail. These primitives are identical in nature to those that appear in
623 offensive security benchmarks and the public Common Vulnerabilities and Exposures advisory record. They are
624 embedded in isolated containers or virtual machines and are framed in a defensive context. The dataset does not
625 contain content directed at any person or group.

626 **Does the dataset relate to people?** No. Scenarios target software artifacts: services, packages, kernel
627 features, configuration files, and accounts on isolated systems. Account names that appear in scenarios (for
628 example `root`, `vagrant`, `admin`) are conventional system identifiers, not references to real individuals.

629 **Does the dataset identify any subpopulations, or contain data that might allow identification of
630 individuals, or contain sensitive data?** No. Not applicable.

631 G.3 Collection process

632 **How was the data associated with each instance acquired?** Each suite was constructed by manual
633 scenario authoring, drawing on the following public sources, listed by suite. The `ccdc` suite decomposes
634 published Collegiate Cyber Defense Competition blue-team hardening materials [10] into per-vulnerability
635 containers. The `meta2` suite is seeded by the OpenVAS scan of the canonical Metasploitable 2 image [16],
636 with the scan output included for auditability under `openvas-scan-reports/`. The `vulnhub` suite is built
637 from per-virtual-machine analysis of fourteen community VulnHub images [21]. The `meta3/ubuntu` and
638 `meta3/windows` suites vendor portions of the Rapid7 Metasploitable 3 Chef cookbook [17] under its BSD-3-
639 Clause license, with full acknowledgment in `meta3/ubuntu/shared/UPSTREAM_LICENSE`. The `meta4` suite
640 was authored by us against public Common Vulnerabilities and Exposures advisories, vendor security bulletins,
641 and proof-of-concept disclosures. The `hivestorm` scoring methodology follows the public Hivestorm scoring
642 engine [5]; planted artifacts and randomization logic were authored by us.

643 **What mechanisms or procedures were used to collect the data?** Scenarios were authored, built with
644 Docker or Vagrant, and validated by the authors through a two-stage check: a positive check that the scenario
645 reproduces the documented exploit on the unhardened image, and a negative check that `verify.sh` returns exit
646 0 after the canonical remediation in `threat.md` is applied. Findings from the OpenVAS scan of Metasploitable
647 2 were re-run against our reproduced container to confirm fidelity to the upstream image.

648 **If the dataset is a sample from a larger set, what was the sampling strategy?** For `ccdc`, `meta2`,
649 and `meta3`, the strategy is exhaustive coverage of distinct findings in the source material with de-duplication.
650 For `vulnhub`, the strategy is per-image decomposition, with each VulnHub image contributing between one and

651 four scenarios depending on the number of independent vulnerabilities it exposes. For `meta4`, the strategy is
652 purposive sampling along the four coverage axes described in Section 4, prioritizing high-impact, operationally
653 common, and post-cutoff Common Vulnerabilities and Exposures. For `hivestorm`, the strategy is platform-mix
654 parity with the source competition, covering Linux distributions (Debian, Ubuntu, CentOS, FreeBSD), Windows
655 Server Core, and Active Directory.

656 **Who was involved in the data collection process?** Anonymized for review

657 **Over what timeframe was the data collected?** Anonymized for review

658 **Were any ethical review processes conducted?** SysRepair-Bench contains no human-subjects data, no
659 personally identifiable information, and no surveillance or behavioral data; institutional review board review is
660 not applicable. Dual-use considerations are discussed in Section 8: the benchmark contains no novel exploit
661 code, every scenario builds from material that already exists in indexed public form, all scenarios run in isolation,
662 and the benchmark rewards closing vulnerabilities while penalizing service disruption.

663 **G.4 Preprocessing, cleaning, and labeling**

664 **Was any preprocessing, cleaning, or labeling of the data done?** Yes. Preprocessing differs by
665 suite. For `ccdc`, the source hardening checklists, which are written for full-system competition images, were
666 decomposed into per-vulnerability container scenarios so that each scenario isolates a single misconfiguration;
667 this prevents cross-contamination of remediation actions across scenarios. For `meta2`, the OpenVAS scan output
668 was de-duplicated by service and finding type, then each finding was rebuilt as a per-vulnerability container
669 reproducing the issue. For `vulnhub`, each upstream image was analyzed for distinct vulnerabilities, each of
670 which was isolated into its own container. For `meta3`, scenarios were ported from the upstream Rapid7 Chef
671 cookbook and re-keyed to the `Dockerfile` plus `threat.md` plus `verify.sh` format. For `meta4`, scenarios
672 were authored from CVE disclosures and vendor advisories. All scenarios were labeled with a remediation
673 category (Access Control, Configuration Hardening, Dependency and Package Management, Network Security
674 and Firewall Policy, Compensating Controls) and a CVSS severity bucket where the underlying CVE provides
675 one.

676 **Was the raw data saved in addition to the preprocessed data?** Yes. The OpenVAS scan re-
677 ports that drove the `meta2` and `meta3/ubuntu` scenario scopes are included in `openvas-scan-reports/`
678 for auditability. The upstream Rapid7 Chef cookbook is preserved alongside its license in
679 `meta3/ubuntu/shared/UPSTREAM_LICENSE`.

680 **Is the software that was used to preprocess, clean, or label the data available?** Yes. Per-suite
681 build scripts, the Inspect AI [20] harness, and the Hivestorm scorer reference implementation are distributed
682 with the benchmark.

683 **G.5 Uses**

684 **Has the dataset been used for any tasks already?** Yes. The baseline experiments reported in Section 7
685 of this paper are the first use, with two families covering three models (MiniMax-M2.7, Qwen3.5-35B-A3B, and
686 Qwen3.5-9B) under a ReAct [25] solver in black-box and report-informed conditions.

687 **Is there a repository that links to any or all papers or systems that use the dataset?** Anonymized
688 for review

689 **What other tasks could the dataset be used for?** The dataset is suited to several downstream tasks
690 beyond agent evaluation. (a) Curriculum design for defensive agent training, where scenarios provide graded
691 difficulty across remediation categories. (b) Studies of agent scaffolding, where the harness allows controlled
692 substitution of solvers, prompts, and tool affordances. (c) Studies of regression risk in automated remediation,
693 where the dual-objective scoring isolates the service-availability constraint as a first-class evaluation axis.

694 **Is there anything about the composition of the dataset or the way it was collected and prepro-
695 cessed that might impact future uses?** Two considerations are material. First, public security material,
696 including the source corpora that seed `ccdc`, `meta2`, `vulnhub`, and `meta3`, is highly indexed and has likely
697 appeared in the pretraining data of large models. Second, scenarios reproduce vulnerabilities by construction
698 and include exploit primitives. Researchers using the benchmark to train defensive agents should be aware that
699 the exploit primitives are also present in the training signal.

700 **Are there tasks for which the dataset should not be used?** SysRepair-Bench should not be used
701 as a training corpus for the same agent that is then evaluated on it. It should not be used to evaluate offensive
702 capability, because while the underlying scenarios overlap with offensive benchmarks, the scoring rewards
703 remediation rather than exploitation. It should not be used to evaluate live cloud remediation, since scenarios run
704 in isolation.

705 NeurIPS Paper Checklist

706 1. Claims

707 Question: Do the main claims made in the abstract and introduction accurately reflect the paper’s
708 contributions and scope?

709 Answer: [Yes]

710 Justification: The abstract and Section 1 enumerate six concrete contributions (313 binary-scored
711 scenarios across six suites, dual-objective scoring extended by compensating-control adequacy, the
712 Meta4 modern-CVE suite, the Hivestorm partial-credit and randomization track, an Inspect AI
713 harness, and baselines on four models in zero-knowledge and one-day variants with two contamination
714 ablations), each delivered by a corresponding section (§3–§7).

715 Guidelines:

- 716 • The answer [N/A] means that the abstract and introduction do not include the claims made in the
717 paper.
- 718 • The abstract and/or introduction should clearly state the claims made, including the contributions
719 made in the paper and important assumptions and limitations. A [No] or [N/A] answer to this
720 question will not be perceived well by the reviewers.
- 721 • The claims made should match theoretical and experimental results, and reflect how much the
722 results can be expected to generalize to other settings.
- 723 • It is fine to include aspirational goals as motivation as long as it is clear that these goals are not
724 attained by the paper.

725 2. Limitations

726 Question: Does the paper discuss the limitations of the work performed by the authors?

727 Answer: [Yes]

728 Justification: Section 8 contains explicit *Scope*, *Contamination*, and *Limitations* paragraphs disclosing
729 Linux skew, the minority share of Windows and Active Directory scenarios, and the deliberate
730 exclusion of source-code modification, zero-days without remediation, hardware/firmware issues, and
731 live cloud-IAM evaluation.

732 Guidelines:

- 733 • The answer [N/A] means that the paper has no limitation while the answer [No] means that the
734 paper has limitations, but those are not discussed in the paper.
- 735 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 736 • The paper should point out any strong assumptions and how robust the results are to violations of
737 these assumptions (e.g., independence assumptions, noiseless settings, model well-specification,
738 asymptotic approximations only holding locally). The authors should reflect on how these
739 assumptions might be violated in practice and what the implications would be.
- 740 • The authors should reflect on the scope of the claims made, e.g., if the approach was only tested
741 on a few datasets or with a few runs. In general, empirical results often depend on implicit
742 assumptions, which should be articulated.
- 743 • The authors should reflect on the factors that influence the performance of the approach. For
744 example, a facial recognition algorithm may perform poorly when image resolution is low or
745 images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide
746 closed captions for online lectures because it fails to handle technical jargon.
- 747 • The authors should discuss the computational efficiency of the proposed algorithms and how
748 they scale with dataset size.
- 749 • If applicable, the authors should discuss possible limitations of their approach to address problems
750 of privacy and fairness.
- 751 • While the authors might fear that complete honesty about limitations might be used by reviewers
752 as grounds for rejection, a worse outcome might be that reviewers discover limitations that
753 aren’t acknowledged in the paper. The authors should use their best judgment and recognize
754 that individual actions in favor of transparency play an important role in developing norms that
755 preserve the integrity of the community. Reviewers will be specifically instructed to not penalize
756 honesty concerning limitations.

757 3. Theory assumptions and proofs

758 Question: For each theoretical result, does the paper provide the full set of assumptions and a complete
759 (and correct) proof?

760 Answer: [N/A]

761 Justification: The paper presents a benchmark and empirical baselines and does not state or rely on
762 formal theorems.

763 Guidelines:

- 764 • The answer [N/A] means that the paper does not include theoretical results.
- 765 • All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- 766 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 767 • The proofs can either appear in the main paper or the supplemental material, but if they appear in
768 the supplemental material, the authors are encouraged to provide a short proof sketch to provide
769 intuition.
- 770 • Inversely, any informal proof provided in the core of the paper should be complemented by
771 formal proofs provided in appendix or supplemental material.
- 772 • Theorems and Lemmas that the proof relies upon should be properly referenced.

773 4. Experimental result reproducibility

774 Question: Does the paper fully disclose all the information needed to reproduce the main experimental
775 results of the paper to the extent that it affects the main claims and/or conclusions of the paper
776 (regardless of whether the code and data are provided or not)?

777 Answer: [Yes]

778 Justification: Each scenario builds deterministically from a versioned `Dockerfile` or `Vagrantfile`
779 (§3); the harness configuration (per-scenario time limit 1800 s, per-command `bash` timeout 180 s,
780 verification timeout 3000 s, `bash-only` tool surface) is specified in §6; models, seeds, solver, and the
781 two knowledge variants are listed in §7; and host requirements (e.g., `vsyscall=emulate` for `meta2/`
782 and the pinned-kernel Vagrant VM for `meta4/kernel-vm/`) are documented in Appendix E.

783 Guidelines:

- 784 • The answer [N/A] means that the paper does not include experiments.
- 785 • If the paper includes experiments, a [No] answer to this question will not be perceived well by
786 the reviewers: Making the paper reproducible is important, regardless of whether the code and
787 data are provided or not.
- 788 • If the contribution is a dataset and/or model, the authors should describe the steps taken to make
789 their results reproducible or verifiable.
- 790 • Depending on the contribution, reproducibility can be accomplished in various ways. For
791 example, if the contribution is a novel architecture, describing the architecture fully might suffice,
792 or if the contribution is a specific model and empirical evaluation, it may be necessary to either
793 make it possible for others to replicate the model with the same dataset, or provide access to
794 the model. In general, releasing code and data is often one good way to accomplish this, but
795 reproducibility can also be provided via detailed instructions for how to replicate the results,
796 access to a hosted model (e.g., in the case of a large language model), releasing of a model
797 checkpoint, or other means that are appropriate to the research performed.
- 798 • While NeurIPS does not require releasing code, the conference does require all submissions
799 to provide some reasonable avenue for reproducibility, which may depend on the nature of the
800 contribution. For example
 - 801 (a) If the contribution is primarily a new algorithm, the paper should make it clear how to
802 reproduce that algorithm.
 - 803 (b) If the contribution is primarily a new model architecture, the paper should describe the
804 architecture clearly and fully.
 - 805 (c) If the contribution is a new model (e.g., a large language model), then there should either be
806 a way to access this model for reproducing the results or a way to reproduce the model (e.g.,
807 with an open-source dataset or instructions for how to construct the dataset).
 - 808 (d) We recognize that reproducibility may be tricky in some cases, in which case authors are
809 welcome to describe the particular way they provide for reproducibility. In the case of
810 closed-source models, it may be that access to the model is limited in some way (e.g.,
811 to registered users), but it should be possible for other researchers to have some path to
812 reproducing or verifying the results.

813 5. Open access to data and code

814 Question: Does the paper provide open access to the data and code, with sufficient instructions to
815 faithfully reproduce the main experimental results, as described in supplemental material?

816 Answer: [Yes]

817 Justification: Scenarios, scan reports, the grading harness, and the evaluation code are released under an
818 open license at the anonymous URL given in the abstract ([https://anonymous.4open.science/
819 r/sysrepair-bench-425F/](https://anonymous.4open.science/r/sysrepair-bench-425F/)); per-suite license details and the BSD-3-Clause acknowledgment for
820 the vendored Rapid7 Chef cookbook will be finalized in Appendix.

821 Guidelines:

- 822 • The answer [N/A] means that paper does not include experiments requiring code.
- 823 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/public/
824 guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 825 • While we encourage the release of code and data, we understand that this might not be possible,
826 so [No] is an acceptable answer. Papers cannot be rejected simply for not including code, unless
827 this is central to the contribution (e.g., for a new open-source benchmark).
- 828 • The instructions should contain the exact command and environment needed to run to reproduce
829 the results. See the NeurIPS code and data submission guidelines ([https://neurips.cc/
830 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 831 • The authors should provide instructions on data access and preparation, including how to access
832 the raw data, preprocessed data, intermediate data, and generated data, etc.
- 833 • The authors should provide scripts to reproduce all experimental results for the new proposed
834 method and baselines. If only a subset of experiments are reproducible, they should state which
835 ones are omitted from the script and why.
- 836 • At submission time, to preserve anonymity, the authors should release anonymized versions (if
837 applicable).
- 838 • Providing as much information as possible in supplemental material (appended to the paper) is
839 recommended, but including URLs to data and code is permitted.

840 6. Experimental setting/details

841 Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters,
842 how they were chosen, type of optimizer) necessary to understand the results?

843 Answer: [Yes]

844 Justification: No model training is performed; baselines use off-the-shelf checkpoints (MiniMax 2.7,
845 Qwen3.6-35B-A3B, Qwen3.5-9B) under a single ReAct scaffold, with inference settings, timeouts,
846 and the tool surface given in §6, the zero-knowledge versus one-day prompt construction described in
847 the same section.

848 Guidelines:

- 849 • The answer [N/A] means that the paper does not include experiments.
- 850 • The experimental setting should be presented in the core of the paper to a level of detail that is
851 necessary to appreciate the results and make sense of them.
- 852 • The full details can be provided either with the code, in appendix, or as supplemental material.

853 7. Experiment statistical significance

854 Question: Does the paper report error bars suitably and correctly defined or other appropriate informa-
855 tion about the statistical significance of the experiments?

856 Answer: [Yes]

857 Justification: We report success@1 and success@5 over 5 random seeds (§7, following the protocol
858 of Zhu et al. [28]) and report the mean with standard error for the Hivestorm partial-credit track in
859 Tables 1.

860 Guidelines:

- 861 • The answer [N/A] means that the paper does not include experiments.
- 862 • The authors should answer [Yes] if the results are accompanied by error bars, confidence intervals,
863 or statistical significance tests, at least for the experiments that support the main claims of the
864 paper.
- 865 • The factors of variability that the error bars are capturing should be clearly stated (for example,
866 train/test split, initialization, random drawing of some parameter, or overall run with given
867 experimental conditions).
- 868 • The method for calculating the error bars should be explained (closed form formula, call to a
869 library function, bootstrap, etc.)
- 870 • The assumptions made should be given (e.g., Normally distributed errors).
- 871 • It should be clear whether the error bar is the standard deviation or the standard error of the
872 mean.

- 873 • It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report
874 a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is
875 not verified.
- 876 • For asymmetric distributions, the authors should be careful not to show in tables or figures
877 symmetric error bars that would yield results that are out of range (e.g., negative error rates).
- 878 • If error bars are reported in tables or plots, the authors should explain in the text how they were
879 calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

881 Question: For each experiment, does the paper provide sufficient information on the computer
882 resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

883 Answer: [Yes]

884 Justification: Host hardware and OS requirements for the harness (Docker, Vagrant,
885 the `vsyscall=emulate` kernel option for `meta2/`, Windows-Containers configuration for
886 `meta3/windows/`, and the pinned-kernel VM for `meta4/kernel-vm/`) are documented in Ap-
887 pendix E, and per-scenario wall-clock and token-cost statistics are reported in Table 2; model inference
888 for the open-weight baselines was served via vendor APIs and required no local accelerators.

889 Guidelines:

- 890 • The answer [N/A] means that the paper does not include experiments.
- 891 • The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud
892 provider, including relevant memory and storage.
- 893 • The paper should provide the amount of compute required for each of the individual experimental
894 runs as well as estimate the total compute.
- 895 • The paper should disclose whether the full research project required more compute than the
896 experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into
897 the paper).

9. Code of ethics

899 Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code
900 of Ethics <https://neurips.cc/public/EthicsGuidelines>?

901 Answer: [Yes]

902 Justification: The benchmark uses publicly available red-team materials and intentionally vulnerable
903 artifacts that already exist in indexed form (§8, *Dual-use*); no human subjects, no personal data, and
904 no novel exploit code are introduced, and all execution is sandboxed in containers or virtual machines
905 with no default network exposure.

906 Guidelines:

- 907 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of Ethics.
- 908 • If the authors answer [No], they should explain the special circumstances that require a deviation
909 from the Code of Ethics.
- 910 • The authors should make sure to preserve anonymity (e.g., if there is a special consideration due
911 to laws or regulations in their jurisdiction).

10. Broader impacts

913 Question: Does the paper discuss both potential positive societal impacts and negative societal impacts
914 of the work performed?

915 Answer: [Yes]

916 Justification: Section 8 (*Dual-use*) frames SysRepair-Bench as a defensive counterweight to existing
917 offensive-security benchmarks (rewarding closure of vulnerabilities and preserving service availability
918 while penalizing destructive collateral) and addresses the negative-impact surface, namely that better
919 defensive evaluation could indirectly inform attackers, by noting that all underlying vulnerability
920 material is already public and that the benchmark contains no novel exploit code.

921 Guidelines:

- 922 • The answer [N/A] means that there is no societal impact of the work performed.
- 923 • If the authors answer [N/A] or [No], they should explain why their work has no societal impact
924 or why the paper does not address societal impact.
- 925 • Examples of negative societal impacts include potential malicious or unintended uses (e.g.,
926 disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deploy-
927 ment of technologies that could make decisions that unfairly impact specific groups), privacy
928 considerations, and security considerations.

- 929
- 930
- 931
- 932
- 933
- 934
- 935
- 936
- 937
- 938
- 939
- 940
- 941
- 942
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate Deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

943 11. Safeguards

944 Question: Does the paper describe safeguards that have been put in place for responsible release of
945 data or models that have a high risk for misuse (e.g., pre-trained language models, image generators,
946 or scraped datasets)?

947 Answer: [Yes]

948 Justification: Scenarios run in isolated Docker containers or Vagrant VMs with no default network
949 exposure (§8), the grading harness is separated from agent code so that leaderboard submissions
950 cannot tamper with scoring, and the per-build identity randomization on Hivestorm doubles as a
951 misuse-resistance feature because memorized artifact names confer no advantage and the suite cannot
952 be repurposed as a fixed exploitation recipe.

953 Guidelines:

- 954
- 955
- 956
- 957
- 958
- 959
- 960
- 961
- The answer [N/A] means that the paper poses no such risks.
 - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

962 12. Licenses for existing assets

963 Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper,
964 properly credited and are the license and terms of use explicitly mentioned and properly respected?

965 Answer: [Yes]

966 Justification: Upstream sources (Metasploitable 2, Metasploitable 3, VulnHub, the CCDC hardening
967 materials, and the Hivestorm scoring engine) are cited in §2 and §3; the BSD-3-Clause acknowledgment
968 for the vendored Rapid7 Chef cookbook.

969 Guidelines:

- 970
- 971
- 972
- 973
- 974
- 975
- 976
- 977
- 978
- 979
- 980
- 981
- 982
- The answer [N/A] means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
 - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
 - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
 - If this information is not available online, the authors are encouraged to reach out to the asset's creators.

983 13. New assets

984 Question: Are new assets introduced in the paper well documented and is the documentation provided
985 alongside the assets?

986 Answer: [Yes]

987 Justification: Two new assets are introduced, the Meta4 suite (§4) and the Hivestorm partial-credit
988 track (§5), with each scenario shipping a `Dockerfile` (or `Vagrantfile`), a `threat.md` metadata
989 file (CVSS, CVE, affected service, canonical remediation), and a `verify.sh` grader (or JSON Lines
990 scorer for Hivestorm); per-suite indices, annotated `verify.sh` walkthroughs, and the full datasheet
991 (per Gebru et al. [3]) are provided in Appendices C, F, and G.

992 Guidelines:

- 993 • The answer [N/A] means that the paper does not release new assets.
- 994 • Researchers should communicate the details of the dataset/code/model as part of their sub-
995 missions via structured templates. This includes details about training, license, limitations,
996 etc.
- 997 • The paper should discuss whether and how consent was obtained from people whose asset is
998 used.
- 999 • At submission time, remember to anonymize your assets (if applicable). You can either create an
1000 anonymized URL or include an anonymized zip file.

1001 **14. Crowdsourcing and research with human subjects**

1002 Question: For crowdsourcing experiments and research with human subjects, does the paper include
1003 the full text of instructions given to participants and screenshots, if applicable, as well as details about
1004 compensation (if any)?

1005 Answer: [N/A]

1006 Justification: No crowdsourcing was used; the only human-in-the-loop step is the failure-trajectory
1007 annotation by two of the paper’s authors (§7, Table 2), and no external participants were recruited or
1008 compensated.

1009 Guidelines:

- 1010 • The answer [N/A] means that the paper does not involve crowdsourcing nor research with human
1011 subjects.
- 1012 • Including this information in the supplemental material is fine, but if the main contribution of the
1013 paper involves human subjects, then as much detail as possible should be included in the main
1014 paper.
- 1015 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other
1016 labor should be paid at least the minimum wage in the country of the data collector.

1017 **15. Institutional review board (IRB) approvals or equivalent for research with human subjects**

1018 Question: Does the paper describe potential risks incurred by study participants, whether such
1019 risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an
1020 equivalent approval/review based on the requirements of your country or institution) were obtained?

1021 Answer: [N/A]

1022 Justification: The work involves no human subjects, no participant recruitment, and no collection of
1023 personal data, so IRB review was not required.

1024 Guidelines:

- 1025 • The answer [N/A] means that the paper does not involve crowdsourcing nor research with human
1026 subjects.
- 1027 • Depending on the country in which research is conducted, IRB approval (or equivalent) may be
1028 required for any human subjects research. If you obtained IRB approval, you should clearly state
1029 this in the paper.
- 1030 • We recognize that the procedures for this may vary significantly between institutions and
1031 locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for
1032 their institution.
- 1033 • For initial submissions, do not include any information that would break anonymity (if applica-
1034 ble), such as the institution conducting the review.

1035 **16. Declaration of LLM usage**

1036 Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard
1037 component of the core methods in this research? Note that if the LLM is used only for writing, editing,
1038 or formatting purposes and does *not* impact the core methodology, scientific rigor, or originality of the
1039 research, declaration is not required.

1040 Answer: [Yes]

1041 Justification: LLMs are the central object of evaluation rather than an authoring aid; the models,
1042 scaffold, and inference configuration are disclosed in §7 (primary baseline MiniMax 2.7 paired with a
1043 ReAct solver, with comparison runs on DeepSeek V4 Pro, Qwen3.6-35B-A3B, and Qwen3.5-9B), and
1044 the rationale for ReAct as the chosen scaffold is given in §2.

1045 Guidelines:

- 1046 • The answer [N/A] means that the core method development in this research does not involve
1047 LLMs as any important, original, or non-standard components.
- 1048 • Please refer to our LLM policy in the NeurIPS handbook for what should or should not be
1049 described.