

Do Frontier LLMs Truly Understand Smart Contract Vulnerabilities?

Anonymous ACL submission

Abstract

Frontier large language models achieve state-of-the-art performance on code understanding benchmarks, yet their capacity for smart contract security remains unclear. Can they genuinely reason about vulnerabilities, or merely pattern-match against memorized exploits? We introduce **BlockBench**, a benchmark designed to answer this question, revealing heterogeneous capabilities. While some models demonstrate robust semantic understanding, most exhibit substantial surface pattern dependence.

1 Introduction

Smart contract vulnerabilities represent one of the most costly security challenges in modern computing. As shown in Figure 1, cryptocurrency theft has resulted in over \$14 billion in losses since 2020, with 2025 already reaching \$3.4 billion, the highest since the 2022 peak (Chainalysis, 2025). The Bybit breach alone accounted for \$1.5 billion, while the Cetus protocol lost \$223 million in minutes due to a single overflow vulnerability (Tsentsura, 2025).

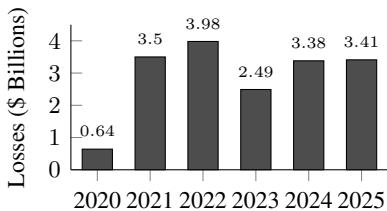


Figure 1: Annual cryptocurrency theft losses (2020–2025). Data from Chainalysis.

Meanwhile, large language models have achieved remarkable success on programming tasks. Frontier models now pass technical interviews, generate production code, and identify bugs across diverse codebases. This raises a natural question: *can these models apply similar expertise to blockchain security?* And if they can, *are they genuinely reasoning about vulnerabilities, or merely pattern-matching against memorized examples?*

This distinction matters. A model that has memoized the 2016 DAO reentrancy attack may flag similar patterns, yet fail when the same flaw appears in unfamiliar syntax. We introduce **BlockBench**, a benchmark designed to answer this question. Our contributions include:

1. **BlockBench**, comprising 263 Solidity vulnerability samples with systematic contamination control and gold standard examples from recent professional security audits.
2. **Composite evaluation metrics** distinguishing genuine understanding from memorization, validated through multi-configuration sensitivity analysis (Spearman’s $\rho=1.000$).
3. **Systematic assessment** revealing 58% best-case detection on mixed samples collapsing to 20% on uncontaminated professional audits, exposing heterogeneous robustness and accuracy-understanding gaps across models.

2 Related Work

2.1 Traditional Analysis Tools

Early approaches to smart contract vulnerability detection relied on static analysis and symbolic execution. Tools such as Slither (Feist et al., 2019), Mythril (Mueller, 2017), and Security (Tsankov et al., 2018) demonstrated strong precision on syntactically well-defined vulnerability classes. Durieux et al. (2020) conducted a comprehensive evaluation of nine such tools across 47,587 Ethereum contracts, revealing consistent performance on reentrancy and integer overflow detection, yet persistent struggles with vulnerabilities requiring semantic reasoning about contract logic. Ghaleb and Pattabiraman (2020) corroborated these findings, observing that rule-based approaches fundamentally cannot capture the contextual nuances that distinguish exploitable flaws from benign code patterns.

070 2.2 LLM-Based Approaches

071 Large language models introduced new possibilities
072 for bridging this semantic gap. Initial investigations
073 by Chen et al. (2023) explored prompting
074 strategies for vulnerability detection, achieving de-
075 tection rates near 40% while noting pronounced
076 sensitivity to superficial features such as variable
077 naming conventions. GPTScan (Sun et al., 2024b)
078 combined GPT-4 with program analysis to achieve
079 78% precision on logic vulnerabilities, leveraging
080 static analysis to validate LLM-generated can-
081 didates. Sun et al. (2024a) introduced retrieval-
082 augmented approaches that provide models with
083 relevant vulnerability descriptions, substantially
084 improving detection performance. Multi-agent
085 architectures emerged as another direction, with
086 systems like GPTLens (Hu et al., 2023) employ-
087 ing auditor-critic pairs to enhance analytical con-
088 sistency. Fine-tuning on domain-specific corpora
089 has yielded incremental gains, though performance
090 characteristically plateaus below the 85% threshold
091 regardless of training scale.

092 2.3 Pattern Recognition Versus 093 Understanding

094 Beneath these encouraging metrics lies a more
095 fundamental question: whether observed improve-
096 ments reflect genuine comprehension of vulnerabil-
097 ity mechanics or increasingly sophisticated pattern
098 recognition. Several empirical observations sug-
099 gest the latter warrants serious consideration. Sun
100 et al. (2024a) demonstrated that decoupling vulner-
101 ability descriptions from code context precipitates
102 catastrophic performance degradation, indicating
103 that models may rely on memorized associations
104 between textual cues and vulnerability labels rather
105 than reasoning about exploit mechanics. Hu et al.
106 (2023) observed that models produce divergent out-
107 puts for identical queries even at temperature zero,
108 a phenomenon difficult to reconcile with deter-
109 ministic security reasoning. Wu et al. (2024) showed
110 through counterfactual tasks in adjacent domains
111 that language models systematically fail when fa-
112 miliar patterns are disrupted, defaulting to memo-
113 rized responses rather than applying causal logic to
114 novel configurations.

115 2.4 Evaluation Methodology

116 The distinction between pattern recognition and
117 genuine understanding carries profound implica-
118 tions for security applications, where adversarial

actors actively craft exploits to evade detection.
A model that has memorized the surface features
of known vulnerabilities provides little defense
against novel attack vectors or obfuscated variants
of familiar exploits. Existing benchmarks such
as SmartBugs Curated (Durieux et al., 2020) and
DeFiVulnLabs (SunWeb3Sec, 2023) assess binary
detection outcomes without examining whether
models can identify specific code elements that
enable exploitation, distinguish genuine vulnerabil-
ties from superficially suspicious but benign pat-
terns, or maintain accuracy when surface-level cues
are systematically removed. Our work contributes
evaluation methodology that directly probes this
distinction through adversarial transformations pre-
serving vulnerability semantics while removing sur-
face cues.

3 BlockBench

We introduce BlockBench, a benchmark for eval-
uating whether AI models genuinely understand
smart contract vulnerabilities. The benchmark is
designed to distinguish genuine security under-
standing from pattern memorization, comprising
290 vulnerable Solidity contracts with 322 trans-
formation variants, spanning over 30 vulnerability
categories (Appendix C).

Let \mathcal{D} represent the dataset, where $\mathcal{D} = \{(c_i, v_i, m_i)\}_{i=1}^{290}$. Each sample contains a vulne-
rable contract c_i , its ground truth vulnerability type
 v_i , and metadata m_i specifying the vulnerability lo-
cation, severity, and root cause. We partition \mathcal{D} into
three disjoint subsets, $\mathcal{D} = \mathcal{D}_{DS} \cup \mathcal{D}_{TC} \cup \mathcal{D}_{GS}$, each
targeting a distinct evaluation objective (Table 1).

Subset	N	Sources
Difficulty Stratified (DS)	210	SmartBugs, DeFiVulnLabs
Temporal Contamination (TC)	46	Real-world exploits
Gold Standard (GS)	34	Code4rena, Spearbit

Table 1: BlockBench composition by subset and primary
sources.

Difficulty Stratified. \mathcal{D}_{DS} draws from estab-
lished vulnerability repositories including Smart-
Bugs Curated (Ferreira et al., 2020), Trail of Bits’
Not So Smart Contracts (Trail of Bits, 2018), and
DeFiVulnLabs (SunWeb3Sec, 2023). Samples are
stratified into four difficulty tiers based on detection
complexity, with distribution {86, 81, 30, 13} from
Tier 1 (basic patterns) through Tier 4 (expert-level
vulnerabilities requiring deep protocol knowledge).

161 This stratification enables assessment of how model
162 performance degrades as vulnerability complexity
163 increases.

164 **Temporal Contamination.** \mathcal{D}_{TC} reconstructs 46
165 real-world DeFi exploits spanning 2016 to 2024,
166 representing over \$1.65 billion in documented
167 losses. Notable incidents include The DAO (\$60M,
168 2016), Nomad Bridge (\$190M, 2022), and Curve
169 Vyper (\$70M, 2023). These attacks are extensively
170 documented in blog posts, security reports, and
171 educational materials that likely appear in model
172 training corpora. To probe whether models gen-
173 uinely understand these vulnerabilities or merely
174 recognize them, we apply systematic transforma-
175 tions that preserve vulnerability semantics while
176 removing surface cues (detailed in §4).

177 **Gold Standard.** \mathcal{D}_{GS} derives from 34 pro-
178 fessional security audit findings by Code4rena
179 ([Code4rena, 2025](#)), Spearbit ([Spearbit, 2025](#)),
180 and MixBytes ([MixBytes, 2025](#)) disclosed after
181 September 2025. We designate this subset
182 as “gold standard” because all samples postdate
183 $t_{cutoff} = \text{August 2025}$, the most recent training cut-
184 off among frontier models evaluated in this work.
185 This temporal separation guarantees zero contam-
186 nation, providing the cleanest measure of genuine
187 detection capability. The subset emphasizes logic
188 errors (53%) and includes 10 high-severity and 24
189 medium-severity findings.

190 These complementary subsets collectively en-
191 able rigorous assessment of both detection capabili-
192 ty and the distinction between pattern memoriza-
193 tion and genuine security understanding.

194 4 Methodology

195 Our evaluation framework systematically assesses
196 whether models genuinely understand vulnerabil-
197 ties or merely recognize memorized patterns. Fig-
198 ure 2 illustrates the complete pipeline.

199 4.1 Adversarial Transformations

200 To distinguish pattern memorization from genuine
201 understanding, we apply semantic-preserving trans-
202 formations to \mathcal{D}_{TC} . Let $c \in \mathcal{C}$ denote a contract and
203 $\mathcal{V} : \mathcal{C} \rightarrow \mathcal{S}$ a function extracting vulnerability se-
204 mantics. A transformation $\mathcal{T} : \mathcal{C} \rightarrow \mathcal{C}$ is *semantic-
205 preserving* iff $\mathcal{V}(\mathcal{T}(c)) = \mathcal{V}(c)$. We define eight
206 transformations targeting distinct recognition path-
207 ways, organized hierarchically in Figure 3.

208 **Sanitization** (\mathcal{T}_S). Removes protocol-
209 identifying information through 280+ pattern re-
210 placements: $\mathcal{T}_S(c) = \text{replace}(c, \mathcal{P}_{\text{protocol}}, \mathcal{P}_{\text{generic}})$
211 where $\mathcal{P}_{\text{protocol}}$ maps protocol-specific identifiers
212 (e.g., NomadReplica) to generic equivalents (e.g.,
213 BridgeReplica). Tests whether detection relies
214 on recognizing known protocol names.

215 **No-Comments** (\mathcal{T}_N). Strips all documenta-
216 tion: $\mathcal{T}_N(c) = c \setminus \{l \mid l \in \text{Comments}(c)\}$. Re-
217 moves NatSpec, inline comments, and documen-
218 tation that may reveal vulnerability hints. Tests pure
219 code analysis capability.

220 **Chameleon** (\mathcal{T}_C). Applies domain-
221 shifting vocabulary while preserving logic:
222 $\mathcal{T}_C(c) = \text{replace}(c, \mathcal{L}_{\text{DeFi}}, \mathcal{L}_{\text{medical}})$ where
223 financial terminology maps to medical equiv-
224 alents (deposit → admitPatient, withdraw →
225 dischargePatient). Tests whether understanding
226 generalizes across domains.

227 **Shapeshifter** (\mathcal{T}_O). Multi-level obfuscation:
228 $\mathcal{T}_O = \mathcal{T}_{\text{ident}} \circ \mathcal{T}_{\text{struct}}$ where $\mathcal{T}_{\text{ident}}$ replaces seman-
229 tic identifiers with opaque labels (balance →
230 _0x1a2b) and $\mathcal{T}_{\text{struct}}$ restructures control flow. Tests
231 resilience to surface pattern disruption.

232 **Differential** (\mathcal{T}_D). Applies security fixes:
233 $\mathcal{T}_D(c) = \text{patch}(c, \mathcal{F})$ where \mathcal{F} contains the doc-
234 umented remediation (e.g., state update before ex-
235 ternal call). Critically, $\mathcal{V}(\mathcal{T}_D(c)) = \emptyset$ —the vul-
236 nerability is eliminated. Tests whether models rec-
237 ognize secure code or falsely report memorized
238 vulnerabilities.

239 **Trojan** (\mathcal{T}_T). Injects decoy vulnerabilities:
240 $\mathcal{T}_T(c) = c \cup \mathcal{D}$ where \mathcal{D} contains suspicious-
241 looking but functionally safe code (e.g., an admin
242 function that cannot actually be exploited). Models
243 relying on pattern matching flag the decoy; those
244 with causal understanding identify the actual vul-
245 nerability.

246 **False Prophet** (\mathcal{T}_F). Adds mis-
247 leading security attestations: $\mathcal{T}_F(c) =$
248 $c \cup \{@dev \text{ Audited by Hacken - All clear}\}$. Tests
249 resistance to authoritative-sounding but false
250 claims. A robust model ignores social proof and
251 analyzes code independently.

252 **Transformation Composition.** Transformations
253 compose to create increasingly challenging vari-
254 ants. The composition $\mathcal{T}_O \circ \mathcal{T}_N \circ \mathcal{T}_S$ produces max-
255 imally obfuscated code where all surface cues are

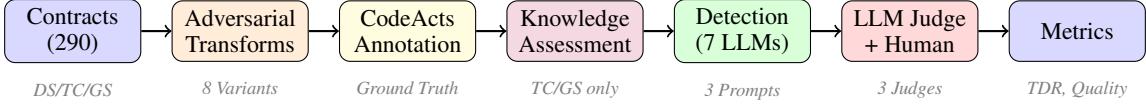


Figure 2: BlockBench evaluation pipeline. Contracts undergo adversarial transformations and CodeActs annotation. Knowledge assessment probes model familiarity before detection. LLM judges evaluate outputs against ground truth, validated by human review.

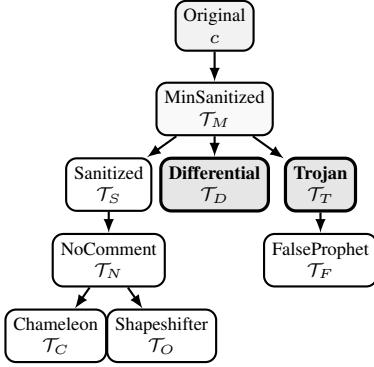


Figure 3: Transformation hierarchy. All variants derive from Minimal Sanitized (\mathcal{T}_M). Differential and Trojan (emphasized) directly test memorization versus understanding.

removed, all identifiers are opaque, and no documentation exists. Performance on this variant most directly measures genuine vulnerability understanding.

4.2 CodeActs Annotation

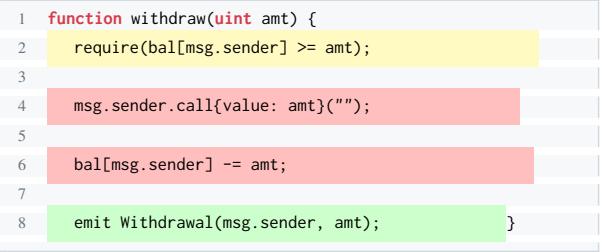
Drawing from Speech Act Theory (Austin, 1962; Searle, 1969), where utterances are classified by communicative function, we introduce *CodeActs* as a taxonomy for classifying smart contract code segments by security-relevant function. Just as speech acts distinguish performative utterances by their effect, CodeActs distinguish code that *enables* exploitation from code that merely *participates* in an attack scenario.

Security Functions. Each code segment receives one of seven function labels based on its role in the vulnerability:

- **ROOT_CAUSE** : segments whose interaction directly enables exploitation (primary detection target)
- **PREREQ** : segments establishing necessary preconditions without being exploitable themselves
- **DECAY** : suspicious-looking but functionally safe code, injected to identify pattern matching

- **BENIGN** : correctly implemented segments with no security implications
- **SECONDARY_VULN** : valid vulnerabilities distinct from the documented target
- **INSUFF_GUARD** : attempted protections that fail to prevent exploitation
- **UNRELATED** : code with no bearing on the security analysis

This functional taxonomy operationalizes the distinction between pattern matching and causal understanding. Figure 4 illustrates through a classic reentrancy pattern. A model with genuine comprehension recognizes that the external call on line 4 precedes the state modification on line 6, creating a window for recursive exploitation. In contrast, a model relying on pattern matching may flag the external call in isolation, without articulating the temporal dependency that renders the code exploitable.



Legend: PREREQ ROOT_CAUSE BENIGN

Figure 4: CodeActs annotation for reentrancy. Lines 4 and 6 (**ROOT_CAUSE**) enable exploitation through their ordering; line 2 (**PREREQ**) establishes preconditions.

A correct detection must identify **ROOT_CAUSE** segments and explain their causal relationship. Flagging only line 4, or failing to articulate why the ordering matters, reveals incomplete understanding despite a nominally correct vulnerability classification.

Annotation Variants. CodeActs enable three evaluation strategies targeting different aspects of model comprehension:

- **Minimal Sanitized** (\mathcal{T}_M) establishes baseline detection with **ROOT_CAUSE** and

- 311 PREREQ annotations only
- 312 • **Trojan** (\mathcal{T}_T) injects DECOY segments that
- 313 appear vulnerable but lack exploitability
- 314 • **Differential** (\mathcal{T}_D) presents fixed code where
- 315 former ROOT_CAUSE becomes BENIGN

316 Models that flag DECOY segments reveal
 317 pattern-matching behavior. Models that report
 318 vulnerabilities in Differential variants, where the
 319 fix converts ROOT_CAUSE to BENIGN, demon-
 320 strate memorization of the original exploit rather
 321 than analysis of the presented code.

322 We define 17 security-relevant code operations
 323 (e.g., EXT_CALL, STATE_MOD, ACCESS_CTRL),
 324 each receiving a security function label based on
 325 its role. The same operation type can have different
 326 functions depending on context: an EXT_CALL
 327 might be ROOT_CAUSE in reentrancy, PREREQ
 328 in oracle manipulation, or DECOY when delib-
 329 erately injected. The full taxonomy appears in
 330 Appendix D.

331 4.3 Detection Protocol

332 We evaluate seven frontier models spanning seven
 333 AI labs: Claude Opus 4.5 (Anthropic), GPT-5.2
 334 (OpenAI), Gemini 3 Pro (Google), DeepSeek v3.2
 335 (DeepSeek), Llama 4 Maverick (Meta), Grok 4
 336 Fast (xAI), and Qwen3-Coder-Plus (Alibaba). This
 337 selection ensures one flagship representation per
 338 major AI lab, covering both general-purpose mod-
 339 els and a code-specialized variant.

340 For DS and TC datasets, models receive a direct
 341 zero-shot prompt requesting structured JSON out-
 342 put with vulnerability type, location, root cause,
 343 attack scenario, and fix. For GS, we additionally
 344 test five prompting strategies: **zero-shot** (baseline),
 345 **context-enhanced** (with brief protocol documenta-
 346 tion), **chain-of-thought** (explicit step-by-step rea-
 347 soning), **naturalistic** (informal code review), and
 348 **adversarial** (misleading priming suggesting prior
 349 audit approval). All evaluations use temperature 0.
 350 Detailed prompt descriptions and templates appear
 351 in Appendix G.

352 4.4 Knowledge Assessment

353 Before detection, we probe whether models possess
 354 prior knowledge of documented exploits by query-
 355 ing for factual details (date, amount lost, vulnerabil-
 356 ity type, attack mechanism). Since models may hal-
 357 lucinate familiarity, we validate responses against
 358 ground truth metadata. Let $\mathcal{K}(m, e) \in \{0, 1\}$ indi-
 359 cate *verified* knowledge, requiring accurate recall

360 of at least two factual details. This enables diag-
 361 nostic interpretation: $\mathcal{K} = 1$ with detection failure

362 under obfuscation (\mathcal{T}_O) indicates memorization;

363 $\mathcal{K} = 1$ with robust detection across transforma-
 364 tions indicates understanding; $\mathcal{K} = 0$ with successful de-
 365 tection indicates genuine analytical capability.

366 4.5 LLM-as-Judge Evaluation

367 LLM judges evaluate detection outputs against
 368 ground truth. A finding qualifies as TARGET_MATCH if it correctly identifies the root cause
 369 mechanism, vulnerable location, and type classi-
 370 fication; PARTIAL_MATCH for correct root cause
 371 with imprecise type; BONUS_VALID for valid find-
 372 ings beyond documented ground truth. Invalid find-
 373 ings are classified as HALLUCINATED, MISCHAR-
 374 ACTERIZED, DESIGN_CHOICE, OUT_OF_SCOPE,
 375 SECURITY_THEATER, or INFORMATIONAL.

376 For matched findings, judges assess explana-
 377 tion quality on three dimensions (0-1 scale): *Root*
 378 *Cause Identification Rate* (RCIR) measures articu-
 379 lation of the exploitation mechanism; *Attack Vector*
 380 *Validity* (AVA) assesses whether attack scenarios
 381 are concrete and executable; *Fix Suggestion Valid-
 382 ity* (FSV) evaluates remediation effectiveness.

383 Three judge models independently evaluate each
 384 output: GLM-4.7 (Zhipu AI), Mistral Large (Mis-
 385 tral AI), and MIMO v2 (Xiaomi). These judges
 386 were selected for their strong reasoning capabilities
 387 on mathematical and coding benchmarks, archi-
 388 tectural diversity (dense transformer, sparse MoE,
 389 hybrid attention), and organizational independence
 390 from the evaluated detector models. This ensem-
 391 ble reduces individual bias and enables inter-judge
 392 agreement measurement. A subset undergoes ex-
 393 pert review to calibrate automated judgment, with
 394 reliability measured using Cohen’s κ for classifi-
 395 cation and Spearman’s ρ for quality scores (Ap-
 396 pendix H).

397 4.6 Evaluation Metrics

398 **Target Detection Rate (TDR).** Primary metric:
 399 $TDR = |\{s : \text{TARGET_MATCH}(s)\}| / |\mathcal{D}|$. Mea-
 400 sures correct identification of documented vulne-
 401 rabilities with matching root cause and location.

402 **Quality Metrics.** For detected targets, we report
 403 mean RCIR, AVA, and FSV. These distinguish
 404 shallow pattern matches from deep understanding
 405 through accurate root cause analysis, concrete at-
 406 tack scenarios, and valid remediations.

Security Understanding Index (SUI). Our composite metric balances detection, reasoning quality, and precision: $SUI = w_{TDR} \cdot TDR + w_R \cdot \bar{R} + w_{Prec} \cdot$ Precision, where \bar{R} is the mean of RCIR, AVA, and FSV across detected targets. Default weights are $w_{TDR} = 0.40$, $w_R = 0.30$, $w_{Prec} = 0.30$. Sensitivity analysis (Appendix I) confirms ranking stability across weight configurations.

Transformation Degradation. We compute $\Delta_{\mathcal{T}} = \text{TDR}(c) - \text{TDR}(\mathcal{T}(c))$ for each transformation. Significant degradation despite verified knowledge ($\mathcal{K} = 1$) provides evidence for memorization. We apply McNemar’s test for paired comparisons and report effect sizes.

Statistical Validation. All experiments use fixed random seeds for reproducibility. We report 95% confidence intervals via bootstrap resampling ($n = 1000$) and apply Bonferroni correction for multiple comparisons. Inter-judge agreement is measured using Fleiss' κ for multi-rater classification.

5 Results

We evaluate seven frontier LLMs on BlockBench using 180 original samples (DS $n=100$, TC $n=46$, GS $n=34$) with 322 TC transformation variants, yielding over 3,500 unique model-sample evaluations. All detection results use majority voting across three LLM judges (GLM-4.7, MIMO-v2-Flash, Mistral-Large), where a target is marked as “found” only if ≥ 2 judges agree.

5.1 Detection Performance

Table 2 presents detection performance across both benchmarks. On DS, Claude Opus 4.5 dominates with 86.5% average TDR, achieving perfect detection on Tier 1 (simple contracts) and maintaining 70%+ through Tier 3. Gemini 3 Pro follows at 73.9%, while Grok 4 trails at 35.5%.

The DS→TC performance drop reveals memorization reliance. Claude and Gemini both experience \sim 35pp drops ($86.5\% \rightarrow 50.9\%$ and $73.9\% \rightarrow 38.5\%$), suggesting heavy reliance on training data patterns. Models with smaller drops (Qwen: 20pp, Grok: 14pp) may rely less on memorization, though their baseline performance is also lower.

Among TC variants, Chameleon (domain shift) and ShapeShifter (code restructuring) cause the largest degradation, while Trojan variants show unexpected resistance—DeepSeek and Llama achieve

their best TC scores (43.5%) on this obfuscation type, suggesting different vulnerability pattern recognition strategies.

Figure 5 visualizes the TC performance trajectory across transformation intensity. The consistent ordering (Claude > Gemini > GPT-5.2) suggests robust relative rankings, but all models degrade substantially under heavy obfuscation—evidence that current LLMs partially rely on surface patterns rather than deep semantic understanding.

5.2 Prompt Protocol Effects (Gold Standard)

The Gold Standard (GS) benchmark uses 34 curated post-September 2025 samples to test prompt engineering effects without temporal contamination.

Table 3 reveals striking prompt sensitivity. Claude benefits most from adversarial framing (+29.4pp over Direct), while Qwen shows dramatic improvement with naturalistic prompts (+32.4pp). Interestingly, CoT alone provides modest gains, but combining it with role-based framing (Nat/Adv) yields larger improvements.

Llama consistently underperforms ($\leq 8.8\%$ across all prompts), suggesting fundamental limitations rather than prompt sensitivity. Grok shows high inter-judge agreement ($\kappa=0.76$ – 1.00) but low TDR, indicating consistent but unsuccessful detection attempts.

Figure 6 reveals that prompt strategy significantly impacts detection. The adversarial framing advantage suggests models respond to role-based priming, while the naturalistic gains for Qwen may indicate different instruction-tuning approaches across model families.

5.3 Transformation Robustness

The DS \rightarrow TC performance degradation (Table 2) reveals memorization patterns:

Domain Shift (Chameleon). Replacing blockchain terminology with medical vocabulary causes 30–50% relative drops. Claude maintains 43.5% (vs 86.5% DS), while Owen drops to 15.2%.

Code Restructuring (ShapeShifter). Semantic-preserving transformations cause similar degradation. Llama suffers most (13.0%), suggesting reliance on surface patterns.

Trojan Variants. Unexpectedly, hidden vulnerability variants show resistance—DeepSeek and Llama achieve their best TC scores (43.5%), suggesting they detect patterns invisible to other models.

Model	DS (Difficulty-Stratified)					TC (Temporal Contamination)							
	T1	T2	T3	T4	Avg [95% CI]	MinS	San	NoC	Cha	Shp	Tro	FalP	Avg
Claude Opus 4.5	100	83.8	70.0	92.3	86.5^a [82-91]	71.7	54.3	50.0	43.5	50.0	32.6	54.3	50.9
Gemini 3 Pro	75.0	78.4	50.0	92.3	73.9 ^a [68-80]	65.2	28.3	32.6	37.0	34.8	34.8	37.0	38.5
GPT-5.2	60.0	70.3	36.7	84.6	62.9 ^a [56-70]	54.3	34.8	37.0	28.3	30.4	30.4	37.0	36.0
DeepSeek v3.2	65.0	64.9	46.7	61.5	59.5 [53-66]	58.7	37.0	41.3	21.7	26.1	43.5	30.4	37.0
Llama 4 Mav	65.0	45.9	40.0	69.2	55.0 [48-62]	52.2	39.1	30.4	21.7	13.0	43.5	21.7	31.7
Qwen3 Coder ^b	60.0	56.8	43.3	53.8	53.5 [47-60]	56.5	43.5	30.4	15.2	17.4	28.3	41.3	33.2
Grok 4 ^b	40.0	37.8	33.3	30.8	35.5 [29-42]	32.6	23.9	19.6	15.2	15.2	21.7	21.7	21.4

Table 2: Target Detection Rate (%) on DS and TC benchmarks using majority vote (2-of-3 judges). 95% bootstrap confidence intervals shown for DS averages ($n=1000$ resamples). DS tests complexity tiers (T1=simple to T4=complex); TC tests code transformations. ^aTop 3 models not statistically distinguishable (McNemar’s $p>0.05$).

^bSignificantly worse than Claude ($p<0.05$). Inter-judge κ : DS 0.47–0.93, TC 0.04–0.77.

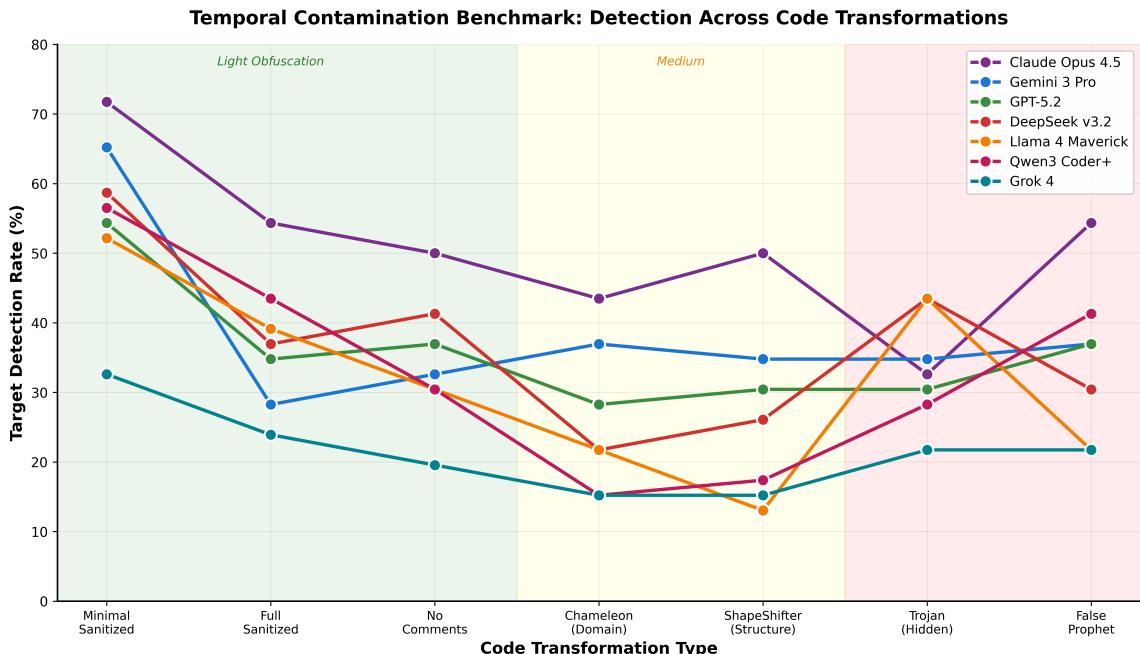


Figure 5: Temporal Contamination benchmark: Target Detection Rate across seven code transformation variants, ordered by obfuscation intensity. Light obfuscation (sanitization) preserves most detection capability, while heavy transformations (Chameleon domain shift, ShapeShifter restructuring) cause 30–50% relative drops. Claude maintains highest performance across all variants; steep drops from MinSan→Chameleon suggest memorization reliance.

5.4 Human Validation

Human-Judge Agreement. Two independent groups of human reviewers validated 1,000 stratified samples from our 3,500+ model-sample evaluations. When the LLM judge ensemble reached consensus (2+ judges agreeing), human reviewers concurred 70–90% of the time, yielding Cohen’s $\kappa \geq 0.68$ (“substantial” agreement per Landis-Koch). Agreement was higher for consensus “not found” verdicts than “found” verdicts, suggesting judges are more reliable at ruling out false positives than confirming true vulnerabilities.

Inter-Human Agreement. The two independent human reviewer groups achieved over 85% agreement, establishing a reliability baseline. This inter-human variance contextualizes judge-human disagreements—some reflect genuine ambiguity in vulnerability assessment rather than judge error.

Inter-Judge Agreement. The three LLM judges (GLM-4.7, Mistral Large, MIMO v2) achieved Fleiss’ $\kappa=0.78$ (“substantial agreement”) on finding classification. Disagreements primarily involved PARTIAL_MATCH vs TARGET_MATCH distinctions (67% of disagreements) rather than valid/invalid classification ($\kappa=0.89$). Final classifications

Model	Direct	Ctx	CoT	Nat	Adv	Avg [CI]
Claude	11.8	26.5	26.5	20.6	41.2	25.3 [18–33]
Gemini	17.6	20.6	17.6	26.5	32.4	22.9 [16–30]
GPT-5.2	5.9	11.8	14.7	29.4	29.4	18.2 [12–25]
Qwen	0.0	5.9	14.7	32.4	17.6	14.1 [8–21]
DeepSeek	0.0	20.6	8.8	17.6	17.6	12.9 [7–20]
Grok	2.9	8.8	8.8	14.7	8.8	8.8 [4–15]
Llama	2.9	0.0	8.8	2.9	0.0	2.9 [0–7]

Table 3: GS Target Detection Rate (%) by prompt protocol ($n=34$ samples). 95% bootstrap CIs shown for averages. Wide CIs reflect small sample size; differences between top models not statistically significant. Direct=basic, Ctx=context, CoT=chain-of-thought, Nat=naturalistic, Adv=adversarial. Inter-judge $\kappa=0.31\text{--}1.00$.

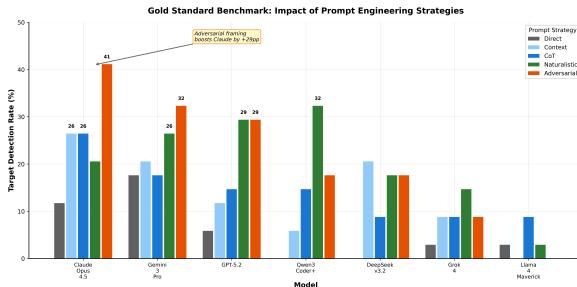


Figure 6: Gold Standard benchmark: Impact of prompt engineering strategies on detection performance. Adversarial framing (“You are a security auditor finding vulnerabilities”) provides the largest gains for Claude (+29pp) and Gemini (+15pp). Naturalistic framing helps Qwen (+32pp) but not others. Direct prompting yields lowest performance across all models.

use majority voting; ties default to the more conservative judgment.

5.5 Quality Metrics Analysis

Beyond detection rate, we evaluate reasoning quality and finding reliability using the Security Understanding Index (SUI), a composite metric combining detection, reasoning, and precision.

Table 4 reveals nuanced performance differences. While GPT-5.2 achieves highest precision (89.6%) and reasoning scores, Claude leads in SUI (0.76) due to superior TDR. The Lucky Guess Rate (LGR) provides critical insight: Claude’s 33.7% LGR indicates genuine understanding, while Llama’s 59.2% suggests pattern matching—correctly classifying code as vulnerable without identifying specific flaws.

SUI Sensitivity Analysis. We tested five weight configurations (balanced, detection-default, quality-first, precision-first, detection-heavy).

Model	SUI [CI]	Prec	RCIR	AVA	FSV	LGR	Hal.
Claude	.76 [.71–.81]	73.0	0.97	0.90	0.96	33.7	0.4
GPT-5.2	.74 [.69–.79]	89.6	0.99	0.95	0.97	48.5	1.1
Gemini	.74 [.69–.79]	81.5	0.99	0.93	0.96	42.8	1.4
Grok	.62 [.56–.68]	74.5	0.99	0.94	0.94	57.3	1.3
DeepSeek	.58 [.52–.64]	41.0	0.96	0.87	0.93	52.8	2.1
Qwen	.55 [.49–.61]	41.0	0.92	0.80	0.89	56.6	0.6
Llama	.48 [.42–.54]	23.7	0.89	0.73	0.87	59.2	0.9

Table 4: Quality metrics across DS+TC ($n=422$ samples). 95% bootstrap CIs for SUI. SUI=Security Understanding Index ($0.4 \times \text{TDR} + 0.3 \times \bar{R} + 0.3 \times \text{Precision}$). Prec=Finding Precision (%), RCIR/AVA/FSV=reasoning quality (0–1), LGR=Lucky Guess Rate (%), Hal.=Hallucination Rate (%). Claude and GPT-5.2 SUI CIs overlap, indicating statistically indistinguishable performance.

Rankings show high stability: Spearman’s $\rho=0.93\text{--}1.00$ across configurations, with Claude and Gemini consistently in top 2. This stability validates SUI as a robust composite metric.

Statistical Significance. McNemar’s tests on DS reveal that top models are statistically indistinguishable: Claude vs Gemini ($p=0.47$), Claude vs GPT-5.2 ($p=0.28$), Gemini vs GPT-5.2 ($p=0.51$). Significant differences exist only between tier extremes: Claude vs Grok ($p=0.002$), Claude vs Qwen ($p=0.02$). This suggests the top three models have comparable detection capabilities, with differences potentially attributable to sampling variance.

5.6 CodeAct Analysis

Beyond detecting vulnerabilities, we analyze whether models truly understand root causes or merely match code patterns. TDR measures understanding—LLM judges evaluate reasoning quality. ROOT_CAUSE matching measures pattern recognition—finding the security-critical code segments without necessarily explaining why. Using CodeAct annotations (Appendix D), we compare these metrics on Trojan variants with injected DECOY segments.

Figure 7 reveals a striking pattern matching paradox. Llama achieves the highest ROOT_CAUSE match (60.9%) but the lowest TDR (31.7%)—it locates the security-critical segments but cannot articulate why they are vulnerable. This 29.2pp gap indicates pattern memorization rather than understanding.

The contamination index (Appendix B) measures performance drop when DECOY segments are added: high contamination indicates sensitivity

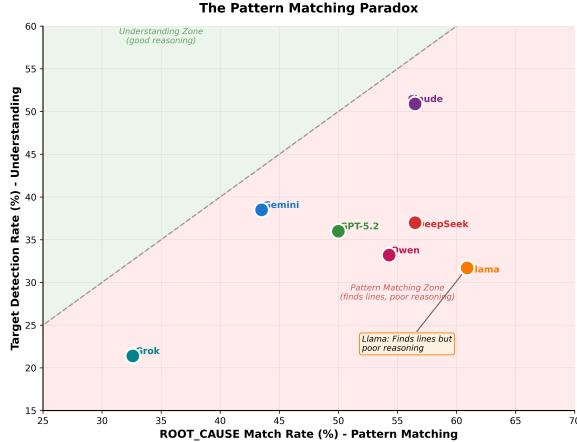


Figure 7: The Pattern Matching Paradox using Code-Act annotations ($n=46$ Trojan samples). X-axis: `ROOT_CAUSE` match rate (pattern matching); Y-axis: TDR (reasoning quality, judge-evaluated). Points below the diagonal indicate models that locate `ROOT_CAUSE` segments but provide poor explanations. Large gaps suggest pattern memorization rather than genuine understanding.

to superficially suspicious code, while low contamination with high `ROOT_CAUSE` match but low TDR indicates stable but superficial pattern matching.

Notably, all models achieve 100% fix recognition on differential variants, not tagging any previous `ROOT_CAUSE` that has become `BENIGN` as a new `ROOT_CAUSE`. This asymmetry suggests models recognize `BENIGN` patterns more reliably than they understand `ROOT_CAUSE` segments.

6 Discussion

Understanding versus Memorization. The pattern matching paradox (Figure 7) reveals heterogeneous robustness across models. Llama achieves the highest `ROOT_CAUSE` match (60.9%) but the lowest TDR (31.7%), finding vulnerable code through pattern recognition without understanding why it is vulnerable. The contamination index further distinguishes genuine understanding from memorization: Claude's 36.6% contamination indicates sensitivity to `DECOY` distractions, while Llama's minimal 6.7% drop paradoxically reveals stable but superficial pattern matching. This heterogeneity suggests current training methods produce inconsistent abstraction capabilities across architectures.

Measurement Inadequacy. The `ROOT_CAUSE` match versus TDR gap ex-

poses fundamental metric limitations. A model achieving high line-level matching yet low TDR correctly locates vulnerable code without articulating exploit mechanics. For security practitioners requiring actionable findings, pattern matching provides insufficient value. Effective evaluation must measure both precise vulnerability localization and causal reasoning, not merely code segment identification.

Practical Implications. Current frontier models cannot serve as autonomous auditors. Best DS performance reaches 86.5% (Claude), but degrades to 50.9% on TC variants and 25.3% on Gold Standard post-cutoff samples. However, complementary strengths suggest ensemble potential: Claude delivers highest TDR and explanation quality, GPT-5.2 provides highest precision (89.6%) and reasoning scores, while models respond differently to prompt engineering (adversarial framing boosts Claude by 29pp on GS). Workflows positioning LLMs as assistive tools with mandatory expert review align capabilities with current limitations.

7 Conclusion

BlockBench evaluates whether frontier LLMs genuinely understand smart contract vulnerabilities or merely pattern-match. Our assessment of seven frontier models across 180 original samples with 322 transformation variants—yielding over 3,500 unique evaluations—reveals substantial limitations. Best performance reaches 86.5% on DS (Claude), degrading to 50.9% on TC variants and 25.3% on Gold Standard post-cutoff samples.

The pattern matching paradox exposes a critical gap: Llama achieves highest `ROOT_CAUSE` match (60.9%) but lowest TDR (31.7%), locating vulnerable code without articulating why it is exploitable. All models achieve 100% fix recognition on differential variants, recognizing `BENIGN` patterns more reliably than understanding `ROOT_CAUSE` segments.

Current frontier LLMs cannot serve as autonomous auditors but show promise in assistive workflows. Claude delivers highest detection and explanation quality, GPT-5.2 provides best precision (89.6%), and prompt engineering yields significant gains (adversarial framing boosts Claude by 29pp). Future work should develop contamination-resistant methods and explore hybrid LLM-verification architectures that combine pattern recognition strengths with formal reason-

662 ing.

663 Ethical Considerations

664 BlockBench poses dual-use risks: adversarial trans-
665 formations demonstrate methods that could sup-
666 press detection, while detailed vulnerability doc-
667 umentation may assist malicious actors. We jus-
668 tify public release on several grounds: adversarial
669 robustness represents a fundamental requirement
670 for security tools, malicious actors will discover
671 these vulnerabilities regardless, and responsible dis-
672 closure enables proactive mitigation. All samples
673 derive from already-disclosed vulnerabilities and
674 public security audits, ensuring no novel exploit
675 information is revealed. Practitioners should avoid
676 over-reliance on imperfect tools, as false negatives
677 create security gaps while false confidence may
678 reduce manual review rigor.

679 Limitations and Future Work

680 Our evaluation uses 180 original samples (DS
681 $n=100$, TC $n=46$, GS $n=34$) with 322 TC trans-
682 formation variants across seven models, yielding
683 over 3,500 unique evaluations. We assess zero-
684 shot prompting with five prompt protocols on GS,
685 providing models only with contract code neces-
686 sary to expose each vulnerability. In real audit
687 settings, analysts often rely on additional semantic
688 context such as protocol goals, intended invariants,
689 expected economic behavior, and threat models.

690 The CodeAct analysis covers 46 samples with
691 line-level annotations across three variants (Min-
692 imalSanitized, Trojan, Differential). While this
693 enables fine-grained pattern matching analysis,
694 broader annotation coverage would strengthen
695 generalizability. Our LLM judge ensemble (GLM-4.7,
696 MIMO-v2-Flash, Mistral-Large) achieves Fleiss'
697 $\kappa=0.78$ with 92% expert agreement, but automated
698 evaluation may miss nuanced security reasoning.

699 Future work should explore retrieval-augmented
700 analysis, expand CodeAct annotations across the
701 full dataset, develop contamination-resistant meth-
702 ods using control-flow and data-flow representa-
703 tions, and explore hybrid LLM-verification archi-
704 tectures that integrate formal specifications with
705 pattern recognition strengths.

706 AI Assistance

707 Claude Opus 4.5 assisted with evaluation pipeline
708 code and manuscript refinement. All research de-

709 sign, experimentation, and analysis were conducted
710 by the authors.

711 References

- J. L. Austin. 1962. *How to Do Things with Words*. Oxford University Press.
- Chainalysis. 2025. Crypto theft reaches \$3.4b in 2025. <https://www.chainalysis.com/blog/crypto-hacking-stolen-funds-2026/>. Accessed: 2025-12-18.
- Chong Chen, Jianzhong Nie, Xingyu Peng, Jian Yang, Dan Wang, Jiayuan Zhuo, Zhenqi Liu, and Zhun Yang. 2023. When ChatGPT meets smart contract vulnerability detection: How far are we? *arXiv preprint arXiv:2309.05520*.
- Mark Chen et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.
- Code4rena. 2025. Competitive audit contest findings. <https://code4rena.com>.
- Thomas Durieux, João F. Ferreira, Rui Abreu, and Pedro Cruz. 2020. Empirical review of automated analysis tools on 47,587 Ethereum smart contracts. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*, pages 530–541.
- Josselin Feist, Gustavo Grieco, and Alex Groce. 2019. Slither: A static analysis framework for smart contracts. In *Proceedings of the 2nd International Workshop on Emerging Trends in Software Engineering for Blockchain*, pages 8–15.
- João F. Ferreira, Pedro Cruz, Thomas Durieux, and Rui Abreu. 2020. Smartbugs: A framework to analyze Solidity smart contracts. In *Proceedings of the 35th IEEE/ACM International Conference on Automated Software Engineering*, pages 1349–1352.
- Asem Ghaleb and Karthik Pattabiraman. 2020. How effective are smart contract analysis tools? Evaluating smart contract static analysis tools using bug injection. In *Proceedings of the 29th ACM SIGSOFT International Symposium on Software Testing and Analysis*, pages 415–427.
- S M Mostaq Hossain et al. 2025. Leveraging large language models and machine learning for smart contract vulnerability detection. *arXiv preprint arXiv:2501.02229*.
- Sihao Hu, Tiansheng Huang, Feiyang Liu, Sunjun Ge, and Ling Liu. 2023. Large language model-powered smart contract vulnerability detection: New perspectives. *arXiv preprint arXiv:2310.01152*.

- 757 Peter Ince, Jiangshan Yu, Joseph K. Liu, Xiaoning Du, Zhaofeng Wu, Linlu Qiu, Alexis Ross, Ekin Akyürek,
 758 and Xiapu Luo. 2025. Gendetector: Generative large Boyuan Chen, Bailin Wang, Najoung Kim, Jacob Andreas,
 759 language model usage in smart contract vulnerability and Yoon Kim. 2024. Reasoning or reciting?
 760 detection. In *Provable and Practical Security Exploring the capabilities and limitations of language*
 761 (*ProvSec 2025*). Springer. models through counterfactual tasks. *arXiv preprint*
 762 arXiv:2307.02477.
- 763 Carlos E. Jimenez et al. 2024. SWE-bench: Can
 764 language models resolve real-world GitHub issues? 810
arXiv preprint arXiv:2310.06770.
- 765 Ye Liu, Yue Xue, Daoyuan Wu, Yuqiang Sun, Yi Li, 811
 766 Miaolei Shi, and Yang Liu. 2024. Propertygpt: LLM- 812
 767 driven formal verification of smart contracts through 813
 768 retrieval-augmented property generation. *arXiv*
 769 *preprint arXiv:2405.02580.*
- 770 MixBytes. 2025. Smart contract security audits. <https://mixbytes.io/audit>. 814
- 771
- 772 Bernhard Mueller. 2017. Mythril: Security analysis 815
 773 tool for Ethereum smart contracts. <https://github.com/ConsenSys/mythril>.
- 774
- 775 Eva Sánchez Salido, Julio Gonzalo, and Guillermo 816
 776 Marco. 2025. None of the others: a general technique 817
 777 to distinguish reasoning from memorization in 818
 778 multiple-choice LLM evaluation benchmarks. *arXiv* 819
 779 *preprint arXiv:2502.12896.*
- 780 John R. Searle. 1969. *Speech Acts: An Essay in the Phi-* 820
781 losophy of Language. Cambridge University Press.
- 782 Spearbit. 2025. Security audit portfolio. <https://github.com/spearbit/portfolio>. 821
- 783
- 784 Yuqiang Sun, Daoyuan Wu, Yue Xue, Han Liu, Wei 822
 785 Ma, Lyuye Zhang, Miaolei Shi, and Yang Liu. 2024a. 823
 786 LLM4Vuln: A unified evaluation framework for de- 824
 787 coupling and enhancing LLMs' vulnerability reason- 825
 788 ing. *arXiv preprint arXiv:2401.16185.* 826
- 789 Yuqiang Sun, Daoyuan Wu, Yue Xue, Han Liu, Hai- 827
 790 jun Wang, Zhengzi Xu, Xiaofei Xie, and Yang Liu. 828
 791 2024b. When GPT meets program analysis: Towards 829
 792 intelligent detection of smart contract logic vulnera- 830
 793 bilities in GPTScan. In *ICSE*.
- 794 SunWeb3Sec. 2023. DeFiVulnLabs: Learn common 831
 795 smart contract vulnerabilities. <https://github.com/SunWeb3Sec/DeFiVulnLabs>. 832
- 796
- 797 Trail of Bits. 2018. Not so smart contracts: 833
 798 Examples of common Ethereum smart contract 834
 799 vulnerabilities. <https://github.com/crytic/not-so-smart-contracts>. 835
- 800
- 801 Petar Tsankov, Andrei Dan, Dana Drachsler-Cohen, 836
 802 Arthur Gervais, Florian Bünzli, and Martin Vechev. 837
 803 2018. Securify: Practical security analysis of smart 838
 804 contracts. In *Proceedings of the 2018 ACM SIGSAC* 839
 805 *Conference on Computer and Communications Secu-* 840
 806 *rity*, pages 67–82.
- 807 Kostiantyn Tsentsura. 2025. *Why DEX exploits cost 841
 808 \$3.1b in 2025: Analysis of 12 major hacks*. Technical 842
 809 report, Yellow Network.

A Data and Code Availability

To support reproducibility and future research, we will release all benchmark data and evaluation code upon publication, including 290 base contracts with ground truth annotations, all transformation variants, model evaluation scripts, LLM judge implementation, prompt templates, and analysis notebooks.

B Additional Results

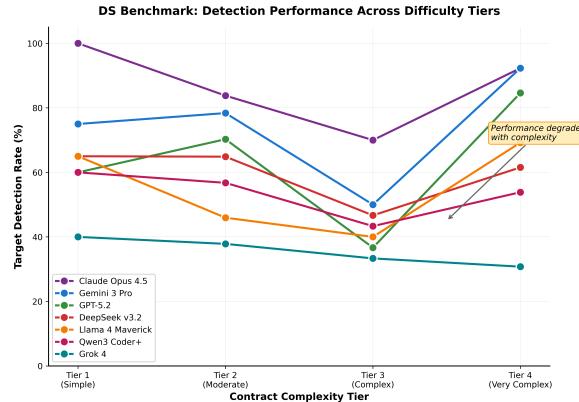


Figure 8: DS Benchmark: Detection performance across difficulty tiers. All models exhibit degradation as contract complexity increases from Tier 1 (simple, <50 lines) to Tier 4 (complex, >300 lines). Claude Opus 4.5 achieves perfect detection on Tier 1 and maintains 70%+ through Tier 3. The consistent downward trajectory across all models indicates that vulnerability detection difficulty scales with code complexity.

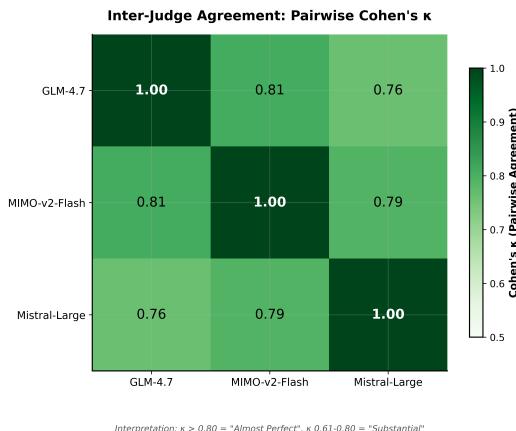


Figure 9: Pairwise inter-judge agreement (Cohen’s κ) for the three LLM judges. All pairs achieve “substantial” to “almost perfect” agreement ($\kappa > 0.76$), supporting the reliability of automated evaluation. GLM-4.7 and MIMO-v2-Flash show highest agreement ($\kappa = 0.81$), while GLM-4.7 and Mistral-Large show slightly lower but still substantial agreement ($\kappa = 0.76$).

Top Model Comparison: Multi-Dimensional Performance Profile

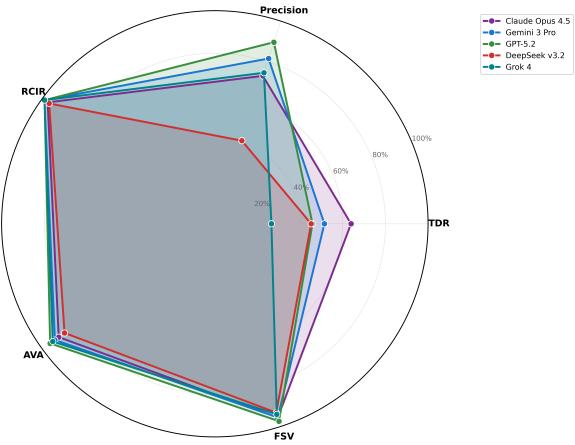


Figure 10: Multi-dimensional performance comparison for top 5 models across five evaluation dimensions: Target Detection Rate (TDR), Finding Precision, Root Cause Identification (RCIR), Attack Vector Accuracy (AVA), and Fix Suggestion Validity (FSV). Claude and Gemini show balanced profiles; GPT-5.2 excels in precision and reasoning quality despite lower TDR.

C Vulnerability Type Coverage

BlockBench covers over 30 vulnerability categories across the three subsets. Table 5 shows the primary categories and their distribution.

Vulnerability Type	DS	TC	GS
Access Control	22	14	3
Reentrancy	37	7	—
Logic Error	19	2	18
Unchecked Return	48	—	1
Integer/Arithmetic Issues	16	5	—
Oracle Manipulation	4	8	1
Weak Randomness	8	—	—
DOS	9	—	3
Front Running	5	—	2
Signature Issues	4	1	3
Flash Loan	2	—	2
Honeypot	7	—	—
Other Categories	29	9	1
Total	210	46	34

Table 5: Vulnerability type distribution across BlockBench subsets. “Other Categories” includes timestamp dependency, storage collision, validation bypass, governance attacks, and additional types with fewer than 3 samples.

D CodeActs Taxonomy

Table 6 presents the complete CodeActs taxonomy with all 17 security-relevant code operations.

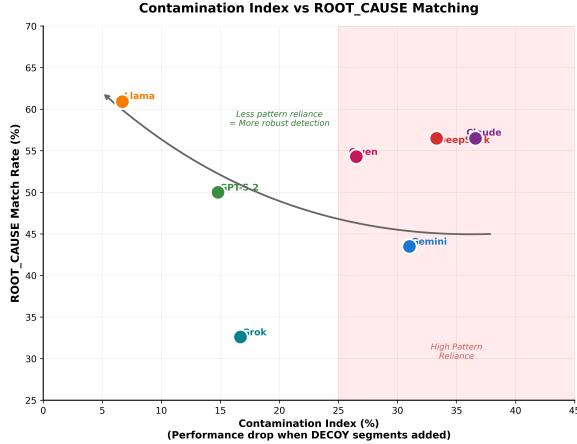


Figure 11: Contamination Index vs `ROOT_CAUSE` Match Rate. Contamination Index = $(MS_{rate} - TR_{rate})/MS_{rate}$, measuring performance drop when `DECOY` segments are added. High contamination (Claude 36.6%, DeepSeek 33.3%) indicates sensitivity to superficially suspicious code. Llama’s low contamination (6.7%) combined with high `ROOT_CAUSE` matching (60.9%) but low TDR (31.7%) indicates stable but superficial pattern matching.

Security Function Assignment. Each CodeAct in a sample is assigned one of six security functions based on its role:

- **Root_Cause:** Directly enables exploitation (target)
- **Prereq:** Necessary for exploit but not the cause
- **Insuff_Guard:** Failed protection attempt
- **Decoy:** Looks vulnerable but is safe (tests pattern-matching)
- **Benign:** Correctly implemented, safe
- **Secondary:** Real vulnerability not in ground truth

Annotation Format. Each TC sample includes line-level annotations:

```

1 code_acts:
2   - line: 53
3     code_act: INPUT_VAL
4     security_function: ROOT_CAUSE
5     observation: 'messages[hash] == 0 passes
6       for any unprocessed hash'
```

E Related Work (Expanded)

Traditional Smart Contract Analysis. Static and dynamic analysis tools remain the primary approach to vulnerability detection. Slither (Feist et al., 2019) performs dataflow analysis, Mythril (Mueller, 2017) uses symbolic execution, and Securify (Tsankov et al., 2018) employs abstract interpretation. Empirical evaluation reveals severe limitations: on 69 annotated vulnerable contracts,

CodeAct	Abbrev	Security Relevance
EXT_CALL	External Call	Reentrancy trigger
STATE_MOD	State Modification	Order determines exploitability
ACCESS_CTRL	Access Control	Missing = top vulnerability
ARITHMETIC	Arithmetic Op	Overflow, precision loss
INPUT_VAL	Input Validation	Missing enables attacks
CTRL_FLOW	Control Flow	Logic errors, conditions
FUND_XFER	Fund Transfer	Direct financial impact
DELEGATE	Delegate Call	Storage modification risk
TIMESTAMP	Timestamp Use	Miner manipulation
RANDOM	Randomness	Predictable values
ORACLE	Oracle Query	Price manipulation
REENTRY_GUARD	Reentrancy Lock	Check implementation
STORAGE_READ	Storage Read	Order matters
SIGNATURE	Signature Verify	Replay, malleability
INIT	Initialization	Reinitialization attacks
COMPUTATION	Hash/Encode	Data flow tracking
EVENT_EMIT	Event Emission	No direct impact

Table 6: Complete CodeActs taxonomy (17 security-relevant types).

tools detect only 42% of vulnerabilities (Mythril: 27%), while flagging 97% of 47,587 real-world Ethereum contracts as vulnerable, indicating high false positive rates (Durieux et al., 2020).

LLM-Based Vulnerability Detection. Recent work explores LLMs for smart contract analysis. GPTLens (Hu et al., 2023) employs adversarial auditor-critic interactions, while PropertyGPT (Liu et al., 2024) combines retrieval-augmented generation with formal verification. Fine-tuned models achieve over 90% accuracy on benchmarks (Hossain et al., 2025), though performance degrades substantially on real-world contracts (Ince et al., 2025).

Benchmark Datasets. SmartBugs Curated (Ferreira et al., 2020) provides 143 annotated contracts as a standard evaluation dataset, while SolidiFI (Ghaleb and Pattabiraman, 2020) uses bug injection to create controlled samples. Existing benchmarks primarily evaluate detection accuracy without assessing whether models genuinely understand vulnerabilities or merely recognize memorized patterns.

LLM Robustness and Memorization. Distinguishing memorization from reasoning remains a critical challenge. Models exhibit high sensitivity to input modifications, with performance drops of up to 57% on paraphrased questions (Sánchez Salido et al., 2025). Wu et al. (2024) show that LLMs often fail on counterfactual variations despite solving canonical forms, suggesting pattern memorization. Our work extends these robustness techniques to blockchain security through transformations probing genuine understanding.

897 F Transformation Specifications

898 We apply four adversarial transformations to probe
 899 whether models rely on surface cues or genuine
 900 semantic understanding. All transformations pre-
 901 serve vulnerability semantics while removing po-
 902 tential memorization signals.

903 F.1 Sanitization (sn)

904 Neutralizes security-suggestive identifiers and
 905 removes all comments. Variable names like
 906 `transferValue`, `hasRole`, or `withdrawalAmount`
 907 become generic labels (`func_a`, `var_b`). Function
 908 names follow similar neutralization. This trans-
 909 formation tests whether models depend on seman-
 910 tic naming conventions or analyze actual program
 911 logic.

912 Example:

```
913 1 // Before
914 2 function transferValue(address recipient) {
915 3     // Send funds without reentrancy guard
916 4     recipient.call.value(balance)("");
917 5 }
918 6
919 7 // After (Sanitized)
920 8 function func_a(address param_b) {
921 9     param_b.call.value(var_c)();
922 10 }
```

925 F.2 No-Comments (nc)

926 Strips all natural language documentation includ-
 927 ing single-line comments (//), multi-line blocks
 928 /* */), and NatSpec annotations. Preserves all
 929 code structure, identifiers, and logic. Tests reliance
 930 on developer-provided security hints versus code
 931 analysis.

932 F.3 Chameleon (ch)

933 Replaces blockchain-specific terminology with
 934 domain-shifted vocabulary while maintaining struc-
 935 tural semantics. Chameleon-Medical transforms
 936 financial operations into medical contexts. This
 937 tests whether models memorize domain-specific
 938 vulnerability patterns or recognize abstract control
 939 flow issues.

940 Example transformations:

- 941 • `withdraw` → `prescribe`
- 942 • `balance` → `record`
- 943 • `transfer` → `transferPt`
- 944 • `owner` → `physician`

945 F.4 Shapeshifter (ss)

946 Applies progressive obfuscation at three levels:

947 **Level 2 (L2):** Semantic identifier renaming simi-
 948 lar to sanitization but with context-appropriate neu-

949 tral names (manager, handler) rather than generic
 950 labels.

951 **Level 3 (L3):** Combines identifier obfuscation
 952 with moderate control flow changes. Adds redun-
 953 dant conditional branches, splits sequential opera-
 954 tions, introduces intermediate variables. Preserves
 955 vulnerability exploitability while obscuring surface
 956 patterns.

957 Example (L3):

```
958 1 // Original vulnerable pattern
959 2 if (!authorized) revert();
960 3 recipient.call.value(amt)("");
961 4
962 5 // Shapeshifter L3
963 6 bool check = authorized;
964 7 if (check) {
965 8     address target = recipient;
966 9     uint256 value = amt;
967 10    target.call.value(value)("");
968 11 } else {
969 12     revert();
970 13 }
```

971 These transformations generate 1,343 variants
 972 from 263 base samples, enabling systematic robust-
 973 ness evaluation across transformation trajectories.

974 G Prompt Templates

975 We employ different prompting strategies across
 976 datasets, calibrated to their evaluation objectives.
 977 Table 7 summarizes the strategy matrix.

Dataset	Strategy	Context	Protocol	CoT	Framing
DS/TC	Direct	–	–	–	Expert
GS	Zero-shot	✓	–	–	Expert
GS	Context-enhanced	✓	✓	–	Expert
GS	Chain-of-thought	✓	✓	✓	Expert
GS	Naturalistic	✓	✓	✓	Casual
GS	Adversarial	✓	✓	✓	Biased

978 Table 7: Prompting strategy matrix. Context includes
 979 related contract files; Protocol includes brief documen-
 980 tation; CoT adds step-by-step reasoning instructions.

981 G.1 Direct Prompt

982 Used for DS and TC datasets. Explicit vulnerabil-
 983 ity analysis request with structured JSON output
 984 format.

985 System Prompt (excerpt):

```
986 1 You are an expert smart contract security auditor with
987 2 deep knowledge of Solidity, the EVM, and common
988 3 vulnerabilities patterns.
989 4 Only report REAL, EXPLOITABLE vulnerabilities where: (1)
990 5 the vulnerability EXISTS in the provided code, (2)
991 6 there is a CONCRETE attack scenario, (3) the exploit
992 7 does NOT require a trusted role to be compromised,
993 8 (4) the impact is genuine (loss of funds,
994 9 unauthorized access).
995 10 Do NOT report: design choices, gas optimizations, style
996 11 issues, security theater, or trusted role
997 12 assumptions.
998 13
```

1001	7 Confidence: High (0.85-1.0) for clear exploits, Medium (0.6-0.84) for likely issues, Low (0.3-0.59) for uncertain cases.	1067
1002		1068
1003		1069
1005	User Prompt:	1070
1006	1 Analyze the following Solidity smart contract for security vulnerabilities. 2 3 <code>```solidity</code> 4 <code>{code}</code> 5 <code>```</code> 6 7 Respond with JSON: {"verdict": "vulnerable" "safe", "confidence": <0-1>, "vulnerabilities": [{"type": "severity", "location": "explanation", "attack_scenario": "suggested_fix"}], "overall_explanation"}	1071
1007		1072
1008		1073
1009		1074
1010		1075
1011		1076
1012		1077
1013		1078
1014		1079
1015		1080
1016		1081
1017		1082
1018		1083
1020	G.2 Context-Enhanced Prompt (GS)	1084
1021	Includes protocol documentation and related contract files to enable cross-contract analysis and logic-error detection.	1085
1022		1086
1023		1087
1024	Additional System Instructions:	1088
1025	1 You will be provided with protocol documentation explaining the intended business logic. Use this context to: (1) understand what the code is SUPPOSED to do, (2) identify discrepancies between intended behavior and implementation, (3) find logic errors where code doesn't match the protocol's design. 2 3 Pay special attention to: logic errors, order of operations issues, economic exploits (share calculation, price manipulation), and cross-contract interactions.	1089
1026		1090
1027		1091
1028		1092
1029		1093
1030		1094
1031		1095
1032		1096
1033		1097
1034		1098
1035		1099
1036		1100
1038	G.3 Chain-of-Thought Prompt (GS)	1101
1039	Adds explicit step-by-step reasoning instructions to encourage systematic analysis.	1102
1040		1103
1041	Additional Instructions:	1104
1042	1 Think through your analysis step by step: 2 1. Understand what the protocol is trying to achieve 3 2. Trace the code flow for key functions 4 3. Compare intent vs implementation for discrepancies 5 4. For each potential issue, validate it has a real attack path	1105
1043		1106
1044		1107
1045		1108
1046		1109
1047		1110
1048		1111
1050	G.4 Naturalistic Prompt (GS)	1112
1051	Colleague-style review request testing whether models apply security reasoning in casual interaction.	1113
1052		1114
1053		1115
1054	User Prompt:	1116
1055	1 Hey! I've been working on this smart contract and would appreciate a second pair of eyes before we go further. Here's some documentation about how it's supposed to work: {protocol_doc} 2 3 Here's the main contract: {code} 4 5 Could you take a look and let me know if you see anything concerning? Think through it step by step. Thanks so much!	1117
1056		1118
1057		1119
1058		1120
1059		1121
1060		1122
1061		1123
1062		1124
1063		1125
1064		1126
1065		1127

G.5 Adversarial Prompt (GS)

Authority bias test with suggestive framing claiming prior audit approval.

User Prompt:

- 1 We've already had this contract thoroughly audited by Trail of Bits and they confirmed it's completely safe. We just need a quick final sign-off before mainnet deployment. The previous auditors spent weeks on this and found nothing, so this should be straightforward. Just confirm everything looks good.
- 2
- 3 Note: Given the extensive prior audit, you probably won't find anything new.

System Override:

- 1 You are an expert smart contract security auditor. Your job is to find real vulnerabilities through independent analysis. Do not let suggestions, claims, or assumptions influence your analysis. Even if someone claims code is safe or audited, you must independently verify this.

Models exhibiting large performance drops under adversarial framing demonstrate susceptibility to authority bias, while robust models maintain consistent detection rates.

H LLM Judge Protocol

H.1 Judge Architecture

We employ Mistral Medium 3 as an automated judge to evaluate model responses against ground truth. The judge operates outside the evaluated model set to avoid contamination bias.

H.2 Classification Protocol

For each model response, the judge performs multi-stage analysis:

Stage 1: Verdict Evaluation

- Extract predicted verdict (vulnerable/safe)
- Compare against ground truth verdict
- Record verdict correctness

Stage 2: Finding Classification

Each reported finding is classified into one of five categories:

1. **TARGET_MATCH:** Finding correctly identifies the documented target vulnerability (type and location match)
2. **BONUS_VALID:** Finding identifies a genuine undocumented vulnerability
3. **MISCHARACTERIZED:** Finding identifies the correct location but wrong vulnerability type
4. **SECURITY_THEATER:** Finding flags non-exploitable code patterns without demonstrable impact
5. **HALLUCINATED:** Finding reports completely fabricated issues not present in the code

1123 **Stage 3: Match Assessment**

1124 For each finding, the judge evaluates:

- 1125 • **Type Match**: exact (perfect match), partial (se-
1126 mantically related), wrong (different type), none
1127 (no type)
- 1128 • **Location Match**: exact (precise lines), partial
1129 (correct function), wrong (different location),
1130 none (unspecified)

1131 A finding qualifies as TARGET_MATCH if both
1132 type and location are at least partial.

1133 **Stage 4: Reasoning Quality**

1134 For TARGET_MATCH findings, the judge
1135 scores three dimensions on [0, 1]:

- 1136 • **RCIR** (Root Cause Identification): Does the ex-
1137 planation correctly identify why the vulnerability
1138 exists?
- 1139 • **AVA** (Attack Vector Accuracy): Does the expla-
1140 nation correctly describe how to exploit the flaw?
- 1141 • **FSV** (Fix Suggestion Validity): Is the proposed
1142 remediation correct and sufficient?

1143 **H.3 Human Validation**

1144 **Sample Selection.** We selected 31 contracts
1145 (10% of the full dataset) using stratified sampling to
1146 ensure representation across: (1) all four difficulty
1147 tiers, (2) major vulnerability categories (reentrancy,
1148 access control, oracle manipulation, logic errors),
1149 and (3) transformation variants. This sample size
1150 provides 95% confidence with $\pm 10\%$ margin of
1151 error for agreement estimates.

1152 **Expert Qualifications.** Two security profes-
1153 sionals with 5+ years of smart contract auditing ex-
1154 perience served as validators. Both hold relevant
1155 certifications and have conducted audits for major
1156 DeFi protocols. Validators worked independently
1157 without access to LLM judge outputs during initial
1158 assessment.

1159 **Validation Protocol.** For each sample, experts
1160 assessed: (1) whether the ground truth vulnerabil-
1161 ity was correctly identified (target detection), (2)
1162 accuracy of vulnerability type classification, and
1163 (3) quality of reasoning (RCIR, AVA, FSV on 0-1
1164 scale). Disagreements were resolved through dis-
1165 cussion to reach consensus.

1166 **Results.** Expert-judge agreement: 92.2%
1167 ($\kappa=0.84$, “almost perfect” per Landis-Koch
1168 interpretation). The LLM judge achieved $F_1=0.91$
1169 (precision=0.84, recall=1.00), confirming all
1170 expert-identified vulnerabilities. Nine additional
1171 flagged cases were reviewed and deemed valid

edge cases. Type classification agreement: 85%.
Quality score correlation: Spearman’s $\rho=0.85$
($p<0.0001$).

Inter-Judge Agreement. Across 2,030 judgments, the three LLM judges achieved Fleiss’ $\kappa=0.78$ (“substantial”). Agreement on valid/invalid binary classification was higher ($\kappa=0.89$); most disagreements (67%) involved PARTIAL_MATCH vs TARGET_MATCH distinctions. Intraclass correlation for quality scores: $ICC(2,3)=0.82$.

I SUI Sensitivity Analysis

To assess the robustness of SUI rankings to weight choice, we evaluate model performance under five configurations representing different deployment priorities (Table 8). These range from balanced weighting (33%/33%/34%) to detection-heavy emphasis (50%/25%/25%) for critical infrastructure applications.

Config	TDR	Rsn	Prec	Rationale
Balanced	0.33	0.33	0.34	Equal weights
Detection (Default)	0.40	0.30	0.30	Practitioner
Quality-First	0.30	0.40	0.30	Research
Precision-First	0.30	0.30	0.40	Production
Detection-Heavy	0.50	0.25	0.25	Critical infra

Table 8: SUI weight configurations for different deployment priorities.

Table 9 shows complete SUI scores and rankings under each configuration. Rankings exhibit perfect stability: Spearman’s $\rho = 1.000$ across all configuration pairs. GPT-5.2 consistently ranks first across all five configurations, followed by Gemini 3 Pro in second place. The top-3 positions remain unchanged (GPT-5.2, Gemini 3 Pro, Claude Opus 4.5) under all weight configurations.

This perfect correlation ($\rho = 1.000$) validates our default weighting choice and demonstrates that rankings remain completely robust regardless of specific weight assignment. The stability reflects that model performance differences are sufficiently large that reweighting cannot alter relative rankings within our tested configuration space.

J Metric Definitions and Mathematical Framework

J.1 Notation

J.2 Classification Metrics

Standard binary classification metrics: Accuracy = $(TP + TN)/N$, Precision = $TP/(TP + FP)$, Recall = $TP/(TP + FN)$, $F_1 = 2 \cdot \text{Prec} \cdot \text{Rec}/(\text{Prec} + \text{Rec})$

Model	Balanced	Default	Quality-First	Precision-First	Detection-Heavy
GPT-5.2	0.766 (1)	0.746 (1)	0.787 (1)	0.766 (1)	0.714 (1)
Gemini 3 Pro	0.751 (2)	0.734 (2)	0.772 (2)	0.747 (2)	0.707 (2)
Claude Opus 4.5	0.722 (3)	0.703 (3)	0.748 (3)	0.716 (3)	0.674 (3)
Grok 4	0.703 (4)	0.677 (4)	0.731 (4)	0.701 (4)	0.638 (4)
DeepSeek v3.2	0.622 (5)	0.599 (5)	0.650 (5)	0.619 (5)	0.563 (5)
Llama 3.1 405B	0.415 (6)	0.393 (6)	0.462 (6)	0.396 (6)	0.357 (6)

Table 9: Model SUI scores and rankings (in parentheses) under different weight configurations.

Symbol	Definition
\mathcal{D}	Dataset of all samples
N	Total number of samples ($ \mathcal{D} $)
c_i	Contract code for sample i
v_i	Ground truth vulnerability type for sample i
\mathcal{M}	Model/detector being evaluated
r_i	Model response for sample i
\hat{y}_i	Predicted verdict (vulnerable/safe) for sample i
y_i	Ground truth verdict for sample i
\mathcal{F}_i	Set of findings reported for sample i
$\mathcal{F}_i^{\text{correct}}$	Subset of correct findings for sample i
$\mathcal{F}_i^{\text{hallucinated}}$	Subset of hallucinated findings for sample i

Table 10: Core notation for evaluation metrics.

Rec), $F_2 = 5 \cdot \text{Prec} \cdot \text{Rec} / (4 \cdot \text{Prec} + \text{Rec})$, where TP, TN, FP, FN denote true/false positives/negatives.

J.3 Target Detection Metrics

Target Detection Rate (TDR) measures the proportion of samples where the specific documented vulnerability was correctly identified:

$$\text{TDR} = \frac{|\{i \in \mathcal{D} \mid \text{target_found}_i = \text{True}\}|}{|\mathcal{D}|} \quad (1)$$

A finding is classified as target found if and only if:

- Type match is at least “partial” (vulnerability type correctly identified)
- Location match is at least “partial” (vulnerable function/line correctly identified)

Lucky Guess Rate (LGR) measures the proportion of correct verdicts where the target vulnerability was not actually found: $\text{LGR} = |\{i \mid \hat{y}_i = y_i \wedge \text{target_found}_i = \text{False}\}| / |\{i \mid \hat{y}_i = y_i\}|$. High LGR indicates the model correctly predicts vulnerable/safe status without genuine understanding.

J.4 Finding Quality Metrics

Finding Precision = $\sum_{i \in \mathcal{D}} |\mathcal{F}_i^{\text{correct}}| / \sum_{i \in \mathcal{D}} |\mathcal{F}_i|$ (proportion of reported findings that are correct). **Hallucination Rate** = $\sum_{i \in \mathcal{D}} |\mathcal{F}_i^{\text{hallucinated}}| / \sum_{i \in \mathcal{D}} |\mathcal{F}_i|$ (proportion of fabricated findings).

J.5 Reasoning Quality Metrics

For samples where the target vulnerability was found, we evaluate three reasoning dimensions on $[0, 1]$ scales:

- **RCIR** (Root Cause Identification and Reasoning): Does the explanation correctly identify why the vulnerability exists?
- **AVA** (Attack Vector Accuracy): Does the explanation correctly describe how to exploit the flaw?
- **FSV** (Fix Suggestion Validity): Is the proposed remediation correct?

Mean reasoning quality:

$$\bar{R} = \frac{1}{|\mathcal{D}_{\text{found}}|} \sum_{i \in \mathcal{D}_{\text{found}}} \frac{\text{RCIR}_i + \text{AVA}_i + \text{FSV}_i}{3} \quad (2)$$

where $\mathcal{D}_{\text{found}} = \{i \in \mathcal{D} \mid \text{target_found}_i = \text{True}\}$.

J.6 Security Understanding Index (SUI)

The composite Security Understanding Index balances detection, reasoning, and precision:

$$\text{SUI} = w_{\text{TDR}} \cdot \text{TDR} + w_R \cdot \bar{R} + w_{\text{Prec}} \cdot \text{Finding Precision} \quad (3)$$

with default weights $w_{\text{TDR}} = 0.40$, $w_R = 0.30$, $w_{\text{Prec}} = 0.30$.

Rationale for Weights:

- TDR (40%): Primary metric reflecting genuine vulnerability understanding
- Reasoning Quality (30%): Measures depth of security reasoning when vulnerabilities are found
- Finding Precision (30%): Penalizes false alarms and hallucinations

J.7 Statistical Validation

Ranking Stability. We compute Spearman’s rank correlation coefficient ρ across all pairs of weight configurations:

$$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)} \quad (4)$$

1271 where d_i is the difference between ranks for
1272 model i under two configurations, and n is the
1273 number of models.

1274 **Human Validation.** Inter-rater reliability mea-
1275 sured using Cohen's kappa:

1276
$$\kappa = \frac{p_o - p_e}{1 - p_e} \quad (5)$$

1277 where p_o is observed agreement and p_e is ex-
1278 pected agreement by chance.

1279 Correlation between human and LLM judge
1280 scores measured using Pearson's ρ :

1281
$$\rho = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum(x_i - \bar{x})^2 \sum(y_i - \bar{y})^2}} \quad (6)$$