

## 981 A CORRECT Evaluation Framework

982 The CORRECT evaluation introduces rationale-based assessment to verify whether models identify  
 983 vulnerabilities for the correct reasons: The rationale refers to the model's step-by-step reasoning  
 984 process and results that lead to its vulnerability detection conclusion. CORRECT evaluates this  
 985 rationale against ground-truth vulnerability information (CVE descriptions, patches, and commit  
 986 messages) to determine correctness: The framework operates as follows:

- 987 (1) **For vulnerable code:** The model must correctly identify the ground-truth vulnerability  
     in its rationale. If the rationale includes the actual vulnerability cause or key elements in  
     vulnerabilities, it is deemed correct; otherwise, the detection is considered a false negative  
     despite potentially correct binary labeling.
- 988 (2) **For patched code:** Two evaluation modes are employed:
  - 989 • *Lenient Mode*: Accepts the model's results if it either correctly identifies the code as  
     non-vulnerable, or if it reports a vulnerability but the rationale does not reference the  
     original (now patched) vulnerability.
  - 990 • *Strict Mode*: This mode extends Lenient Mode by adding an explicit iterative check.  
     When the model flags a patched function as vulnerable but does not reference the  
     ground-truth vulnerability, the framework instructs the model to disregard the previously  
     reported issues and re-evaluate whether the function still contains a vulnerability.  
     If during this iterative process the model again reports the ground-truth vulnerabil-  
     ity, the case is marked as a false positive, since the model fails to recognize that the  
     vulnerability has already been fixed.

1000 These scenarios are summarized in Table 4, which contrasts the evaluation outcomes of Lenient and  
 1001 Strict modes across different prediction cases. In practice, this rationale assessment is implemented  
 1002 by another LLM acting as a judge, which takes the model's explanation as input and decides whether  
 1003 it correctly reflects the ground-truth vulnerability.

1004 Table 4. Evaluation outcomes under Lenient and Strict modes for different prediction scenarios.

1009 <b>Ground Truth</b>	<b>Prediction</b>	<b>Rationale</b>	<b>Lenient</b>	<b>Strict</b>
1010 Vulnerable	Vulnerable	Correct	TP	TP
1011 Vulnerable	Vulnerable	Incorrect	FN	FN
1012 Vulnerable	Non-vulnerable	–	FN	FN
1013 Patched	Non-vulnerable	–	TN	TN
1014 Patched	Vulnerable	References original vuln	FP	FP
1015 Patched	Vulnerable	Other issues only	TN	Feedback→TN/FP

1016 **Choice of Evaluation Mode.** While CORRECT provides both Lenient and Strict modes, we  
 1017 adopt Lenient Mode in our evaluation based on empirical evidence and practical considerations.

1018 Empirical results in the CORRECT paper, including extensive statistics reported in its appendix,  
 1019 show that the benefits of *Strict Mode* are marginal: only 10.3% of cases are modified after the  
 1020 first feedback round, dropping to 3.6% in round 2, and merely 1.5% by round 4. This diminishing  
 1021 return suggests that the iterative refinement process yields small improvements while significantly  
 1022 increasing computational costs.

1023 In addition, applying *Strict Mode* to existing methods is particularly problematic. Training-based  
 1024 approaches such as **ReVD**, and multi-agent frameworks such as **VulTrial**, rely on carefully designed  
 1025 instructions. Forcing these models to ignore previously reported vulnerabilities and repeatedly

1030 re-analyze the code disrupts their prompt structures, often leading to poor performance. Supporting  
 1031 such feedback loops would require retraining models from scratch with feedback-aware objectives,  
 1032 which introduces substantial overhead in both data construction and training optimization.

1033 Given these issues, we adopt *Lenient Mode* for all evaluations. This ensures a fair comparison  
 1034 across baselines while still keeping evaluation rigor through rationale correctness assessment.  
 1035

## 1036 **B Prompt Template.**

### 1037 **Prompt: System Prompt for General Specifications**

1039 A structured threat modeling analysis process where security experts conduct systematic security analysis based on  
 1040 provided information. The expert must:

1041 **1. Understand Code Context** (within <understand> tags)

1042 Thoroughly analyze and describe the system context without revealing the vulnerability itself:

1043 **System Identification**

- 1044 • **What system:** Clearly identify the software system, library, or application
- 1045 • **Domain/Subsystem:** Specify the particular domain or subsystem where the code operates
- 1046 • **Module/Component:** Identify the specific module, component, or functional unit

1047 **Functional Analysis**

- 1048 • **Core functionality:** Describe what this system/module is designed to do in detail: 1. 2. 3.

1049 **2. Security Domain Classification** (within <classification> tags)

1050 Classify vulnerabilities according to 10 core security domains:

1051 **Core Security Domains:**

- 1052 (1) **MEM:** Memory Safety [Buffer errors, pointer issues, use-after-free, allocation problems, etc.]
- 1053 (2) **STATE:** State Management [Inconsistent states, object lifecycle, concurrency issues, etc.]
- 1054 (3) **INPUT:** Input Validation [Parsing logic, data validation, type checking, encoding, etc.]
- 1055 (4) **LOGIC:** Program Logic [Arithmetic errors, type confusion, logical mistakes, etc.]
- 1056 (5) **SEC:** Security Features [Authentication, cryptography, permissions, policy enforcement]
- 1057 (6) **IO:** I/O Interaction [Filesystem operations, networking, device interaction, etc.]
- 1058 (7) **CONF:** Configuration Environment [Configuration parsing, environmental variables, etc.]
- 1059 (8) **TIMING:** Timing & Concurrency [Race conditions, synchronization issues, TOCTOU, etc.]
- 1060 (9) **PROTOCOL:** Protocol Communication [Message parsing/formatting, session handling, etc.]
- (10) **HARDWARE:** Hardware & Low-level [Low-level interfaces, architectural specifics, etc.]

1061 **3. Security Specification** (within <spec> tags)

1062 Security Specification helps understand how vulnerable code violates developer's original constraints and how patches  
 1063 implement fixes.

1064 **Example:**

1065 **Input Information**

- 1066 • **Repository:** ksmbd
- **Commit Message:** ksmbd: Fix dangling pointer in krb\_authenticate
- **CVE Description:** In the Linux kernel, the following vulnerability has been resolved: ksmbd: Fix dangling pointer  
 in krb\_authenticate...
- **CWE Type:** CWE-416 (Use After Free)

1067 **Code Diff:**

```
1068 -if (sess->state == SMB2_SESSION_VALID)
1069 +if (sess->state == SMB2_SESSION_VALID) {
1070     ksmbd_free_user(sess->user);
1071     sess->user = NULL;
1072 }
```

1073 **Expected Output Format:**

1074 <understand>

1075 **System Identification**

- 1076 • **What system:** ksmbd - in-kernel SMB server implementation for Linux
- **Domain/Subsystem:** SMB/CIFS network file sharing protocol implementation

```

1079 • Module/Component: Kernel component receives SMB requests, uses netlink IPC...
1080 </understand>
1081 <classification>
1082 <primary>MEM.LIFECYCLE</primary>
1083 <tags>[STATE.CONSTANCY, SEC.AUTHENTICATION, PROTOCOL.SMB]</tags>
1084 <reasoning>The root cause is the failure to manage the lifecycle...</reasoning>
1085 </classification>
1086 <spec>HS-MEM-001: Pointer release operations require atomic cleanup with immediate
1087 nullification</spec>
1088 - Reasoning: Dangling pointer vulnerability → freed but not nullified → atomic release-nullification prevents use-after-free


---


1089 Current Analysis Target:
1090 Repository: {repository}
1091 Commit Message: {commit_message}
1092 CVE Description: {cve_description}
1093 CWE Type: {cwe_type}
1094 Vulnerable Code:
1095 {vuln}
1096 Solution:
1097 {fixed}
1098 Please conduct analysis following the above format.

```

### Prompt: Detailed Vulnerability Cases in General Specifications

```

1100 A structured threat modeling analysis process where security experts conduct systematic security analysis based on
1101 provided information.
1102 Analysis Framework
1103 1. System Understanding (provided context)
1104 {understand}
1105 2. Security Specifications (provided rules)
1106 {specification}
1107 3. System-Level Threat Modeling (within <model> tags)
1108 Analyze vulnerability at system design level:
1109 • Trust Boundaries: Identify where system components transition between trusted/untrusted states
1110 • Attack Surfaces: Focus on realistic attack vectors that led to this specific vulnerability
1111 • CWE Analysis: Trace complete vulnerability chain (e.g., initial CWE-X triggers subsequent CWE-Y, where at least
1112 one matches: {cwe_type})
1113 4. Code-Level Analysis
1114 Vulnerability Context (within <vuln> tags)
1115 Provide a granular, narrative explanation of the vulnerability:
1116 (1) Entry Point & Preconditions: Describe how the attack is initiated and what system state is required
1117 (2) Vulnerable Code Path Analysis: Step-by-step trace of execution flow, naming key functions and variables.
1118 Pinpoint The Flaw and its Consequence
1119 (3) Specification Violation Mapping: Link code path steps to specific HS- specifications they violate
1120 Fix Implementation (within <solution> tags)
1121 Explain how the patch enforces security specifications:
1122 • Specific code changes and their security impact
1123 • How fixes restore compliance with violated specifications


---


1124 Example Output Format:
1125 <model>
1126 • trust_boundaries: User-Kernel boundary during SMB2 session setup; Intra-kernel function contract violation

```

```

1128 • attack_surfaces: Malicious SMB2 SESSION_SETUP request; Error path exploitation
1129 • cwe_analysis: Primary CWE-416 (Use After Free) enabled by state management violation
1130 </model>
1131 <vuln>
1132 (1) Entry Point: Privileged user sends Netlink message with crafted CIPSOV4 tags
1133 (2) Code Path: Loop processes tags → The Flaw: Off-by-one error in bounds check → Consequence: Stack buffer
1134 overflow
1135 (3) Violations: HS-MEM-001 (incorrect bounds check), HS-STATE-002 (incomplete initialization)
1136 </vuln>
1137 <solution>
1138 Change 1: Bounds Check Correction
1139 -if (iter > CIPSO_V4_TAG_MAXCNT)
1140 +if (iter >= CIPSO_V4_TAG_MAXCNT)
1141 Compliance: Changes exclusive to inclusive comparison, preventing array overflow
1142 Change 2: Complete Array Initialization
1143 -doi_def->tags[iter] = CIPSO_V4_TAG_INVALID;
1144 +while (iter < CIPSO_V4_TAG_MAXCNT)
1145 +    doi_def->tags[iter++] = CIPSO_V4_TAG_INVALID;
1146 Compliance: Ensures all array elements initialized to safe values
1147 </solution>
1148
Input Information:
1149 • CVE: {cve_description}
1150 • CWE: {cwe_type}
1151 • Commit: {commit_message}
1152 • Vulnerable Code: {vuln}
1153 • Fixed Code: {fixed}
1154 • Code Context: {code_context}
1155 Please conduct analysis following the above framework.

```

## Prompt: VulInstruct Knowledge Scoring Mechanism

You are a security expert. Please evaluate the relevance between the following code and VulInstruct vulnerability cases.

### Target Code

{code\_snippet}

### VulInstruct Cases to Evaluate

{chr(10).join(cases\_for\_evaluation)}

Please score the relevance of each case to the target code (1-10 points):

### Scoring Criteria:

- **10 points:** Highly relevant, vulnerability type, trigger conditions, and code patterns are almost identical
- **8-9 points:** Strong relevance, main vulnerability features are similar, can provide valuable reference
- **6-7 points:** Moderate relevance, some features are similar, has certain reference value
- **4-5 points:** Weak relevance, only few similarities
- **1-3 points:** Very low relevance, basically no reference value

Please strictly follow the HTML format for output:

```

1169 <vulinstruct_evaluation>
1170 <case_1_score>6</case_1_score>
1171 <case_1_reasoning>Scoring reason</case_1_reasoning>
1172 <case_2_score>8</case_2_score>
1173 <case_2_reasoning>Scoring reason</case_2_reasoning>
1174 ...

```

```
1177 </vulninstruct_evaluation>
```

## 1179 Prompt: Domain-specific Specification Extraction

1181 You are a security expert. Analyze these related vulnerabilities and extract reusable security specifications.

### 1182 Related Historical Vulnerabilities

```
{chr(10).join(nvd_descriptions)}
```

### 1183 Task: Extract Attack-Derived Specifications

1184 For each vulnerability pattern you identify:

- 1185 (1) Identify the recurring attack mechanism across these CVEs
- 1186 (2) Convert it to a positive security specification that would prevent such attacks
- 1187 (3) Format as defensive requirements developers must implement

### 1188 Output Format:

```
1189 <attack_specifications>
1190   <specification_1>
1191     <attack_pattern>
1192       Description of recurring attack mechanism
1193       in cve-xxx and cve-xxx in detail
1194     </attack_pattern>
1195     <defensive_spec>
1196       AS-DOMAIN-001: Security rule that describes
1197       the code behavior that prevents this attack
1198     </defensive_spec>
1199     <implementation_hint>
1200       Specific checks or validations needed
1201     </implementation_hint>
1202   </specification_1>
1203   <specification_2>...</specification_2>
1204 </attack_specifications>
```

## 1203 Prompt: Vulnerability Detection

1205 You are a senior code security expert. Please perform systematic multi-layer security analysis on the following code.

1206 **Analysis Mode:** [Determined by knowledge relevance scoring]

- 1207 • **Autonomous Analysis:** Low relevance with knowledge base, perform independent analysis
- 1208 • **Knowledge-Assisted:** High relevance knowledge filtered through LLM evaluation as reference

### 1209 Input Components:

- 1210 • Code Snippet: {code\_snippet}
- 1211 • Code Context: {code\_context}
- 1212 • LLM-filtered Security Knowledge: {selected\_knowledge}

### 1213 Multi-Layer Vulnerability Analysis Framework

1214 1. **Surface Symptom Analysis** Identify direct suspicious operations.

1215 2. **Root Cause Investigation** Trace deeper causes that give rise to the surface symptoms, focusing on data/control flow, completeness of input validation, adequacy of error handling, and potential attacker exploitation paths.

1216 3. **Architectural & Contextual Analysis** Examine broader design-level factors and domain-specific assumptions in the application logic.

1217 **Comprehensive Security Assessment.** Based on the above LLM-filtered Security Knowledge and three-layer analysis mode framework, please provide your professional judgment:

- 1218 (i) Analysis Process: [Please describe your three-layer analysis process in detail, including discovered issues and reasoning chains]
- 1219 (ii) Key Findings: [List the most important security findings]
- 1220 (iii) Final Conclusion:

1221 Please strictly follow the format below for output:

### 1222 Output Format:

1224

1226 Table 5. Failure analysis of baseline methods. Left: iteration distribution of Vul-RAG. Right: categorized  
 1227 failure cases of ReVD.

1228

Iteration Range	Ratio	Vuln.	Secure
1–3	53.1%	12.5%	87.5%
4–6	22.6%	18.0%	82.0%
7–9	6.2%	27.5%	72.5%
10 (final)	5.4%	50.0%	50.0%
10 (no decision)	12.7%	0.0%	100%

1236 (a) Vul-RAG iteration distribution where Vuln and  
 1237 Secure represents the final binary prediction.

Failure Type	Ratio
Zero reasoning	4%
Wrong vulnerability type	30%
Similar type confusion	20%
Mechanism misinterpretation	14%
Over-generalization	26%
Other errors	6%

1238 (b) ReVD failure categories.

1239

1240

```

1241 <vulnerability_assessment>\\
1242 Please strictly follow the format below for output:
1243   <has_vulnerability>yes/no</has_vulnerability>
1244   <confidence>0-1</confidence>
1245   <suspected_root_cause>Core findings summary</suspected_root_cause>
1246 </vulnerability_assessment>
```

#### Format Description:

- has\_vulnerability: “yes” or “no”
- confidence: Confidence level between 0.0 and 1.0
- If a fixing solution has been applied, you may judge “no”
- Focus on analysis quality, avoid over-sensitivity

1251

1252

## C Failure Analysis in Vul-RAG and ReVD

1253 We conduct a targeted analysis of the key performance factors in three representative approaches:  
 1254 the retrieval mechanism in Vul-RAG, the reasoning capability in ReVD, and the specifications in  
 1255 VulInstruct.

1256 *Failure Analysis of Vul-RAG.* Vul-RAG conducts vulnerability detection in an iterative retrieval  
 1257 manner. Given a target code snippet, the model first retrieves the top- $k$  most relevant vulnerabil-  
 1258 ity-patch cases ( $k = 10$  in our experiments). For each case, it checks (i) whether the same type of  
 1259 vulnerability exists in the target code and (ii) whether the corresponding patch has already been  
 1260 applied. The outcome is a binary signal: (1, 0) indicates a vulnerability without a patch (*classified as  
 1261 vulnerable*), (0, 1) indicates a patch without vulnerability (*classified as secure*), while (0, 0) or (1, 1)  
 1262 are inconclusive and trigger the next iteration. The process continues until a conclusive decision  
 1263 is made; if no decisive signal is found after all iterations, the output is treated as *no decision*. As  
 1264 shown in Table 5, over half of the samples (53.1%) terminate within the first three iterations, while  
 1265 12.7% end without any decision.

1266 As shown in Table 5(a), over half of the samples (53.1%) terminate within three iterations, while  
 1267 12.7% never reach a decision. A common failure is prematurely classifying samples as secure once a  
 1268 single patch match is observed, even if other vulnerabilities remain. This reliance on cross-project  
 1269 case matching proves unreliable: few vulnerabilities recur in exactly the same form, resulting in  
 1270 only 3.4% reasoning correctness under CORRECT.

1271 *Failure Analysis of ReVD.* We manually analyzed 50 vulnerable cases where ReVD produced the  
 1272 correct binary label but failed in reasoning. Four recurring error types emerged (Table 5(b)): (i)  
 1273 **Zero reasoning** where no explanation was provided, (ii) **Incorrect vulnerability categorization**

1274

1275 with either wrong CWE families or confusion between closely related types, (iii) **Incomplete**  
 1276 **mechanism understanding** where the type was identified but the related code or the exploitation  
 1277 results was wrong, and (iv) **Over-generalization** where only vague statements like code quality  
 1278 concerns replaced specific root causes. Although ReVD leverages large-scale reasoning data, its  
 1279 generated explanations often lack the precision and depth which highlights the key limitation of  
 1280 reasoning-distillation approaches: they improve surface-level consistency but do not guarantee  
 1281 faithful identification of vulnerability root causes.

## 1282 D Manual analysis of successful cases in VulnInstruct

1284 We analyze VulnInstruct's specification-guided approach using DeepSeek-R1, examining how speci-  
 1285 fications enhance vulnerability detection. The method achieves 16.7% improvement in detecting  
 1286 actual vulnerabilities ( $\text{FN} \rightarrow \text{TP}$ ) and 8.0% reduction in false positives ( $\text{FP} \rightarrow \text{TN}$ ). To understand these  
 1287 improvements, we randomly sampled 10 cases from each category and performed a manual analysis,  
 1288 identifying four distinct repair mechanisms through which specifications enhance detection capa-  
 1289 bilities. Table 6 shows four recurring *repair mechanisms*: (i) **Missing Security Dimension (20%**  
 1290 **of cases)** which occurs when models lack awareness of entire vulnerabilities. For example, in CVE-  
 1291 2021-37848, the model only checked `strcmp` for logical errors. Our specification addressed this  
 1292 by mandating constant-time comparisons for security-sensitive operations. (ii) **Domain-Specific**  
 1293 **Blindness (30% of cases)** occurs when models detect potential defects but underestimate their  
 1294 severity in specific contexts. For instance, in CVE-2022-24214, the model identified an integer  
 1295 overflow in DNS code but classified it as low risk, failing to recognize that DNS TXT records are  
 1296 attacker-controlled and high-risk vectors. Our domain specification AS-DOMAIN-1 provided this  
 1297 essential context, upgrading the assessment from "possible issue" to "high-severity vulnerability."  
 1298 (iii) **Deep Reasoning Enhancement (25% of cases)** improves models' ability to connect isolated  
 1299 risks into complete analysis of vulnerability mechanism. We provide a detailed case study in Fig-  
 1300 ure 6. (iv) **Secure Pattern Validation (25% of cases, exclusively  $\text{FP} \rightarrow \text{TN}$ )** enables models to  
 1301 recognize secure implementations and avoid false positives. For example, in CVE-2022-21654, the  
 1302 model mistakenly flagged the use of `EVP_sha256` as insecure. Our specification HS-CRYPTO-003  
 1303 confirmed this as correct cryptographic practice, reinforcing that knowledge of secure patterns is  
 1304 as critical as vulnerability detection.

1306 Table 6. Repair mechanisms by which specifications improve vulnerability detection

1308 Repair Mechanism	FN → TP	FP → TN	Total	Primary Knowledge Sources
1309 Missing Security Dimension	4	0	4 (20%)	General + Domain-specific
1310 Domain-Specific Blindness	3	3	6 (30%)	Domain-specific + Detailed cases
1311 Deep Reasoning Enhancement	3	2	5 (25%)	Domain-specific + Detailed cases
1312 Secure Pattern Validation	0	5	5 (25%)	General + Detailed cases

## 1314 E Using VulnInstruct Finding real-world vulnerability

1315 **Case Study: From CVE-2021-32056 to a New Access Control Bypass.** We provide a detailed  
 1316 case study to illustrate how our specification-guided auditing framework can facilitate the discov-  
 1317 ery of new real-world vulnerabilities. Starting from CVE-2021-32056 in the Cyrus IMAP Server  
 1318 (`cyrusimap/cyrus-imapd`), a randomly selected vulnerability case from the CORRECT dataset, we  
 1319 successfully identified and reported a previously unknown high-severity flaw in the same codebase.  
 1320 CVE-2021-32056 allowed authenticated users to bypass intended access restrictions on server  
 1321

1324 annotations, leading to potential replication stalls and service disruptions. The root cause was an  
 1325 improperly scoped permission check in `imap/annotate.c`, where the `maywrite` check was nested  
 1326 inside a conditional block and skipped when the mailbox pointer was `NULL`. From this case, our  
 1327 framework distilled three reusable security specifications, including the specification HS-SEC-001  
 1328 (annotation write operations must always enforce strict privilege-based access control).

1329 Guided by this specification, we constructed an auditing agent workflow that simulated the  
 1330 workflow of a manual security audit: starting from the known flaw, generalizing its pattern, and  
 1331 searching the codebase for other instances where security checks might be incorrectly scoped.

1332 Our workflow first prompted a LLM to generate candidate `git grep` commands that could reveal  
 1333 similar patterns elsewhere in the repository. We employed Gemini-2.5-Pro, whose relatively large  
 1334 context window allowed us to supply extracted code fragments without requiring a more elaborate  
 1335 dynamic windowing design. The generated commands reflected the model's reasoning about how  
 1336 the extracted specification might be violated in other parts of the codebase. In particular, the model  
 1337 focused on core database operations such as `store` and `delete`, which are security-critical under  
 1338 HS-SEC-001, and combined them with contextual information from the repository's application  
 1339 context (e.g., session management, mailbox operations). In doing so, the model hypothesized  
 1340 potential validation points where access control checks might be inconsistently applied. The next  
 1341 stage of our workflow was to feed the retrieved code snippets back into the model for analysis under  
 1342 the extracted specifications. This allowed the model to reason about whether the conditional checks  
 1343 in each match were relevant to security enforcement. Among four generated queries produced,  
 1344 one was especially effective, as shown below:

```
1345 git grep -p --all-match \
1346   -e 'if (mailbox)' \
1347   -e 'cyrusdb_store\|cyrusdb_delete' \
1348   -- '*.c'
```

1349 This query rediscovered the already patched `write_entry` function in `annotate.c`, thereby val-  
 1350 iduating the search strategy, while simultaneously surfacing additional candidate sites. Furthermore,  
 1351 the model dismiss some false positives, such as matches in `imap/tls.c`, where the conditionals  
 1352 guarded resource management logic rather than access control, and therefore did not violate  
 1353 HS-SEC-001. At the same time, the model highlighted `imap/mboxlist.c` as a high-risk location,  
 1354 noting that several conditional branches in this file surrounded critical database operations such as  
 1355 `cyrusdb_store` and `cyrusdb_delete`. To further investigate, we selected this candidate and applied  
 1356 our Automatic Context Extraction Tool to retrieve the surrounding function bodies together with a  
 1357 depth-3 call chain, ensuring that the model could reason about how access control checks were prop-  
 1358 agated across related functions. Given this enriched context, the model conducted a more detailed  
 1359 analysis and identified problematic control-flow paths. Guided by this assessment, we performed a  
 1360 closer inspection of the extracted functions—`mboxlist_update`, `mboxlist_update_entry_full`,  
 1361 and `mboxlist_renamemailbox`—which ultimately led to the discovery of a severe privilege bypass.  
 1362 The most critical finding emerged in `mboxlist_renamemailbox`, a function with complex multi-  
 1363 path control flow. Here, we discovered that a `goto` statement transferred execution directly to the  
 1364 database update section, bypassing any privilege checks:  
 1365

```
1366 if (mbentry->mbtype & MBTYPE_INTERMEDIATE) {
1367   // ... destination checks ...
1368   goto dbupdate; // BYPASSES permission check!
1369 }
1370 myrights = cyrus_acl_myrights(auth_state, mbentry->acl);
1371 if (!isadmin && !(myrights & ACL_DELETEMBOX)) {
1372 }
```

```

1373     return IMAP_PERMISSION_DENIED;
1374 }

```

1375 As a consequence, an authenticated user without ACL\_DELETEMBX rights could still rename  
 1376 intermediate mailboxes. This represents a direct violation of HS-SEC-001: security checks were  
 1377 present but misplaced, enabling a privilege bypass. Conceptually, this flaw mirrors CVE-2021-  
 1378 32056—the same fundamental security specification was violated—but here the error manifested  
 1379 through control-flow misordering (`goto`) rather than conditional scoping.  
 1380

1381 The potential consequences were significant: unauthorized mailbox renaming could disrupt  
 1382 shared mailbox hierarchies, shift shared folders into private namespaces, and cause denial-of-service  
 1383 for legitimate users. We responsibly disclosed the issue to the Cyrus IMAP development team, who  
 1384 confirmed the vulnerability, reproduced it via integration tests, and patched it by relocating the  
 1385 permission check before the special-case handling logic.  
 1386

## F Automatic Context Extraction Tool

1387 We implement a unified Commit URL parser that normalizes repository and patch commit metadata  
 1388 across heterogeneous hosts. Using a small set of hand-written matching rules, the parser recogni-  
 1389 zes both modern and legacy interfaces (e.g., GitHub, GitLab, Bitbucket, as well as cgit/gitweb  
 1390 deployments such as kernel.org and GNU Savannah), and deterministically maps each commit URL  
 1391 to a canonical triple `{repo_name, commit_hash, clone_url}`. For example, kernel.org cgit links are  
 1392 remapped to stable GitHub mirrors to enable uniform downstream handling.  
 1393

1394 Building on this, we resolve—when available via the host’s API—the canonical parent repository  
 1395 for each commit, thereby obtaining the repositories for both the patched and vulnerable versions.  
 1396 Following CORRECT, we derive target functions from the change lines and then use `joern` and  
 1397 `cflow` to extract the four context types described in Section 4.1.1. For large repositories (e.g.,  
 1398 Linux, TensorFlow), we further augment the default flow with lightweight subsystem/component  
 1399 identification and module-boundary detection, guided by directory structure cues together with  
 1400 staged filtering (includes, calls, symbol-to-file mapping) and shallow depth limits. These heuristics  
 1401 make CPG construction incremental and tractable, while avoiding uncontrolled expansion.  
 1402

1403 Building on this, we resolve—when available via the host’s API—the canonical parent reposi-  
 1404 tory for each commit, thereby obtaining the repositories for both the patched and vulnerable  
 1405 versions. Following CORRECT, we derive target functions from the change lines and then use  
 1406 `joern` and `cflow` to extract the four context types described in Section 4.1.1. For large repositories  
 1407 (e.g., Linux, TensorFlow), we further augment the default flow with lightweight subsystem/com-  
 1408 ponent identification and module-boundary detection, *seeded by a small set of manually curated  
 1409 anchors* (Linux: `net/fs/kernel`; TensorFlow: `core/python/compiler`) and simple boundary rules  
 1410 (`find_module_root`, `is_module_boundary`). The analysis employs staged filtering (header-include  
 1411 analysis → call-graph edges → symbol-to-file mapping) and adaptive file selection (targets+module,  
 1412 callers/callees, include neighbors, paired `.h/.c`), together with conservative caps (up to 3,000 files;  
 1413 depth ≤ 3) and graceful timeouts/recovery. These heuristics keep CPG construction incremental  
 1414 and tractable while preventing uncontrolled expansion.  
 1415

1416 Our final preprocessed dataset provides substantially richer context than PrimeVul, which only  
 1417 supplies the full file in which a vulnerable function resides. In contrast, our dataset captures the four  
 1418 distinct types of context described in Section 4.1.1, thereby enabling more fine-grained program  
 1419 analysis and downstream tasks.  
 1420

1421  
 1418  
 1419  
 1420  
 1421