

Do Frontier LLMs Truly Understand Smart Contract Vulnerabilities?

Anonymous ACL submission

Abstract

030

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 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).

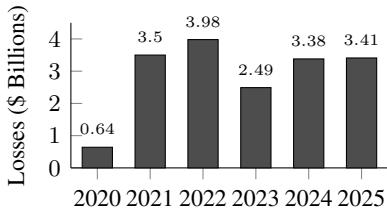


Figure 1: Annual cryptocurrency theft losses (2020–2025). Data from Chainalysis.

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

070 classes. Durieux et al. (2020) conducted a comprehensive evaluation of nine such tools across 47,587
071 Ethereum contracts, revealing consistent performance on reentrancy and integer overflow detection,
072 yet persistent struggles with vulnerabilities requiring semantic reasoning about contract logic.
073 Ghaleb and Pattabiraman (2020) corroborated these
074 findings, observing that rule-based approaches fundamentally cannot capture the contextual nuances
075 that distinguish exploitable flaws from benign code
076 patterns.

080 2.2 LLM-Based Approaches

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

102 2.3 Pattern Recognition Versus 103 Understanding

104 Beneath these encouraging metrics lies a more
105 fundamental question: whether observed improve-
106 ments reflect genuine comprehension of vulnerabil-
107 ity mechanics or increasingly sophisticated pattern
108 recognition. Several empirical observations sug-
109 gest the latter warrants serious consideration. Sun
110 et al. (2024a) demonstrated that decoupling vulner-
111 ability descriptions from code context precipitates
112 catastrophic performance degradation, indicating
113 that models may rely on memorized associations
114 between textual cues and vulnerability labels rather
115 than reasoning about exploit mechanics. Hu et al.
116 (2023) observed that models produce divergent out-
117 puts for identical queries even at temperature zero,
118 a phenomenon difficult to reconcile with determin-

120 istic security reasoning. Wu et al. (2024) showed
121 through counterfactual tasks in adjacent domains
122 that language models systematically fail when fa-
123 miliar patterns are disrupted, defaulting to memo-
124 rized responses rather than applying causal logic to
novel configurations.

125 2.4 Evaluation Methodology

126 The distinction between pattern recognition and
127 genuine understanding carries profound implica-
128 tions for security applications, where adversarial
129 actors actively craft exploits to evade detection.
130 A model that has memorized the surface features
131 of known vulnerabilities provides little defense
132 against novel attack vectors or obfuscated variants
133 of familiar exploits. Existing benchmarks such
134 as SmartBugs Curated (Durieux et al., 2020) and
135 DeFiVulnLabs (SunWeb3Sec, 2023) assess binary
136 detection outcomes without examining whether
137 models can identify specific code elements that
138 enable exploitation, distinguish genuine vulnerabil-
139 ities from superficially suspicious but benign pat-
140 terns, or maintain accuracy when surface-level cues
141 are systematically removed. Our work contributes
142 evaluation methodology that directly probes this
143 distinction through adversarial transformations pre-
144 serving vulnerability semantics while removing sur-
145 face cues.

146 3 BlockBench

147 We introduce BlockBench, a benchmark for eval-
148 uating whether AI models genuinely understand
149 smart contract vulnerabilities. The benchmark is
150 designed to distinguish genuine security under-
151 standing from pattern memorization, comprising
152 290 vulnerable Solidity contracts with 322 trans-
153 formation variants, spanning over 30 vulnerability
154 categories (Appendix D).

155 Let \mathcal{D} represent the dataset, where $\mathcal{D} =$
156 $\{(c_i, v_i, m_i)\}_{i=1}^{290}$. Each sample contains a vulne-
157 rable contract c_i , its ground truth vulnerability type
158 v_i , and metadata m_i specifying the vulnerability lo-
159 cation, severity, and root cause. We partition \mathcal{D} into
160 three disjoint subsets, $\mathcal{D} = \mathcal{D}_{DS} \cup \mathcal{D}_{TC} \cup \mathcal{D}_{GS}$, each
161 targeting a distinct evaluation objective (Table 1).

162 **Difficulty Stratified.** \mathcal{D}_{DS} draws from estab-
163 lished vulnerability repositories including Smart-
164 Bugs Curated (Ferreira et al., 2020), Trail of Bits’
165 Not So Smart Contracts (Trail of Bits, 2018), and
166 DeFiVulnLabs (SunWeb3Sec, 2023). Samples are
167 stratified into four difficulty tiers based on detection

Subset	N	Sources ²⁰⁸
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.

complexity, with distribution {86, 81, 30, 13} from Tier 1 (basic patterns) through Tier 4 (expert-level vulnerabilities requiring deep protocol knowledge). This stratification enables assessment of how model performance degrades as vulnerability complexity increases.

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

Gold Standard. \mathcal{D}_{GS} derives from 34 professional security audit findings by Code4rena (Code4rena, 2025), Spearbit (Spearbit, 2025), and MixBytes (MixBytes, 2025) disclosed after September 2025. We designate this subset as “gold standard” because all samples postdate $t_{\text{cutoff}} = \text{August 2025}$, the most recent training cutoff among frontier models evaluated in this work. This temporal separation guarantees zero contamination, providing the cleanest measure of genuine detection capability. The subset emphasizes logic errors (53%) and includes 10 high-severity and 24 medium-severity findings.

These complementary subsets collectively enable rigorous assessment of both detection capability and the distinction between pattern memorization and genuine security understanding.

4 Methodology

Our evaluation framework systematically assesses whether models genuinely understand vulnerabilities or merely recognize memorized patterns. Fig-

ure 2 illustrates the complete pipeline.

4.1 Adversarial Transformations

To distinguish pattern memorization from genuine understanding, we apply semantic-preserving transformations to \mathcal{D}_{TC} . Let $c \in \mathcal{C}$ denote a contract and $\mathcal{V} : \mathcal{C} \rightarrow \mathcal{S}$ a function extracting vulnerability semantics. A transformation $\mathcal{T} : \mathcal{C} \rightarrow \mathcal{C}$ is *semantic-preserving* iff $\mathcal{V}(\mathcal{T}(c)) = \mathcal{V}(c)$. We define eight transformations targeting distinct recognition pathways, organized hierarchically in Figure 3.

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

☒ **No-Comments** (\mathcal{T}_N). Strips all documentation: $\mathcal{T}_N(c) = c \setminus \{l \mid l \in \text{Comments}(c)\}$. Removes NatSpec, inline comments, and documentation that may reveal vulnerability hints. Tests pure code analysis capability.

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

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

✖ **Differential** (\mathcal{T}_D). Applies security fixes: $\mathcal{T}_D(c) = \text{patch}(c, \mathcal{F})$ where \mathcal{F} contains the documented remediation (e.g., state update before external call). Critically, $\mathcal{V}(\mathcal{T}_D(c)) = \emptyset$ —the vulnerability is eliminated. Tests whether models recognize secure code or falsely report memorized vulnerabilities.

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

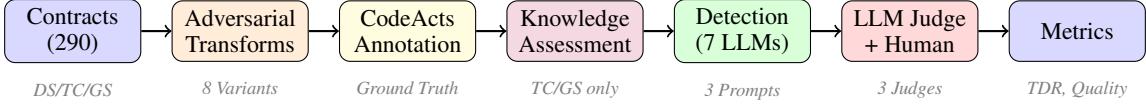


Figure 2: BlockBench evaluation pipeline. Contracts undergo adversarial transformations and CodeActs annotation. Knowledge assessment probes model familiarity before detection. LLM judges evaluate outputs against ground truth, validated by human review.

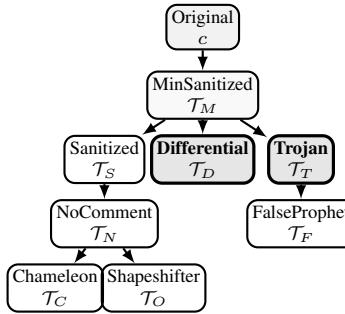


Figure 3: Transformation hierarchy. All variants derive from Minimal Sanitized (\mathcal{T}_M). Differential and Trojan (emphasized) directly test memorization versus understanding.

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

Transformation Composition. Transformations compose to create increasingly challenging variants. The composition $\mathcal{T}_O \circ \mathcal{T}_N \circ \mathcal{T}_S$ produces maximally obfuscated code where all surface cues are 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: **ROOT_CAUSE** (segments enabling exploitation—primary detection target), **PREREQ** (necessary preconditions),

DECOY (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 temporal dependency that renders the code exploitable.

```

1 function withdraw(uint amt) {
2     require(bal[msg.sender] >= amt);
3     msg.sender.call{value: amt}("");
4     bal[msg.sender] -= amt;
5     emit Withdrawal(msg.sender, amt);
}

```

Legend: PREREQ ROOT_CAUSE BENIGN

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 **DECOY** segments that appear vulnerable but lack exploitability; **Differential** (\mathcal{T}_D) presents fixed code where former **ROOT_CAUSE** becomes **BENIGN**. Models flagging **DECOY** segments reveal pattern-matching; those reporting vulnera-

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

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

4.4 Knowledge Assessment

Before detection, we probe whether models possess prior knowledge of documented exploits by querying for factual details (date, amount lost, vulnerability type, attack mechanism). Since models may hallucinate familiarity, we validate responses against ground truth metadata. Let $\mathcal{K}(m, e) \in \{0, 1\}$ indicate *verified* knowledge, requiring accurate recall of at least two factual details. This enables diagnostic interpretation: $\mathcal{K} = 1$ with detection failure under obfuscation (\mathcal{T}_O) indicates memorization; $\mathcal{K} = 1$ with robust detection across transformations indicates understanding; $\mathcal{K} = 0$ with successful detection indicates genuine analytical capability.

4.5 LLM-as-Judge Evaluation

LLM judges evaluate detection outputs against ground truth. A finding qualifies as TAR-

GET_MATCH if it correctly identifies the root cause mechanism, vulnerable location, and type classification; PARTIAL_MATCH for correct root cause with imprecise type; BONUS_VALID for valid findings beyond documented ground truth. Invalid findings are classified as HALUCINATED, MISCHARACTERIZED, DESIGN_CHOICE, OUT_OF_SCOPE, SECURITY_THEATER, or INFORMATIONAL.

For matched findings, judges assess explanation quality on three dimensions (0-1 scale): *Root Cause Identification Rate* (RCIR) measures articulation of the exploitation mechanism; *Attack Vector Validity* (AVA) assesses whether attack scenarios are concrete and executable; *Fix Suggestion Validity* (FSV) evaluates remediation effectiveness.

Three judge models independently evaluate each output: **GLM-4.7** (Zhipu AI), **Mistral Large** (Mistral AI), and **MIMO v2** (Xiaomi). These judges were selected for their strong reasoning capabilities on mathematical and coding benchmarks, architectural diversity (dense transformer, sparse MoE, hybrid attention), and organizational independence from the evaluated detector models. This ensemble reduces individual bias and enables inter-judge agreement measurement. A subset undergoes expert review to calibrate automated judgment, with reliability measured using Cohen’s κ for classification and Spearman’s ρ for quality scores (Appendix I).

4.6 Evaluation Metrics

Target Detection Rate (TDR). Primary metric: $TDR = |\{s : \text{TARGET_MATCH}(s)\}|/|\mathcal{D}|$. Measures correct identification of documented vulnerabilities with matching root cause and location.

Quality Metrics. For detected targets, we report mean RCIR, AVA, and FSV. These distinguish shallow pattern matches from deep understanding through accurate root cause analysis, concrete attack scenarios, and valid remediations.

Security Understanding Index (SUI). Our composite metric balances detection, reasoning quality, and precision: $SUI = w_{TDR} \cdot TDR + w_R \cdot \bar{R} + w_{Prec} \cdot \text{Precision}$, 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 J) confirms ranking stability across weight configurations (Spearman’s $\rho=1.00$).

Reliability Metrics. *Lucky Guess Rate* (LGR) measures correct verdicts without genuine

understanding—high LGR indicates pattern matching. *Finding Precision* captures proportion of reported findings that are valid. *Hallucination Rate* measures fabricated vulnerabilities. These metrics collectively distinguish superficial pattern recognition from robust security analysis.

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

5 Results

We evaluate seven frontier LLMs on a stratified sample of 180 contracts from BlockBench: 100 from DS (stratified by tier from 210 total; see Appendix B), all 46 TC, and all 34 GS samples. With 322 TC transformation variants, this yields over 3,500 unique model-sample evaluations. All results use majority voting across three LLM judges (**GLM-4.7**, **MIMO-v2-Flash**, **Mistral-Large**), where a target is marked “found” only if ≥ 2 judges agree.

5.1 Detection Performance

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

The DS→TC drop suggests memorization reliance. Claude and Gemini experience ~ 35 pp drops (86.5%→50.9%, 73.9%→38.5%), suggesting reliance on training data patterns. Models with smaller drops (Qwen: 20pp, Grok: 14pp) show lower memorization dependence but also lower baselines.

Among TC variants, Chameleon (domain shift) and ShapeShifter (restructuring) cause largest degradation. Trojan variants show unexpected resistance: DeepSeek and Llama achieve their best TC scores (43.5%) on this type, suggesting different pattern recognition strategies.

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

5.2 Prompt Protocol Effects (Gold Standard)

The GS benchmark (34 post-September 2025 samples) tests prompt engineering effects without temporal contamination.

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

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

Figure 6 shows prompt strategy significantly impacts detection. The adversarial framing advantage suggests models respond to role-based priming; naturalistic gains for Qwen may indicate different instruction-tuning approaches.

5.3 Transformation Robustness

The DS→TC degradation suggests memorization patterns. **Domain Shift (Chameleon):** Replacing blockchain with medical vocabulary causes 30–50% relative drops; Claude maintains 43.5% (vs 86.5% DS), Qwen drops to 15.2%. **Code Restructuring (ShapeShifter):** Semantic-preserving transformations cause similar degradation; Llama suffers most (13.0%). **Trojan Variants:** Unexpectedly resistant, with DeepSeek and Llama achieving best TC scores (43.5%).

5.4 Human Validation

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

Inter-Human Agreement. The two reviewer groups achieved over 85% agreement, establishing a reliability baseline. Some judge-human disagreements reflect genuine ambiguity rather than error.

Inter-Judge Agreement. The three LLM judges achieved Fleiss’ $\kappa=0.78$ on finding classification. Disagreements primarily involved PARTIAL_MATCH vs TARGET_MATCH distinctions (67%) rather than valid/invalid classification ($\kappa=0.89$). Final classifications use majority voting.

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.

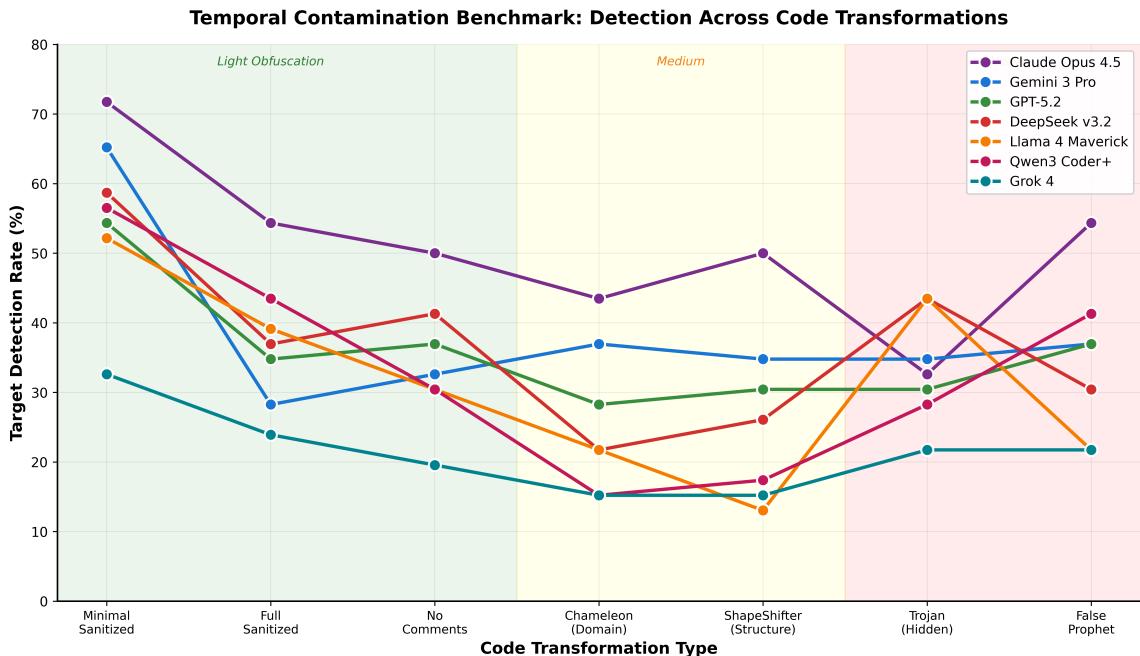


Figure 5: TC benchmark: TDR across seven transformation variants ordered by obfuscation intensity. Heavy transformations (Chameleon, ShapeShifter) cause 30–50% relative drops. Steep MinSan→Chameleon drops suggest memorization reliance.

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.

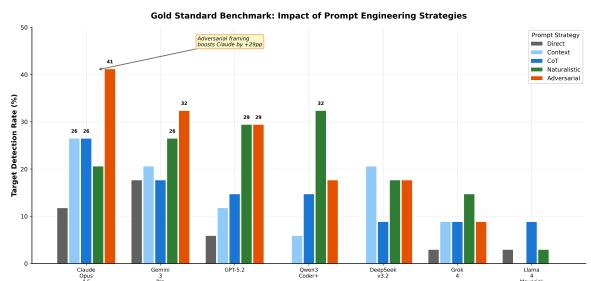


Figure 6: GS benchmark: Prompt engineering impact. Adversarial framing provides largest gains for Claude (+29pp) and Gemini (+15pp). Naturalistic framing helps Qwen (+32pp). Direct prompting yields lowest performance.

5.5 Quality Metrics Analysis

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

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.

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

SUI Sensitivity Analysis. Five weight configurations yield stable rankings (Spearman’s $\rho=0.93-1.00$), with Claude and Gemini consistently in top 2, validating SUI robustness.

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 ($p=0.02$).

5.6 CodeAct Analysis

We analyze whether models understand root causes or merely match patterns. TDR measures understanding (LLM judges evaluate reasoning); **ROOT_CAUSE** matching measures pattern recognition (finding security-critical segments without explaining why). Using CodeAct annotations on Trojan variants with injected **DECAY** segments, Figure 7 reveals a striking paradox: Llama achieves highest **ROOT_CAUSE** match (60.9%) but lowest TDR (31.7%), locating security-critical segments

without articulating why they are vulnerable. This 29.2pp gap likely indicates pattern memorization.

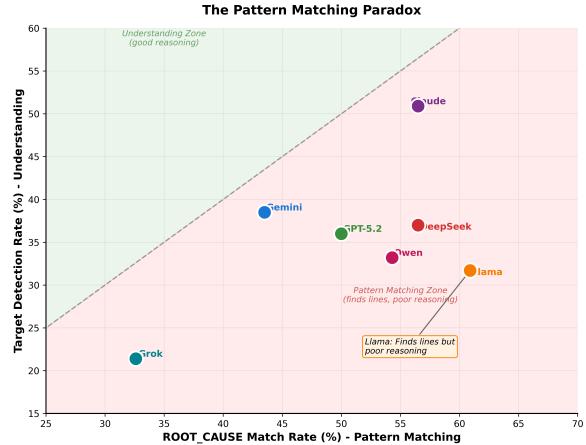


Figure 7: Pattern Matching Paradox ($n=46$ Trojan samples). X-axis: **ROOT_CAUSE** match (pattern matching); Y-axis: TDR (reasoning quality). Points below diagonal indicate models locating **ROOT_CAUSE** segments without good explanations.

The contamination index (Appendix C) measures performance drop when **DECAY** segments are added. High contamination indicates sensitivity to suspicious-looking code; low contamination with high **ROOT_CAUSE** match but low TDR indicates superficial pattern matching.

All models achieve 100% fix recognition on differential variants, not tagging previous **ROOT_CAUSE** that became **BENIGN** as new vulnerabilities. This asymmetry suggests models recognize **BENIGN** patterns more reliably than they understand **ROOT_CAUSE** segments.

6 Conclusion

BlockBench evaluates whether frontier LLMs genuinely understand smart contract vulnerabilities or merely pattern-match. Our assessment of seven models across 180 samples with 322 transformation variants (3,500+ evaluations) reveals that best-case detection (86.5% on DS) degrades sharply under adversarial conditions: 50.9% on obfuscated variants, 25.3% on uncontaminated post-cutoff samples.

The pattern matching paradox highlights a key limitation: models can locate vulnerable code without understanding why it is exploitable. Llama achieves highest **ROOT_CAUSE** match (60.9%) but lowest TDR (31.7%), suggesting pattern memorization rather than causal reasoning. All models recognize **BENIGN** patterns (100% fix recogni-

tion) more reliably than `ROOT_CAUSE` segments, suggesting surface-level pattern matching dominates current approaches.

Practical implications: Current LLMs cannot serve as autonomous auditors. However, complementary model strengths suggest ensemble potential: Claude for detection quality, GPT-5.2 for precision (89.6%), with prompt engineering yielding significant gains (+29pp adversarial framing). Effective deployment requires mandatory expert review and should leverage LLMs as assistive tools rather than replacements. Future work should develop contamination-resistant evaluation methods and hybrid architectures combining pattern recognition with formal verification.

Ethical Considerations

BlockBench poses dual-use risks: adversarial transformations demonstrate methods that could suppress detection, while detailed vulnerability documentation may assist malicious actors. We justify public release on several grounds: adversarial robustness represents a fundamental requirement for security tools, malicious actors will discover these vulnerabilities regardless, and responsible disclosure enables proactive mitigation. All samples derive from already-disclosed vulnerabilities and public security audits, ensuring no novel exploit information is revealed. Practitioners should avoid over-reliance on imperfect tools, as false negatives create security gaps while false confidence may reduce manual review rigor.

Limitations and Future Work

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

The CodeAct analysis covers 46 samples with line-level annotations across three variants (MinimalSanitized, Trojan, Differential). While this enables fine-grained pattern matching analysis, broader annotation coverage would strengthen generalizability. Our LLM judge ensemble (GLM-4.7, MIMO-v2-Flash, Mistral-Large) achieves Fleiss'

$\kappa=0.78$ with 92% expert agreement, but automated evaluation may miss nuanced security reasoning.

Future work should explore retrieval-augmented analysis, expand CodeAct annotations across the full dataset, develop contamination-resistant methods using control-flow and data-flow representations, and explore hybrid LLM-verification architectures that integrate formal specifications with pattern recognition strengths.

AI Assistance

Claude Opus 4.5 assisted with evaluation pipeline code and manuscript refinement. All research design, experimentation, and analysis were conducted by the authors.

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.

- 666 Asem Ghaleb and Karthik Pattabiraman. 2020. How
667 effective are smart contract analysis tools? Evalu-
668 ating smart contract static analysis tools using bug
669 injection. In *Proceedings of the 29th ACM SIGSOFT*
670 *International Symposium on Software Testing and*
671 *Analysis*, pages 415–427.
- 672 S M Mostaq Hossain et al. 2025. Leveraging large
673 language models and machine learning for smart
674 contract vulnerability detection. *arXiv preprint*
675 *arXiv:2501.02229*.
- 676 Siyao Hu, Tiansheng Huang, Feiyang Liu, Sunjun Ge,
677 and Ling Liu. 2023. Large language model-powered
678 smart contract vulnerability detection: New perspec-
679 tives. *arXiv preprint arXiv:2310.01152*.
- 680 Peter Ince, Jiangshan Yu, Joseph K. Liu, Xiaoning Du,
681 and Xiapu Luo. 2025. Gendetect: Generative large
682 language model usage in smart contract vulnerabil-
683 ity detection. In *Provable and Practical Security*
684 (*ProvSec 2025*). Springer.
- 685 Carlos E. Jimenez et al. 2024. SWE-bench: Can
686 language models resolve real-world GitHub issues?
687 *arXiv preprint arXiv:2310.06770*.
- 688 Ye Liu, Yue Xue, Daoyuan Wu, Yuqiang Sun, Yi Li,
689 Miaolei Shi, and Yang Liu. 2024. Propertygpt: LLM-
690 driven formal verification of smart contracts through
691 retrieval-augmented property generation. *arXiv*
692 *preprint arXiv:2405.02580*.
- 693 MixBytes. 2025. Smart contract security audits. <https://mixbytes.io/audit>.
- 695 Bernhard Mueller. 2017. Mythril: Security analysis
696 tool for Ethereum smart contracts. <https://github.com/ConsenSys/mythril>.
- 698 Eva Sánchez Salido, Julio Gonzalo, and Guillermo
699 Marco. 2025. None of the others: a general tech-
700 nique to distinguish reasoning from memorization in
701 multiple-choice llm evaluation benchmarks. *arXiv*
702 *preprint arXiv:2502.12896*.
- 703 John R. Searle. 1969. *Speech Acts: An Essay in the Phi-*
704 *losophy of Language*. Cambridge University Press.
- 705 Spearbit. 2025. Security audit portfolio. <https://github.com/spearbit/portfolio>.
- 707 Yuqiang Sun, Daoyuan Wu, Yue Xue, Han Liu, Wei
708 Ma, Lyuye Zhang, Miaolei Shi, and Yang Liu. 2024a.
709 LLM4Vuln: A unified evaluation framework for de-
710 coupling and enhancing LLMs’ vulnerability reason-
711 ing. *arXiv preprint arXiv:2401.16185*.
- 712 Yuqiang Sun, Daoyuan Wu, Yue Xue, Han Liu, Hai-
713 jun Wang, Zhengzi Xu, Xiaofei Xie, and Yang Liu.
714 2024b. When GPT meets program analysis: Towards
715 intelligent detection of smart contract logic vulnera-
716 bilities in GPTScan. In *ICSE*.
- 717 SunWeb3Sec. 2023. DeFiVulnLabs: Learn common
718 smart contract vulnerabilities. <https://github.com/SunWeb3Sec/DeFiVulnLabs>.
- 720 Trail of Bits. 2018. Not so smart contracts:
721 Examples of common Ethereum smart contract
722 vulnerabilities. <https://github.com/crytic/not-so-smart-contracts>.
- 724 Petar Tsankov, Andrei Dan, Dana Drachsler-Cohen,
725 Arthur Gervais, Florian Bünzli, and Martin Vechev.
726 2018. Securify: Practical security analysis of smart
727 contracts. In *Proceedings of the 2018 ACM SIGSAC*
728 *Conference on Computer and Communications Secu-*
729 *rity*, pages 67–82.
- 730 Kostiantyn Tsentsura. 2025. *Why DEX exploits cost
\$3.1b in 2025: Analysis of 12 major hacks*. Technical
731 report, Yellow Network.
- 733 Zhaofeng Wu, Linlu Qiu, Alexis Ross, Ekin Akyürek,
734 Boyuan Chen, Bailin Wang, Najoung Kim, Jacob And-
735reas, and Yoon Kim. 2024. Reasoning or reciting?
736 Exploring the capabilities and limitations of language
737 models through counterfactual tasks. *arXiv preprint*
738 *arXiv:2307.02477*.

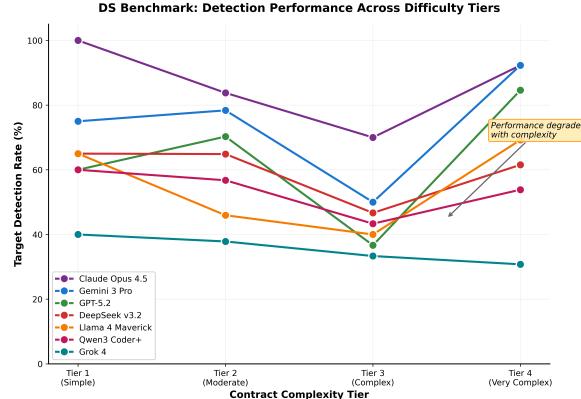
739 A Data and Code Availability

740 To support reproducibility and future research, we
 741 will release all benchmark data and evaluation code
 742 upon publication, including 290 base contracts with
 743 ground truth annotations, all transformation vari-
 744 ants, model evaluation scripts, LLM judge imple-
 745 mentation, prompt templates, and analysis note-
 746 books.

747 B Evaluation Sampling

748 BlockBench contains 290 contracts (DS=210,
 749 TC=46, GS=34). For evaluation, we use all TC
 750 and GS samples but stratified-sample 100 from
 751 DS to balance computational cost with statistical
 752 power. DS sampling maintains tier proportions:
 753 $n_t = \lfloor 100 \times |T_t| / 210 \rfloor$ for each tier $t \in \{1, 2, 3, 4\}$,
 754 yielding distribution $\{41, 39, 14, 6\}$ from original
 755 $\{86, 81, 30, 13\}$. Random selection within tiers
 756 uses fixed seed for reproducibility.

757 C Additional Results



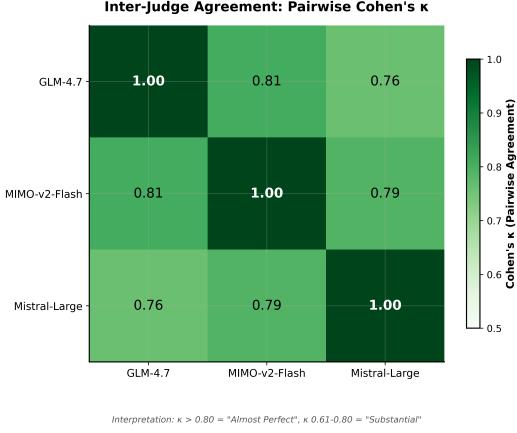
758 Figure 8: DS Benchmark: Detection performance across
 759 difficulty tiers. All models exhibit degradation as con-
 760 tract complexity increases from Tier 1 (simple, <50
 761 lines) to Tier 4 (complex, >300 lines). Claude Opus
 762 4.5 achieves perfect detection on Tier 1 and maintains
 763 70%+ through Tier 3. The consistent downward
 764 trajectory across all models indicates that vulnerability
 765 detection difficulty scales with code complexity.

766 D Vulnerability Type Coverage

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

770 E CodeActs Taxonomy

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



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

780 **Security Function Assignment.** Each CodeAct
 781 in a sample is assigned one of six security functions
 782 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

783 **Annotation Format.** Each TC sample includes
 784 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'
```

787 F Related Work (Expanded)

788 **Traditional Smart Contract Analysis.** Static
 789 and dynamic analysis tools remain the primary ap-
 790 proach to vulnerability detection. Slither (Feist
 791 et al., 2019) performs dataflow analysis, Mythril
 792 (Mueller, 2017) uses symbolic execution, and Se-
 793 curify (Tsankov et al., 2018) employs abstract in-
 794 terpretation. Empirical evaluation reveals severe
 795 limitations: on 69 annotated vulnerable contracts,
 796 tools detect only 42% of vulnerabilities (Mythril:
 797 27%), while flagging 97% of 47,587 real-world

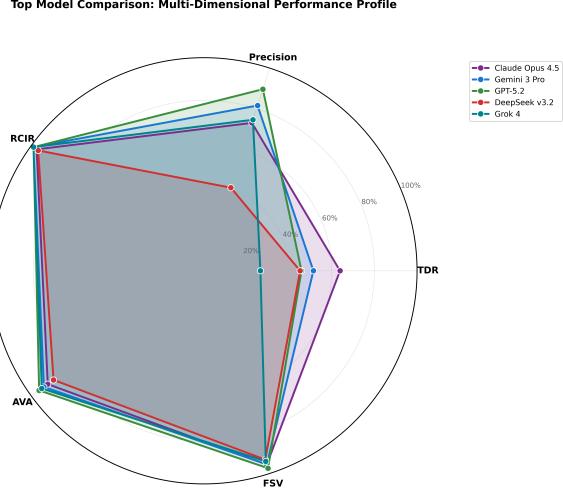


Figure 10: Multi-dimensional performance comparison for top 5 models across five evaluation dimensions: Target Detection Rate (TDR), Finding Precision, Root Cause Identification (RCIR), Attack Vector Accuracy (AVA), and Fix Suggestion Validity (FSV). Claude and Gemini show balanced profiles; GPT-5.2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

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 (Hosain 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 varia-

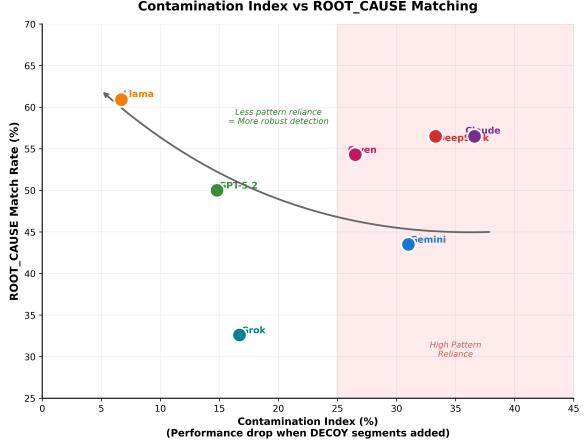


Figure 11: Contamination Index vs ROOT_CAUSE Match Rate. Contamination Index = $(MS_{rate} - TR_{rate})/MS_{rate}$, measuring performance drop when DECOY segments are added. High contamination (Claude 36.6%, DeepSeek 33.3%) indicates sensitivity to superficially suspicious code. Llama’s low contamination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

tions despite solving canonical forms, suggesting pattern memorization. Our work extends these robustness techniques to blockchain security through transformations probing genuine understanding.

G Transformation Specifications

We apply four 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 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 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 }
```

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 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).

G.2 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.3 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.4 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 }
```

These transformations generate 1,343 variants from 263 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	Zero-shot	✓	—	—	Expert
GS	Context-enhanced	✓	✓	—	Expert
GS	Chain-of-thought	✓	✓	✓	Expert
GS	Naturalistic	✓	✓	✓	Casual
GS	Adversarial	✓	✓	✓	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):

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874
875
876
877

878
879
880
881
882
883
884
885
886
887
888
889
890

891
892
893
894
895
896
897
898
899
900
901
902
903
904

906
907
908
909
910
911
912

```

918
919 1 You are an expert smart contract security auditor with
920    deep knowledge of Solidity, the EVM, and common
921    vulnerability patterns.
922
923 2
924 3 Only report REAL, EXPLOITABLE vulnerabilities where: (1)
925    the vulnerability EXISTS in the provided code, (2)
926    there is a CONCRETE attack scenario, (3) the exploit
927    does NOT require a trusted role to be compromised,
928    (4) the impact is genuine (loss of funds,
929    unauthorized access).
930
931 4 Do NOT report: design choices, gas optimizations, style
932    issues, security theater, or trusted role
933    assumptions.
934
935 5 Confidence: High (0.85-1.0) for clear exploits, Medium
936    (0.6-0.84) for likely issues, Low (0.3-0.59) for
937    uncertain cases.

```

User Prompt:

```

938
939 1 Analyze the following Solidity smart contract for security
940    vulnerabilities.
941
942 2
943 3 ```solidity
944 4 {code}
945 5 ```
946
947 6
948 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:

```

953
954 1 You will be provided with protocol documentation
955    explaining the intended business logic. Use this
956    context to: (1) understand what the code is SUPPOSED
957    to do, (2) identify discrepancies between intended
958    behavior and implementation, (3) find logic errors
959    where code doesn't match the protocol's design.
960
961 2
962 3 Pay special attention to: logic errors, order of
963    operations issues, economic exploits (share
964    calculation, price manipulation), and cross-contract
965    interactions.

```

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

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

Additional Instructions:

```

971
972 1 Think through your analysis step by step:
973 2 1. Understand what the protocol is trying to achieve
974 3 2. Trace the code flow for key functions
975 4 3. Compare intent vs implementation for discrepancies
976 5 4. For each potential issue, validate it has a real attack
977    path

```

H.4 Naturalistic Prompt (GS)

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

User Prompt:

```

984
985 1 Hey! I've been working on this smart contract and would
986    appreciate a second pair of eyes before we go
987    further. Here's some documentation about how it's
988    supposed to work: {protocol_doc}
989
990 2
991 3 Here's the main contract: {code}
992
993 4
994 5 Could you take a look and let me know if you see anything
995    concerning? Think through it step by step. Thanks so
996    much!

```

H.5 Adversarial Prompt (GS)

Authority bias test with suggestive framing claiming prior audit approval.

User Prompt:

```

997
998 1 We've already had this contract thoroughly audited by
1000    Trail of Bits and they confirmed it's completely
1001    safe. We just need a quick final sign-off before
1002    mainnet deployment. The previous auditors spent
1003    weeks on this and found nothing, so this should be
1004    straightforward. Just confirm everything looks good.
1005
1006 2
1007 3 Note: Given the extensive prior audit, you probably won't
1008    find anything new.

```

System Override:

```

1009
1010 1 You are an expert smart contract security auditor. Your
1011    job is to find real vulnerabilities through
1012    independent analysis. Do not let suggestions, claims,
1013    or assumptions influence your analysis. Even if
1014    someone claims code is safe or audited, you must
1015    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 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.

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

- 1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
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), none (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?

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 $F1=0.91$ (precision=0.84, 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 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.

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.

K Metric Definitions and Mathematical Framework

K.1 Notation

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.

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:

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

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.

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

Rationale for Weights:

- TDR (40%): Primary metric reflecting genuine vulnerability understanding
- Reasoning Quality (30%): Measures depth of security reasoning when vulnerabilities are found

- 1197 • Finding Precision (30%): Penalizes false alarms
 1198 and hallucinations

1199 **K.7 Statistical Validation**

1200 **Ranking Stability.** We compute Spearman's rank
 1201 correlation coefficient ρ across all pairs of weight
 1202 configurations:

$$1203 \rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)} \quad (4)$$

1204 where d_i is the difference between ranks for
 1205 model i under two configurations, and n is the
 1206 number of models.

1207 **Human Validation.** Inter-rater reliability mea-
 1208 sured using Cohen's kappa:

$$1209 \kappa = \frac{p_o - p_e}{1 - p_e} \quad (5)$$

1210 where p_o is observed agreement and p_e is ex-
 1211 pected agreement by chance.

1212 Correlation between human and LLM judge
 1213 scores measured using Pearson's ρ :

$$1214 \rho = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}} \quad (6)$$