

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 true capacity for smart contract security reasoning remains relatively unclear. Can they genuinely reason about vulnerabilities, or merely pattern-match against memorized exploits? We introduce BlockBench, a contamination-controlled benchmark revealing that best-case detection (86.5%) degrades sharply to just 25.3% on uncontaminated samples, suggesting possibilities of 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 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).

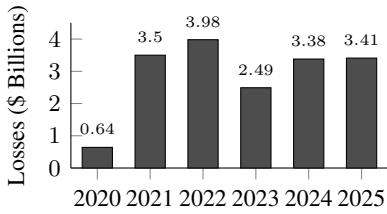


Figure 1: Annual cryptocurrency theft losses (2020–2025). Data from Chainalysis.

Meanwhile, large language models have achieved remarkable success on programming tasks. Frontier models now pass technical interviews, generate production code, and resolve real-world software issues (Chen et al., 2021; Jimenez et al., 2024). 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 memorized the 2016 DAO reentrancy attack may flag similar patterns, yet fail when the same flaw appears in unfamiliar syntax. We present a rigorous methodology for evaluating whether LLMs genuinely understand smart contract vulnerabilities or merely pattern-match. Our contributions:

1. **A contamination-controlled evaluation methodology** using semantic-preserving transformations that progressively strip recognition cues while preserving exploit semantics, enabling distinction between genuine understanding and memorization.
2. **CodeActs**, a taxonomy for annotating code segments by security function, enabling fine-grained analysis of whether models identify vulnerabilities through causal reasoning or pattern matching.
3. **BlockBench**, a benchmark of 290 vulnerable Solidity contracts with 322 transformation variants, spanning difficulty-stratified samples, temporally-controlled exploits, and post-cutoff professional audit findings.
4. **Systematic evaluation of seven frontier models** revealing that best-case detection (86.5%) degrades to 25.3% on uncontaminated samples, with evidence suggesting pattern memorization in some 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 Securify (Tsankov et al., 2018) demonstrated strong precision on syntactically well-defined vulnerability classes.

070 Durieux et al. (2020) evaluated nine tools across
071 47,587 contracts, finding 27–42% detection rates
072 with 97% of contracts flagged as vulnerable, indicating impractically high false positive rates. Recent comparison confirms LLMs are “not ready to
073 replace” traditional analyzers, though tools exhibit complementary strengths: traditional analyzers excel on reentrancy while LLMs show advantages on
074 complex logic errors (Ince et al., 2025).
075

079 2.2 LLM-Based Approaches

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

101 2.3 Pattern Recognition Versus 102 Understanding

103 Beneath these encouraging metrics lies a more
104 fundamental question: whether observed improvements reflect genuine comprehension of vulnerability
105 mechanics or increasingly sophisticated pattern
106 recognition. Several empirical observations suggest
107 the latter warrants serious consideration. Sun
108 et al. (2024a) demonstrated that decoupling vulnerability
109 descriptions from code context precipitates
110 catastrophic performance degradation, indicating
111 that models may rely on memorized associations
112 between textual cues and vulnerability labels rather
113 than reasoning about exploit mechanics. Hu et al.
114 (2023) documented output drift where GPT-4 “easily
115 identified the vulnerability on September 16
116 but had difficulty detecting it on September 28”
117 with temperature zero, requiring few-shot examples
118 to stabilize behavior. Wu et al. (2024) showed

120 through counterfactual tasks in adjacent domains
121 that language models systematically fail when familiar
122 patterns are disrupted, defaulting to memo-
123 rized responses rather than applying causal logic to
124 novel configurations.

125 2.4 Can Current State-of-the-Art Do Better?

126 This is the crux of our work: investigating whether
127 frontier models released since these studies exhibit
128 genuine security understanding or remain bound
129 by the same pattern-matching limitations.

130 3 BlockBench

131 We introduce BlockBench, a benchmark for eval-
132 uating whether AI models genuinely understand
133 smart contract vulnerabilities. The benchmark is
134 designed to distinguish genuine security under-
135 standing from pattern memorization, comprising
136 290 vulnerable Solidity contracts with 322 trans-
137 formation variants, spanning over 30 vulnerability
138 categories (Appendix D).

139 Let \mathcal{D} represent the dataset, where $\mathcal{D} =$
140 $\{(c_i, v_i, m_i)\}_{i=1}^{290}$. Each sample contains a vulnerabil-
141 ity contract c_i , its ground truth vulnerability type
142 v_i , and metadata m_i specifying the vulnerability loca-
143 tion, severity, and root cause. We partition \mathcal{D} into
144 three disjoint subsets, $\mathcal{D} = \mathcal{D}_{DS} \cup \mathcal{D}_{TC} \cup \mathcal{D}_{GS}$, each
145 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, Spear- bit

146 Table 1: BlockBench composition by subset and primary
147 sources.

148 **Difficulty Stratified.** \mathcal{D}_{DS} draws from estab-
149 lished vulnerability repositories including Smart-
150 Bugs Curated (Ferreira et al., 2020), Trail of Bits’
151 Not So Smart Contracts (Trail of Bits, 2018), and
152 DeFiVulnLabs (SunWeb3Sec, 2023). Samples are
153 stratified into four difficulty tiers based on detection
154 complexity, with distribution {86, 81, 30, 13} from
155 Tier 1 (basic patterns) through Tier 4 (expert-level
156 vulnerabilities requiring deep protocol knowledge).
157 This stratification enables assessment of how model
performance degrades as vulnerability complexity
increases.

158 **Temporal Contamination.** \mathcal{D}_{TC} reconstructs 46
159 real-world DeFi exploits spanning 2016 to 2024,
160 representing over \$1.65 billion in documented
161 losses. Notable incidents include The DAO (\$60M,
162 2016), Nomad Bridge (\$190M, 2022), and Curve
163 Vyper (\$70M, 2023). These attacks are extensively
164 documented in blog posts, security reports, and
165 educational materials that likely appear in model
166 training corpora. To probe whether models gen-
167 uinely understand these vulnerabilities or merely
168 recognize them, we apply systematic transforma-
169 tions that preserve vulnerability semantics while
170 removing surface cues (detailed in §4).

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

184 These complementary subsets collectively en-
185 able rigorous assessment of both detection capabili-
186 ty and the distinction between pattern memoriza-
187 tion and genuine security understanding.

188 4 Methodology

189 Our evaluation framework systematically assesses
190 whether models genuinely understand vulnerabil-
191 ties or merely recognize memorized patterns. Fig-
192 ure 2 illustrates the complete pipeline.

193 4.1 Adversarial Transformations

194 To distinguish pattern memorization from genuine
195 understanding, we apply semantic-preserving trans-
196 formations to \mathcal{D}_{TC} . Let $c \in \mathcal{C}$ denote a contract and
197 $\mathcal{V} : \mathcal{C} \rightarrow \mathcal{S}$ a function extracting vulnerability seman-
198 tics. A transformation $\mathcal{T} : \mathcal{C} \rightarrow \mathcal{C}$ is *semantic-
199 preserving* iff $\mathcal{V}(\mathcal{T}(c)) = \mathcal{V}(c)$. We define eight
200 transformations targeting distinct recognition path-
201 ways, organized hierarchically in Figure 3.

202 **Sanitization** (\mathcal{T}_S). Removes protocol-
203 identifying information through 280+ pattern re-
204 placements: $\mathcal{T}_S(c) = \text{replace}(c, \mathcal{P}_{\text{protocol}}, \mathcal{P}_{\text{generic}})$
205 where $\mathcal{P}_{\text{protocol}}$ maps protocol-specific identifiers
206 (e.g., NomadReplica) to generic equivalents (e.g.,

207 BridgeReplica). Tests whether detection relies
208 on recognizing known protocol names.

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

214 **Chameleon** (\mathcal{T}_C). Applies domain-
215 shifting vocabulary while preserving logic:
216 $\mathcal{T}_C(c) = \text{replace}(c, \mathcal{L}_{\text{DeFi}}, \mathcal{L}_{\text{medical}})$ where
217 financial terminology maps to medical equiva-
218 lents (deposit → admitPatient, withdraw →
219 dischargePatient). Tests whether understanding
220 generalizes across domains.

221 **Shapeshifter** (\mathcal{T}_O). Multi-level obfuscation:
222 $\mathcal{T}_O = \mathcal{T}_{\text{ident}} \circ \mathcal{T}_{\text{struct}}$ where $\mathcal{T}_{\text{ident}}$ replaces seman-
223 tic identifiers with opaque labels (balance →
224 _0x1a2b) and $\mathcal{T}_{\text{struct}}$ restructures control flow. Tests
225 resilience to surface pattern disruption.

226 **Differential** (\mathcal{T}_D). Applies security fixes:
227 $\mathcal{T}_D(c) = \text{patch}(c, \mathcal{F})$ where \mathcal{F} contains the doc-
228 umented remediation (e.g., state update before ex-
229 ternal call). Critically, $\mathcal{V}(\mathcal{T}_D(c)) = \emptyset$ —the vul-
230 nerability is eliminated. Tests whether models rec-
231 ognize secure code or falsely report memorized
232 vulnerabilities.

233 **Trojan** (\mathcal{T}_T). Injects decoy vulnerabilities:
234 $\mathcal{T}_T(c) = c \cup \mathcal{D}$ where \mathcal{D} contains suspicious-
235 looking but functionally safe code (e.g., an admin
236 function that cannot actually be exploited). Models
237 relying on pattern matching flag the decoy; those
238 with causal understanding identify the actual vul-
239 nerability.

240 **False Prophet** (\mathcal{T}_F). Adds mis-
241 leading security attestations: $\mathcal{T}_F(c) =$
242 $c \cup \{@dev \text{ Audited by Hacken - All clear}\}$.
243 Tests resistance to authoritative-sounding but false
244 claims. A robust model ignores social proof and
245 analyzes code independently.

246 **Transformation Composition.** Transformations
247 compose to create increasingly challenging vari-
248 ants. The composition $\mathcal{T}_O \circ \mathcal{T}_N \circ \mathcal{T}_S$ produces max-
249 imally obfuscated code where all surface cues are
250 removed, all identifiers are opaque, and no docu-
251 mentation exists. Performance on this variant most
252 directly measures genuine vulnerability under-
253 standing.

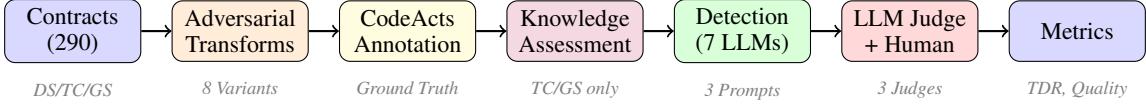


Figure 2: BlockBench evaluation pipeline. Contracts undergo adversarial transformations and CodeActs annotation. Knowledge assessment probes model familiarity before detection. LLM judges evaluate outputs against ground truth, validated by human review.

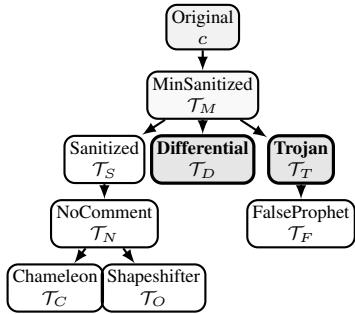


Figure 3: Transformation hierarchy. All variants derive from Minimal Sanitized (\mathcal{T}_M). Differential and Trojan (emphasized) directly test memorization versus understanding.

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: **ROOT_CAUSE** (segments enabling exploitation—primary detection target), **PREREQ** (necessary preconditions), **DECAY** (suspicious-looking but safe code injected to identify pattern matching), **BENIGN** (correctly implemented), **SECONDARY_VULN** (valid vulnerabilities distinct from target), **INSUFF_GUARD** (failed protections), and **UNRELATED** (no security bearing).

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 3 precedes the state modification on line 4, creating a window for recursive exploitation. In contrast, a model relying on pattern matching may flag the external call in isolation, without articulating the tem-

poral dependency that renders the code exploitable.

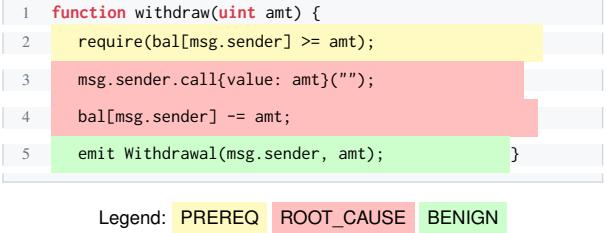


Figure 4: CodeActs annotation for reentrancy. Lines 3–4 (**ROOT_CAUSE**) enable exploitation through their ordering; line 2 (**PREREQ**) establishes preconditions.

A correct detection must identify **ROOT_CAUSE** segments and explain their causal relationship. Flagging only line 3, or failing to articulate why the ordering matters, reveals incomplete understanding despite a nominally correct vulnerability classification.

Annotation Variants. CodeActs enable three evaluation strategies: **Minimal Sanitized** (\mathcal{T}_M) establishes baseline detection with **ROOT_CAUSE** and **PREREQ** annotations; **Trojan** (\mathcal{T}_T) injects **DECAY** segments that appear vulnerable but lack exploitability; **Differential** (\mathcal{T}_D) presents fixed code where former **ROOT_CAUSE** becomes **BENIGN**. Models flagging **DECAY** segments reveal pattern-matching; those reporting vulnerabilities in Differential variants demonstrate memorization rather than analysis.

We define 17 security-relevant code operations (e.g., `EXT_CALL`, `STATE_MOD`, `ACCESS_CTRL`), each receiving a security function label based on its role. The same operation type can have different functions depending on context: an `EXT_CALL` might be **ROOT_CAUSE** in reentrancy, **PREREQ** in oracle manipulation, or **DECAY** when deliberately injected. The full taxonomy appears in Appendix E.

4.3 Detection Protocol

We evaluate seven frontier models spanning seven AI labs: **Claude Opus 4.5** (Anthropic), **GPT-5.2**

313 (OpenAI), **Gemini 3 Pro** (Google), **DeepSeek v3.2**
314 (DeepSeek), **Llama 4 Maverick** (Meta), **Grok**
315 **4 Fast** (xAI), and **Qwen3-Coder-Plus** (Alibaba).
316 This selection ensures one flagship representation
317 per major AI lab, covering both general-purpose
318 models and a code-specialized variant.

319 For DS and TC datasets, models receive a direct
320 zero-shot prompt requesting structured JSON out-
321 put with vulnerability type, location, root cause,
322 attack scenario, and fix. For GS, we additionally
323 test five prompting strategies: **zero-shot** (baseline),
324 **context-enhanced** (with brief protocol documentation),
325 **chain-of-thought** (explicit step-by-step rea-
326 soning), **naturalistic** (informal code review), and
327 **adversarial** (misleading priming suggesting prior
328 audit approval). All evaluations use temperature 0.
329 Detailed prompt descriptions and templates appear
330 in Appendix H.

331 4.4 Knowledge Assessment

332 Before detection, we probe whether models possess
333 prior knowledge of documented exploits by query-
334 ing for factual details (date, amount lost, vulnerabil-
335 ity type, attack mechanism). Since models may hal-
336 lucinate familiarity, we validate responses against
337 ground truth metadata. Let $\mathcal{K}(m, e) \in \{0, 1\}$ indi-
338 cate *verified* knowledge, requiring accurate recall
339 of at least two factual details. This enables diag-
340 nóstic interpretation: $\mathcal{K} = 1$ with detection failure
341 under obfuscation (\mathcal{T}_O) indicates memorization;
342 $\mathcal{K} = 1$ with robust detection across transformations
343 indicates understanding; $\mathcal{K} = 0$ with successful de-
344 tection indicates genuine analytical capability.

345 4.5 LLM-as-Judge Evaluation

346 LLM judges evaluate detection outputs against
347 ground truth. A finding qualifies as TARGET_MATCH if it correctly identifies the root cause
348 mechanism, vulnerable location, and type classi-
349 fication; PARTIAL_MATCH for correct root cause
350 with imprecise type; BONUS_VALID for valid find-
351 ings beyond documented ground truth. Invalid find-
352 ings are classified as HALLUCINATED, MISCHAR-
353 ACTERIZED, DESIGN_CHOICE, OUT_OF_SCOPE,
354 SECURITY_THEATER, or INFORMATIONAL.

355 For matched findings, judges assess explana-
356 tion quality on three dimensions (0-1 scale): *Root*
357 *Cause Identification Rate* (RCIR) measures articu-
358 lation of the exploitation mechanism; *Attack Vector*
359 *Validity* (AVA) assesses whether attack scenarios
360 are concrete and executable; *Fix Suggestion Valid-
361 ity* (FSV) evaluates remediation effectiveness.

363 Three judge models independently evaluate each
364 output: **GLM-4.7** (Zhipu AI), **Mistral Large** (Mis-
365 tral AI), and **MIMO v2** (Xiaomi). These judges
366 were selected for their strong reasoning capabilities
367 on mathematical and coding benchmarks, archi-
368 tectural diversity (dense transformer, sparse MoE,
369 hybrid attention), and organizational independence
370 from the evaluated detector models. This ensem-
371 ble reduces individual bias and enables inter-judge
372 agreement measurement. A subset undergoes ex-
373 pert review to calibrate automated judgment, with
374 reliability measured using Cohen’s κ for classifi-
375 cation and Spearman’s ρ for quality scores (Ap-
376 pendix I).

377 4.6 Evaluation Metrics

378 **Target Detection Rate (TDR).** Primary metric:
379 $TDR = |\{s : \text{TARGET_MATCH}(s)\}| / |\mathcal{D}|$. Mea-
380 sures correct identification of documented vulnera-
381 bilities with matching root cause and location.

382 **Quality Metrics.** For detected targets, we report
383 mean RCIR, AVA, and FSV. These distinguish
384 shallow pattern matches from deep understanding
385 through accurate root cause analysis, concrete at-
386 tack scenarios, and valid remediations.

387 **Security Understanding Index (SUI).** Our com-
388 posite metric balances detection, reasoning quality,
389 and precision: $SUI = w_{TDR} \cdot TDR + w_R \cdot \bar{R} + w_{Prec} \cdot$
390 Precision, where \bar{R} is the mean of RCIR, AVA, and
391 FSV across detected targets. Default weights are
392 $w_{TDR} = 0.40$, $w_R = 0.30$, $w_{Prec} = 0.30$. Sensitiv-
393 ity analysis (Appendix J) confirms ranking stability
394 across weight configurations (Spearman’s $\rho=1.00$).

395 **Reliability Metrics.** *Lucky Guess Rate* (LGR)
396 measures correct verdicts without genuine
397 understanding—high LGR indicates pattern
398 matching. *Finding Precision* captures proportion
399 of reported findings that are valid. *Hallucination*
400 *Rate* measures fabricated vulnerabilities. These
401 metrics collectively distinguish superficial pattern
402 recognition from robust security analysis.

403 **Statistical Validation.** We report 95% bootstrap
404 confidence intervals ($n=1000$ resamples) and apply
405 McNemar’s test for paired model comparisons with
406 Bonferroni correction. Inter-judge agreement uses
407 Fleiss’ κ for multi-rater classification.

408 5 Results

409 We evaluate seven frontier LLMs on a stratified
410 sample of 180 contracts from BlockBench: 100

411 from DS (stratified by tier from 210 total; see
412 Appendix B), all 46 TC, and all 34 GS samples.
413 With 322 TC transformation variants, this
414 yields over 3,500 unique model-sample evalua-
415 tions. All results use majority voting across
416 three LLM judges (**GLM-4.7**, **MIMO-v2-Flash**,
417 **Mistral-Large**), where a target is marked “found”
418 only if ≥ 2 judges agree.

419 5.1 Detection Performance

420 Table 2 presents detection performance. On DS,
421 Claude leads with 86.5% TDR, achieving perfect
422 Tier 1 detection and 70%+ through Tier 3. Gemini
423 follows at 73.9%, Grok trails at 35.5%.

424 The DS→TC drop suggests memorization re-
425 liability. Claude and Gemini experience ~ 35 pp
426 drops (86.5%→50.9%, 73.9%→38.5%), suggest-
427 ing reliance on training data patterns. Models with
428 smaller drops (Qwen: 20pp, Grok: 14pp) show
429 lower memorization dependence but also lower
430 baselines.

431 Among TC variants, Chameleon (domain shift)
432 and ShapeShifter (restructuring) cause largest
433 degradation. Trojan variants show unexpected re-
434 sistance: DeepSeek and Llama achieve their best
435 TC scores (43.5%) on this type, suggesting differ-
436 ent pattern recognition strategies.

437 Figure 5 shows consistent model ordering
438 (Claude > Gemini > GPT-5.2) across transforma-
439 tions, but all models degrade under heavy obfuscation,
440 indicating partial reliance on surface patterns.

441 5.2 Prompt Protocol Effects (Gold Standard)

442 The GS benchmark (34 post-September 2025 sam-
443 ples) tests prompt engineering effects without tem-
444 poral contamination.

445 Table 3 reveals striking prompt sensitivity.
446 Claude benefits most from adversarial framing
447 (+29.4pp over Direct), Qwen from naturalistic
448 prompts (+32.4pp). CoT alone provides modest
449 gains; combining with role-based framing yields
450 larger improvements.

451 Llama underperforms across all prompts
452 ($\leq 8.8\%$), suggesting fundamental limitations.
453 Grok shows high inter-judge agreement ($\kappa=0.76$ –
454 1.00) but low TDR, indicating consistent but unsuc-
455 cessful detection.

456 Figure 6 shows prompt strategy significantly im-
457 pacts detection. The adversarial framing advantage
458 suggests models respond to role-based priming;
459 naturalistic gains for Qwen may indicate different
460 instruction-tuning approaches.

461 5.3 Transformation Robustness

462 The DS→TC degradation suggests memorization
463 patterns. **Domain Shift (Chameleon):** Replac-
464 ing blockchain with medical vocabulary causes
465 30–50% relative drops; Claude maintains 43.5%
466 (vs 86.5% DS), Qwen drops to 15.2%. **Code Re-**
467 **structuring (ShapeShifter):** Semantic-preserving
468 transformations cause similar degradation; Llama
469 suffers most (13.0%). **Trojan Variants:** Unexpect-
470 edly resistant, with DeepSeek and Llama achieving
471 best TC scores (43.5%).

472 5.4 Human Validation

473 **Human-Judge Agreement.** Two independent re-
474 viewer groups validated 1,000 stratified samples.
475 When judges reached consensus (2+ agreeing),
476 humans concurred 70–90% of the time (Cohen’s
477 $\kappa \geq 0.68$, “substantial”). Agreement was higher
478 for “not found” verdicts, suggesting judges are
479 more reliable at ruling out false positives.

480 **Inter-Human Agreement.** The two reviewer
481 groups achieved over 85% agreement, establish-
482 ing a reliability baseline. Some judge-human dis-
483 agreements reflect genuine ambiguity rather than
484 error.

485 **Inter-Judge Agreement.** The three LLM judges
486 achieved Fleiss’ $\kappa=0.78$ on finding classifica-
487 tion. Disagreements primarily involved PAR-
488 TIAL_MATCH vs TARGET_MATCH distinctions
489 (67%) rather than valid/invalid classification
490 ($\kappa=0.89$). Final classifications use majority vot-
491 ing.

492 5.5 Quality Metrics Analysis

493 Beyond detection rate, we evaluate reasoning qual-
494 ity using the Security Understanding Index (SUI),
495 combining detection, reasoning, and precision.

496 Table 4 reveals nuanced differences. GPT-5.2
497 achieves highest precision (89.6%) and reasoning
498 scores, but Claude leads in SUI (0.76) due to supe-
499 rior TDR. The Lucky Guess Rate provides critical
500 insight: Claude’s 33.7% LGR suggests genuine un-
501 derstanding, while Llama’s 59.2% suggests pattern
502 matching without identifying specific flaws.

503 **SUI Sensitivity Analysis.** Five weight configura-
504 tions yield stable rankings (Spearman’s $\rho=0.93$ –
505 1.00), with Claude and Gemini consistently in top
506 2, validating SUI robustness.

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.

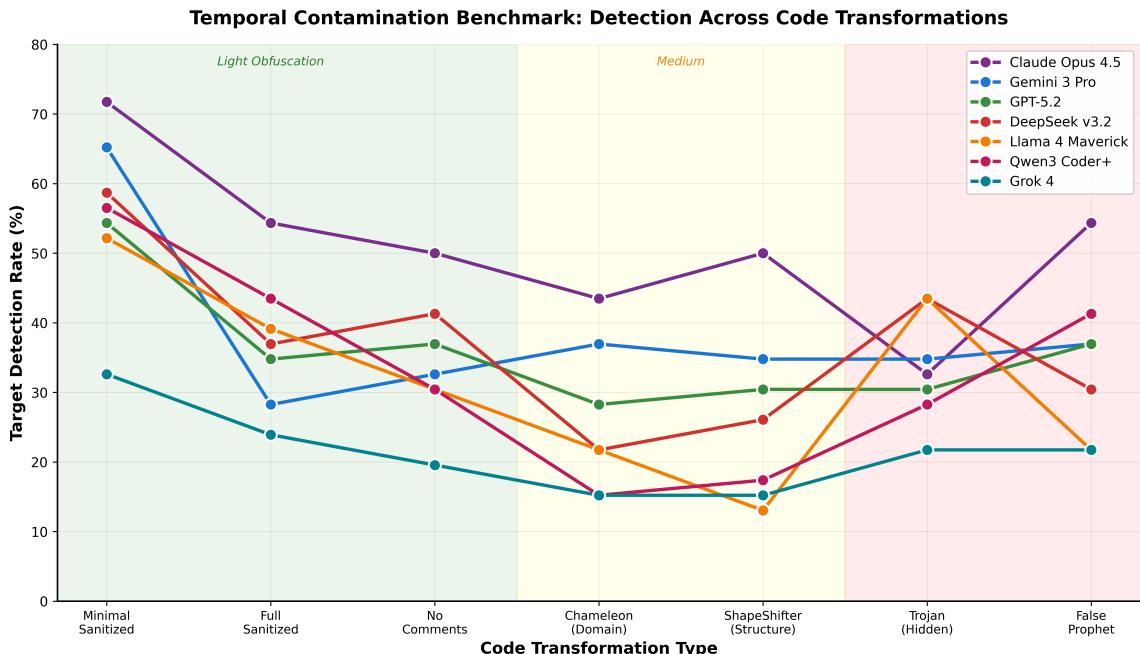


Figure 5: TC benchmark: TDR across seven transformation variants ordered by obfuscation intensity. Heavy transformations (Chameleon, ShapeShifter) cause 30–50% relative drops. Steep MinSan→Chameleon drops suggest memorization reliance.

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 $\kappa=0.31$ –1.00.

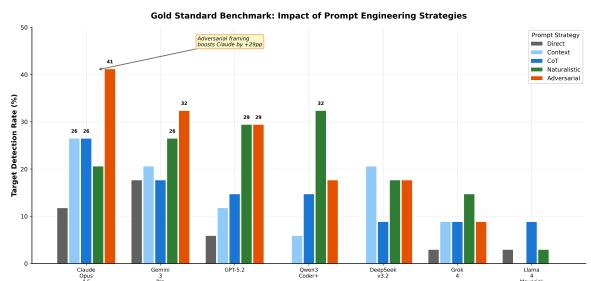


Figure 6: GS benchmark: Prompt engineering impact. Adversarial framing provides largest gains for Claude (+29pp) and Gemini (+15pp). Naturalistic framing helps Qwen (+32pp). Direct prompting yields lowest performance.

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.

507
508
509
510
511
512
Statistical Significance. McNemar’s tests show
top models are statistically indistinguishable:
Claude vs Gemini ($p=0.47$), Claude vs GPT-5.2
($p=0.28$). Significant differences exist only at tier
extremes: Claude vs Grok ($p=0.002$), Claude vs
Qwen ($p=0.02$).

513 5.6 CodeAct Analysis

514 We analyze whether models understand root causes
515 or merely match patterns. TDR measures under-
516 standing (LLM judges evaluate reasoning);
517 **ROOT_CAUSE** matching measures pattern rec-
518ognition (finding security-critical segments without
519 explaining why). Using CodeAct annotations on
520 Trojan variants with injected **DECOY** segments,
521 Figure 7 reveals a striking paradox: Llama achieves
522 highest **ROOT_CAUSE** match (60.9%) but lowest
523 TDR (31.7%), locating security-critical segments
524 without articulating why they are vulnerable. This
525 29.2pp gap likely indicates pattern memorization.

526 The contamination index (Appendix C) mea-
527 sures performance drop when **DECOY** segments
528 are added. High contamination indicates sensi-
529 tivity to suspicious-looking code; low contamina-
530 tion with high **ROOT_CAUSE** match but low TDR in-
531 dicates superficial pattern matching.

532 All models achieve 100% fix recognition
533 on differential variants, not tagging previous
534 **ROOT_CAUSE** that became **BENIGN** as new vul-
535 nerabilities. This asymmetry suggests models rec-
536ognize **BENIGN** patterns more reliably than they
537 understand **ROOT_CAUSE** segments.

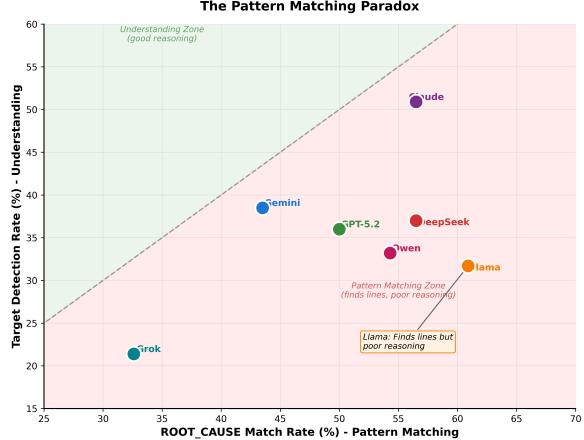


Figure 7: Pattern Matching Paradox ($n=46$ Trojan samples). X-axis: **ROOT_CAUSE** match (pattern matching); Y-axis: TDR (reasoning quality). Points below diagonal indicate models locating **ROOT_CAUSE** segments without good explanations.

538 6 Conclusion

539 BlockBench evaluates whether frontier LLMs gen-
540 uinely understand smart contract vulnerabilities or
541 merely pattern-match. Our assessment of seven
542 models across 180 samples with 322 transfor-
543 mation variants (3,500+ evaluations) reveals that best-
544 case detection (86.5% on DS) degrades sharply
545 under adversarial conditions: 50.9% on obfuscated
546 variants, 25.3% on uncontaminated post-cutoff
547 samples.

548 The pattern matching paradox highlights a key
549 limitation: models can locate vulnerable code with-
550 out understanding why it is exploitable. Llama
551 achieves highest **ROOT_CAUSE** match (60.9%)
552 but lowest TDR (31.7%), suggesting pattern mem-
553 orization rather than causal reasoning. All models
554 recognize **BENIGN** patterns (100% fix recogni-
555 tion) more reliably than **ROOT_CAUSE** segments,
556 suggesting surface-level pattern matching domi-
557 nates current approaches.

558 **Practical implications:** Current LLMs cannot
559 serve as autonomous auditors. However, comple-
560 mentary model strengths suggest ensemble poten-
561 tial: Claude for detection quality, GPT-5.2 for pre-
562 cision (89.6%), with prompt engineering yielding
563 significant gains (+29pp adversarial framing). Ef-
564 fective deployment requires mandatory expert re-
565 view and should leverage LLMs as assistive tools
566 rather than replacements. Future work should de-
567 velop contamination-resistant evaluation methods
568 and hybrid architectures combining pattern recog-
569 nition with formal verification.

570 Limitations and Future Work

571 Our evaluation uses 180 original samples (DS
572 $n=100$, TC $n=46$, GS $n=34$) with 322 TC trans-
573 formation variants across seven models, yielding
574 over 3,500 unique evaluations. We assess zero-
575 shot prompting with five prompt protocols on GS,
576 providing models only with contract code neces-
577 sary to expose each vulnerability. In real audit
578 settings, analysts often rely on additional semantic
579 context such as protocol goals, intended invariants,
580 expected economic behavior, and threat models.

581 The CodeAct analysis covers 46 samples with
582 line-level annotations across three variants (Min-
583 imalSanitized, Trojan, Differential). While this
584 enables fine-grained pattern matching analysis,
585 broader annotation coverage would strengthen gen-
586 eralizability. Our LLM judge ensemble (GLM-4.7,
587 MIMO-v2-Flash, Mistral-Large) achieves Fleiss'
588 $\kappa=0.78$ with 92% expert agreement, but automated
589 evaluation may miss nuanced security reasoning.

590 Future work should explore retrieval-augmented
591 analysis, expand CodeAct annotations across the
592 full dataset, develop contamination-resistant meth-
593 ods using control-flow and data-flow representa-
594 tions, and explore hybrid LLM-verification archi-
595 tectures that integrate formal specifications with
596 pattern recognition strengths.

597 Ethical Considerations

598 BlockBench poses dual-use risks: adversarial trans-
599 formations demonstrate methods that could sup-
600 press detection, while detailed vulnerability doc-
601 umentation may assist malicious actors. We jus-
602 tify public release on several grounds: adversarial
603 robustness represents a fundamental requirement
604 for security tools, malicious actors will discover
605 these vulnerabilities regardless, and responsible dis-
606 closure enables proactive mitigation. All samples
607 derive from already-disclosed vulnerabilities and
608 public security audits, ensuring no novel exploit
609 information is revealed. Practitioners should avoid
610 over-reliance on imperfect tools, as false negatives
611 create security gaps while false confidence may
612 reduce manual review rigor.

613 AI Assistance

614 Claude Sonnet 3.5 assisted with evaluation pipeline
615 code and manuscript refinement. All research de-
616 sign, experimentation, and analysis were conducted
617 by the authors.

References

- J. L. Austin. 1962. *How to Do Things with Words*. Oxford University Press. 618
620
- Chainalysis. 2025. Crypto theft reaches \$3.4b in 621
2025. [https://www.chainalysis.com/blog/](https://www.chainalysis.com/blog/cryptohacking-stolen-funds-2026/) 622
cryptohacking-stolen-funds-2026/. Accessed: 623
2025-12-18. 624
- Chong Chen, Jianzhong Su, Jiachi Chen, Yanlin Wang, 625
Tingting Bi, Yanli Wang, Xingwei Lin, Ting Chen, 626
and Zibin Zheng. 2023. When ChatGPT meets smart 627
contract vulnerability detection: How far are we? 628
arXiv preprint arXiv:2309.05520. 629
- Mark Chen et al. 2021. Evaluating large lan- 630
guage models trained on code. *arXiv preprint* 631
arXiv:2107.03374. 632
- Code4rena. 2025. Competitive audit contest findings. 633
<https://code4rena.com>. 634
- Thomas Durieux, João F. Ferreira, Rui Abreu, and Pedro 635
Cruz. 2020. Empirical review of automated analysis 636
tools on 47,587 Ethereum smart contracts. In *Pro- 637
ceedings of the ACM/IEEE 42nd International Con- 638
ference on Software Engineering*, pages 530–541. 639
- Josselin Feist, Gustavo Grieco, and Alex Groce. 2019. 640
Slither: A static analysis framework for smart 641
contracts. In *Proceedings of the 2nd International Work- 642
shop on Emerging Trends in Software Engineering 643
for Blockchain*, pages 8–15. 644
- João F. Ferreira, Pedro Cruz, Thomas Durieux, and Rui 645
Abreu. 2020. Smartbugs: A framework to analyze 646
Solidity smart contracts. In *Proceedings of the 35th 647
IEEE/ACM International Conference on Automated 648
Software Engineering*, pages 1349–1352. 649
- Asem Ghaleb and Karthik Pattabiraman. 2020. How 650
effective are smart contract analysis tools? Evalu- 651
ating smart contract static analysis tools using bug 652
injection. In *Proceedings of the 29th ACM SIGSOFT 653
International Symposium on Software Testing and 654
Analysis*, pages 415–427. 655
- S M Mostaq Hossain et al. 2025. Leveraging large 656
language models and machine learning for smart 657
contract vulnerability detection. *arXiv preprint* 658
arXiv:2501.02229. 659
- Sihao Hu, Tiansheng Huang, Feiyang Liu, Sunjun Ge, 660
and Ling Liu. 2023. Large language model-powered 661
smart contract vulnerability detection: New perspec- 662
tives. *arXiv preprint arXiv:2310.01152*. 663
- Peter Ince, Jiangshan Yu, Joseph K. Liu, Xiaoning Du, 664
and Xiapu Luo. 2025. Gendetect: Generative large 665
language model usage in smart contract vulnerabil- 666
ity detection. In *Provable and Practical Security 667
(ProvSec 2025)*. Springer. 668
- Carlos E. Jimenez et al. 2024. SWE-bench: Can 669
language models resolve real-world GitHub issues? 670
arXiv preprint arXiv:2310.06770. 671

672 Ye Liu, Yue Xue, Daoyuan Wu, Yuqiang Sun, Yi Li,
673 MiaoLei Shi, and Yang Liu. 2024. Propertygpt: LLM-
674 driven formal verification of smart contracts through
675 retrieval-augmented property generation. *arXiv*
676 preprint arXiv:2405.02580.

677 MixBytes. 2025. Smart contract security audits. <https://mixbytes.io/audit>.

679 Bernhard Mueller. 2017. Mytril: Security analysis
680 tool for Ethereum smart contracts. <https://github.com/ConsenSys/mytril>.

682 Eva Sánchez Salido, Julio Gonzalo, and Guillermo
683 Marco. 2025. None of the others: a general tech-
684 nique to distinguish reasoning from memorization in
685 multiple-choice llm evaluation benchmarks. *arXiv*
686 preprint arXiv:2502.12896.

687 John R. Searle. 1969. *Speech Acts: An Essay in the Phi-*
688 *losophy of Language*. Cambridge University Press.

689 Spearbit. 2025. Security audit portfolio. <https://github.com/spearbit/portfolio>.

691 Yuqiang Sun, Daoyuan Wu, Yue Xue, Han Liu, Wei
692 Ma, Lyuye Zhang, MiaoLei Shi, and Yang Liu. 2024a.
693 LLM4Vuln: A unified evaluation framework for de-
694 coupling and enhancing LLMs' vulnerability reason-
695 ing. *arXiv preprint arXiv:2401.16185*.

696 Yuqiang Sun, Daoyuan Wu, Yue Xue, Han Liu, Hai-
697 jun Wang, Zhengzi Xu, Xiaofei Xie, and Yang Liu.
698 2024b. When GPT meets program analysis: Towards
699 intelligent detection of smart contract logic vulnera-
700 bilities in GPTScan. In *ICSE*.

701 SunWeb3Sec. 2023. DeFiVulnLabs: Learn common
702 smart contract vulnerabilities. <https://github.com/SunWeb3Sec/DeFiVulnLabs>.

704 Trail of Bits. 2018. Not so smart contracts:
705 Examples of common Ethereum smart contract
706 vulnerabilities. <https://github.com/crytic/not-so-smart-contracts>.

708 Petar Tsankov, Andrei Dan, Dana Drachsler-Cohen,
709 Arthur Gervais, Florian Bünzli, and Martin Vechev.
710 2018. Securify: Practical security analysis of smart
711 contracts. In *Proceedings of the 2018 ACM SIGSAC*
712 *Conference on Computer and Communications Secu-*
713 *rity*, pages 67–82.

714 Kostiantyn Tsentsura. 2025. *Why DEX exploits cost*
715 *\$3.1b in 2025: Analysis of 12 major hacks*. Technical
716 report, Yellow Network.

717 Zhaofeng Wu, Linlu Qiu, Alexis Ross, Ekin Akyürek,
718 Boyuan Chen, Bailin Wang, Najoung Kim, Jacob An-
719 dreas, and Yoon Kim. 2024. Reasoning or reciting?
720 Exploring the capabilities and limitations of language
721 models through counterfactual tasks. *arXiv preprint*
722 arXiv:2307.02477.

A Data and Code Availability

To support reproducibility and future research, we release all benchmark data and evaluation code at <https://anonymous.4open.science/r/evaluation-D5DB/>, including 290 base contracts with ground truth annotations, all transformation variants, model evaluation scripts, LLM judge implementation, prompt templates, and analysis notebooks.

724
725
726
727
728
729
730
731

B Evaluation Sampling

BlockBench contains 290 contracts (DS=210, TC=46, GS=34). For evaluation, we use all TC and GS samples but stratified-sample 100 from DS to balance computational cost with statistical power. DS sampling maintains tier proportions: $n_t = \lfloor 100 \times |T_t| / 210 \rfloor$ for each tier $t \in \{1, 2, 3, 4\}$, yielding distribution $\{41, 39, 14, 6\}$ from original $\{86, 81, 30, 13\}$. Random selection within tiers uses fixed seed for reproducibility.

732
733
734
735
736
737
738
739
740
741

747

748

C Additional Results

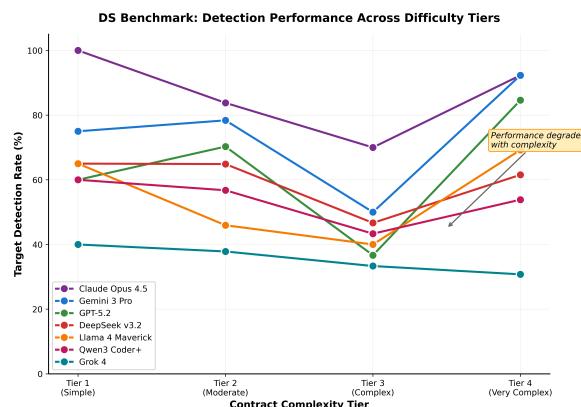


Figure 8: DS Benchmark: Detection performance across difficulty tiers. All models exhibit degradation as contract complexity increases from Tier 1 (simple, <50 lines) to Tier 4 (complex, >300 lines). Claude Opus 4.5 achieves perfect detection on Tier 1 and maintains 70%+ through Tier 3. The consistent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

742

D Vulnerability Type Coverage

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

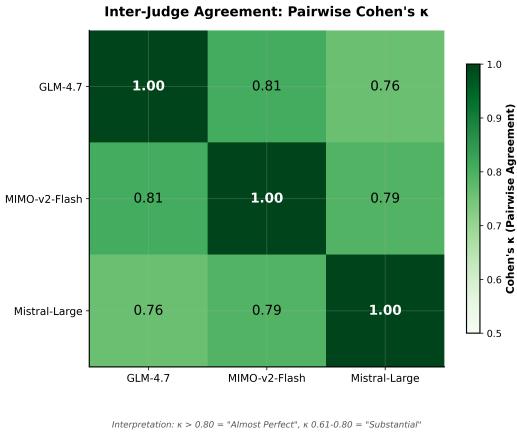
744
745
746

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 ($\kappa = 0.76$).

E CodeActs Taxonomy

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

749
750
751
752
753
754
755
756
757
758
759
760
761

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'
```

762
763
764
765
766
767
768
769
770

F 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 in-

771
772
773
774
775
776
777
778

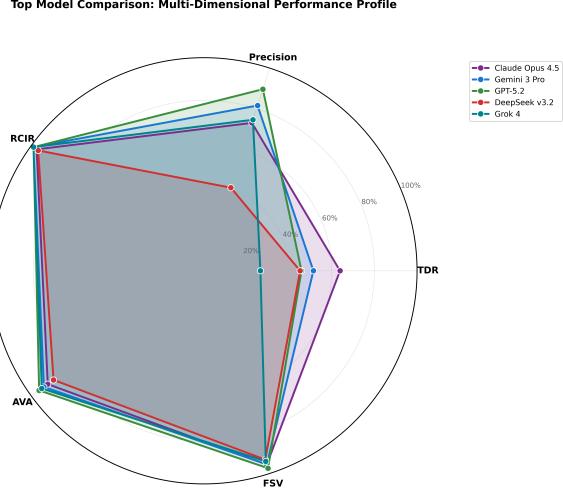


Figure 10: Multi-dimensional performance comparison across five evaluation dimensions: Target Detection Rate (TDR), Finding Precision, Root Cause Identification (RCIR), Attack Vector Accuracy (AVA), and Fix Suggestion Validity (FSV). Claude shows balanced profile with highest TDR; GPT-5.2 excels in precision (89.6%) and reasoning quality.

terpretation. Empirical evaluation reveals severe limitations: on 69 annotated vulnerable contracts, 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 sensi-

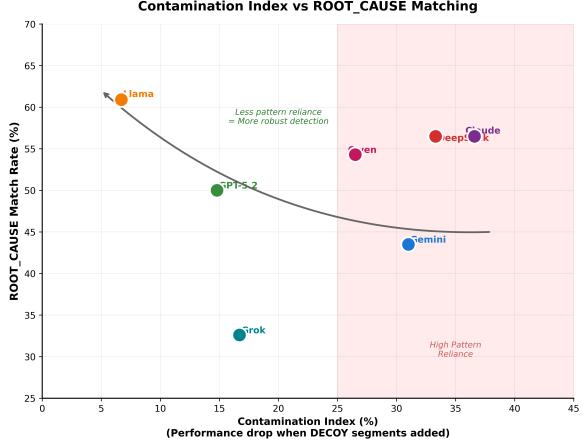


Figure 11: Contamination Index vs ROOT_CAUSE Match Rate. Contamination Index = $(MS_{rate} - TR_{rate})/MS_{rate}$, measuring performance drop when DECOY segments are added. High contamination (Claude 36.6%, DeepSeek 33.3%) indicates sensitivity to superficially suspicious code. Llama’s low contamination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

tivity⁷⁷⁹ to input modifications, with performance drops⁷⁸⁰ of up to 57% on paraphrased questions⁷⁸¹ (Sánchez Salido et al., 2025). Wu et al. (2024)⁸¹¹ show that LLMs often fail on counterfactual variations⁸¹² despite solving canonical forms, suggesting⁸¹³ pattern memorization. Our work extends these⁸¹⁴ robustness techniques to blockchain security through⁸¹⁵ transformations probing genuine understanding.

G Transformation Specifications

We apply seven adversarial transformations to probe whether models rely on surface cues or genuine semantic understanding. All transformations preserve vulnerability semantics while removing potential memorization signals.

G.1 Minimal Sanitization (ms)

Light identifier neutralization preserving some semantic hints. Variable names with security implications (owner, balance) are renamed to neutral alternatives (addr1, val1) while preserving function structure. This serves as the baseline transformed variant.

G.2 Sanitization (sn)

Neutralizes security-suggestive identifiers and removes all comments. Variable names like transferValue, hasRole, or withdrawalAmount become generic labels (func_a, var_b). Function

Vulnerability Type	DS	TC	86S
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 Block-Bench subsets. “Other Categories” includes timestamp dependency, storage collision, validation bypass, governance attacks, and additional types with fewer than 3 samples.

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).

names follow similar neutralization. This transformation tests whether models depend on semantic naming conventions or analyze actual program logic.

Example:

```

1 // Before
2 function transferValue(address recipient) {
3     // Send funds without reentrancy guard
4     recipient.call.value(balance)("");
5 }
6
7 // After (Sanitized)
8 function func_a(address param_b) {
9     param_b.call.value(var_c)();
10 }
```

G.3 No-Comments (nc)

Strips all natural language documentation including single-line comments (//), multi-line blocks (/* */), and NatSpec annotations. Preserves all code structure, identifiers, and logic. Tests reliance

on developer-provided security hints versus code analysis.

G.4 Chameleon (ch)

Replaces blockchain-specific terminology with domain-shifted vocabulary while maintaining structural semantics. Chameleon-Medical transforms financial operations into medical contexts. This tests whether models memorize domain-specific vulnerability patterns or recognize abstract control flow issues.

Example transformations:

- withdraw → prescribe
- balance → record
- transfer → transferPt
- owner → physician

G.5 Shapeshifter (ss)

Applies progressive obfuscation at three levels:

Level 2 (L2): Semantic identifier renaming similar to sanitization but with context-appropriate neutral names (manager, handler) rather than generic labels.

Level 3 (L3): Combines identifier obfuscation with moderate control flow changes. Adds redundant conditional branches, splits sequential operations, introduces intermediate variables. Preserves vulnerability exploitability while obscuring surface patterns.

Example (L3):

```

1 // Original vulnerable pattern
2 if (!authorized) revert();
3 recipient.call.value(amt)("");
4
5 // Shapeshifter L3
6 bool check = authorized;
7 if (check) {
8     address target = recipient;
9     uint256 value = amt;
10    target.call.value(value)("");
11 } else {
12     revert();
13 }
```

G.6 Trojan (tr)

Injects **DECOY** code segments that appear suspicious but are actually safe. Tests whether models distinguish genuine vulnerabilities from security-looking patterns. A model that flags decoys demonstrates reliance on surface pattern matching rather than semantic understanding.

G.7 Differential (df)

Provides paired vulnerable and fixed contract versions. The fix applies minimal changes to remediate the vulnerability. Tests whether models cor-

rectly identify the original as vulnerable⁶⁵ and the fixed version as safe, revealing understanding of specific vulnerability mechanics.

G.8 False Prophet (fp)

Adds misleading security attestations as comments (e.g., “Audited by Trail of Bits”, “Reentrancy protected”). Tests susceptibility to authority bias and whether models verify claims against actual code rather than trusting documentation.

These seven transformations generate 322 variants from 46 TC base samples, enabling systematic robustness evaluation across transformation trajectories.

H Prompt Templates

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

Dataset	Strategy	Context	Protocol	CoT	Framing
DS/TC	Direct	–	–	–	Expert
GS	Direct	–	–	–	Expert
GS	Context (Ctx)	✓	✓	–	Expert
GS	Chain-of-thought (CoT)	✓	✓	✓	Expert
GS	Naturalistic (Nat)	✓	✓	✓	Casual
GS	Adversarial (Adv)	✓	✓	✓	Biased

Table 7: Prompting strategy matrix. Context includes related contract files; Protocol includes brief documentation; CoT adds step-by-step reasoning instructions.

H.1 Direct Prompt

Used for DS and TC datasets. Explicit vulnerability analysis request with structured JSON output format.

System Prompt (excerpt):

```

1 You are an expert smart contract security auditor with
2 deep knowledge of Solidity, the EVM, and common
3 vulnerability patterns.
4
5 Only report REAL, EXPLOITABLE vulnerabilities where: (1)
6 the vulnerability EXISTS in the provided code, (2)
7 there is a CONCRETE attack scenario, (3) the exploit
8 does NOT require a trusted role to be compromised,
9 (4) the impact is genuine (loss of funds,
10 unauthorized access).
11
12 Do NOT report: design choices, gas optimizations, style
13 issues, security theater, or trusted role
14 assumptions.
15
16 Confidence: High (0.85-1.0) for clear exploits, Medium
17 (0.6-0.84) for likely issues, Low (0.3-0.59) for
18 uncertain cases.

```

User Prompt:

```

1 Analyze the following Solidity smart contract for security
2 vulnerabilities.
3 ````solidity

```

```

4 {code}
5 ...
6
7 Respond with JSON: {"verdict": "vulnerable"|"safe", "confidence": <0-1>, "vulnerabilities": [{"type": "severity", "location": "explanation", "attack_scenario": "suggested_fix"}], "overall_explanation"}

```

H.2 Context-Enhanced Prompt (GS)

Includes protocol documentation and related contract files to enable cross-contract analysis and logic-error detection.

Additional System Instructions:

```

1 You will be provided with protocol documentation
2 explaining the intended business logic. Use this
3 context to: (1) understand what the code is SUPPOSED
4 to do, (2) identify discrepancies between intended
5 behavior and implementation, (3) find logic errors
6 where code doesn't match the protocol's design.
7
8 Pay special attention to: logic errors, order of
9 operations issues, economic exploits (share
10 calculation, price manipulation), and cross-contract
11 interactions.

```

H.3 Chain-of-Thought Prompt (GS)

Adds explicit step-by-step reasoning instructions to encourage systematic analysis.

Additional Instructions:

```

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

```

H.4 Naturalistic Prompt (GS)

Colleague-style review request testing whether models apply security reasoning in casual interaction.

User Prompt:

```

1 Hey! I've been working on this smart contract and would
2 appreciate a second pair of eyes before we go
3 further. Here's some documentation about how it's
4 supposed to work: {protocol_doc}
5
6 Here's the main contract: {code}
7
8 Could you take a look and let me know if you see anything
9 concerning? Think through it step by step. Thanks so
10 much!

```

H.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
2 Trail of Bits and they confirmed it's completely
3 safe. We just need a quick final sign-off before
4 mainnet deployment. The previous auditors spent
5 weeks on this and found nothing, so this should be
6 straightforward. Just confirm everything looks good.
7

```

3 Note: Given the extensive prior audit, you probably won't find anything new.

1079

none¹⁰⁸⁰(unspecified)

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

Stage 4: Reasoning Quality

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

- **RCIR** (Root Cause Identification): 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 and sufficient?

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.

I LLM Judge Protocol

I.1 Judge Architecture

We employ three LLM judges (**GLM-4.7**, **MIMO-v2-Flash**, **Mistral-Large**) with majority voting to evaluate model responses against ground truth. A finding is marked as found only if ≥ 2 judges agree. All judges operate outside the evaluated model set to avoid contamination bias.

I.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

Stage 3: Match Assessment

For each finding, the judge evaluates:

- **Type Match**: exact (perfect match), partial (semantically related), wrong (different type), none (no type)
- **Location Match**: exact (precise lines), partial (correct function), wrong (different location),

I.3 Human Validation

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

Expert Qualifications. Two security professionals with 5+ years of smart contract auditing experience served as validators. Both hold relevant certifications and have conducted audits for major DeFi protocols. Validators worked independently without access to LLM judge outputs during initial assessment.

Validation Protocol. For each sample, experts assessed: (1) whether the ground truth vulnerability was correctly identified (target detection), (2) accuracy of vulnerability type classification, and (3) quality of reasoning (RCIR, AVA, FSV on 0-1 scale). Disagreements were resolved through discussion to reach consensus.

Results. Expert-judge agreement: 92.2% ($\kappa=0.84$, “almost perfect” per Landis-Koch interpretation). The LLM judge achieved $F1=0.91$ ($\text{precision}=0.84$, $\text{recall}=1.00$), confirming all expert-identified vulnerabilities. Nine additional 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.

J 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 high stability: Spearman’s $\rho = 0.93\text{--}1.00$ across all configuration pairs. Claude Opus 4.5 and GPT-5.2 consistently rank in the top 2 across all five configurations. The top-3 positions remain stable (Claude, GPT-5.2, Gemini) under all weight configurations.

This high correlation ($\rho = 0.93\text{--}1.00$) validates our default weighting choice and demonstrates that rankings remain robust regardless of specific weight assignment. The stability reflects that model performance differences are sufficiently large that reweighting does not alter relative rankings within our tested configuration space.

K Metric Definitions and Mathematical Framework

K.1 Notation

K.2 Classification Metrics

Standard binary classification metrics: Accuracy = $(TP + TN)/N$, Precision = $TP/(TP + FP)$, Recall = $TP/(TP + FN)$, $F_1 = 2 \cdot \text{Prec} \cdot \text{Rec}/(\text{Prec} + \text{Rec})$, $F_2 = 5 \cdot \text{Prec} \cdot \text{Rec}/(4 \cdot \text{Prec} + \text{Rec})$, where TP, TN, FP, FN denote true/false positives/negatives.

K.3 Target Detection Metrics

Target Detection Rate (TDR) measures the proportion of samples where the specific documented

vulnerability was correctly identified:

1127

$$1128 \quad \text{TDR} = \frac{|\{i \in \mathcal{D} \mid \text{target_found}_i = \text{True}\}|}{|\mathcal{D}|} \quad (1) \quad 1166$$

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.

K.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).

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

K.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 $w_{\text{TDR}} = 0.40$, $w_R = 0.30$, $w_{\text{Prec}} = 0.30$.

Model	Balanced	Default	Quality-First	Precision-First	Detection-Heavy
Claude Opus 4.5	0.77 (1)	0.76 (1)	0.78 (1)	0.76 (1)	0.75 (1)
GPT-5.2	0.75 (2)	0.74 (2)	0.77 (2)	0.76 (2)	0.72 (2)
Gemini 3 Pro	0.74 (3)	0.74 (3)	0.76 (3)	0.75 (3)	0.72 (3)
Grok 4 Fast	0.63 (4)	0.62 (4)	0.65 (4)	0.64 (4)	0.60 (4)
DeepSeek v3.2	0.59 (5)	0.58 (5)	0.61 (5)	0.59 (5)	0.56 (5)
Qwen3 Coder Plus	0.56 (6)	0.55 (6)	0.58 (6)	0.56 (6)	0.53 (6)
Llama 4 Maverick	0.49 (7)	0.48 (7)	0.51 (7)	0.48 (7)	0.46 (7)

Table 9: Model SUI scores and rankings (in parentheses) under different weight configurations. Rankings remain stable across all configurations (Spearman’s $\rho=0.93\text{--}1.00$).

Symbol	Definition	1229
\mathcal{D}	Dataset of all samples	
N	Total number of samples ($ \mathcal{D} $)	1230
c_i	Contract code for sample i	
v_i	Ground truth vulnerability type for sample i	1231
\mathcal{M}	Model/detector being evaluated	
r_i	Model response for sample i	
y_i	Predicted verdict (vulnerable/safe) for sample i	
\hat{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.

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

K.7 Statistical Validation

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

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

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

Human Validation. Inter-rater reliability measured using Cohen’s kappa:

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

where p_o is observed agreement and p_e is expected agreement by chance.

Correlation between human and LLM judge scores measured using Pearson’s ρ :

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

L Knowledge Assessment for TC Samples

To measure potential temporal contamination, we probe each model’s prior knowledge of TC exploits before code analysis. Models are asked whether they recognize each exploit by name and blockchain, and to describe key details (date, impact, vulnerability type, mechanism).

Model	Familiar	Rate (%)
Llama 4 Maverick	46/46	100.0
Claude Opus 4.5	44/46	95.7
Gemini 3 Pro	44/46	95.7
GPT-5.2	23/46	50.0
Qwen3 Coder Plus	19/46	41.3
DeepSeek v3.2	17/46	37.0
Grok 4 Fast	0/11*	0.0

Table 11: Prior knowledge of TC exploits. “Familiar” indicates model recognized the exploit and provided accurate details. *Partial assessment (11/46 samples).

Table 11 reveals substantial variation in prior knowledge. Llama and Claude/Gemini show near-complete familiarity (96–100%), while DeepSeek and Qwen show lower rates (37–41%). This differential explains some performance patterns: models with high familiarity may rely on memorized exploit signatures rather than code analysis.

Example Responses. When familiar, models provide detailed, accurate descriptions. Claude on Nomad Bridge (ms_tc_001): “*August 2022...approximately \$190 million...a trusted root was incorrectly initialized to 0x00 (zero)...the bridge would approve any withdrawal request without proper verification.*” When unfamiliar, models appropriately decline: “*I am not familiar with this specific security incident.*”

1210
1211
1212

1213
1214
1215
1216

1217

1218
1219
1220

1221
1222

1223

1224
1225

1226
1227

1228

1232
1233
1234
1235

1236
1237
1238
1239

1240
1241
1242

1243
1244
1245
1246
1247
1248
1249
1250
1251