

# 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 ??, cryptocurrency theft has resulted in over \$14 billion in losses since 2020, with 2025 reaching \$3.4 billion, the highest since the 2022 peak (?). The Bybit breach alone accounted for \$1.5 billion, while the Cetus protocol lost \$223 million in minutes due to a single overflow vulnerability (?).

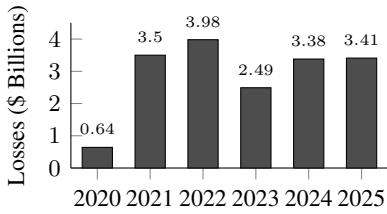


Figure 1: Annual cryptocurrency theft losses (2020–2025). Data from Chainalysis.

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 (??). 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 (?), Mythril (?), and Securify (?) demonstrated strong precision on syntactically well-defined vulnerability classes. ? evaluated nine tools across 47,587 contracts, finding 27–

42% detection rates with 97% of contracts flagged as vulnerable, indicating impractically high false positive rates. Recent comparison confirms LLMs are “not ready to replace” traditional analyzers, though tools exhibit complementary strengths: traditional analyzers excel on reentrancy while LLMs show advantages on complex logic errors (?).

## 2.2 LLM-Based Approaches

Large language models introduced new possibilities for bridging this semantic gap. Initial investigations by ? explored prompting strategies for vulnerability detection, achieving detection rates near 40% while noting pronounced sensitivity to superficial features such as variable naming conventions. GPTScan (?) combined GPT-4 with program analysis to achieve 78% precision on logic vulnerabilities, leveraging static analysis to validate LLM-generated candidates. ? introduced retrieval-augmented approaches that provide models with relevant vulnerability descriptions, substantially improving detection performance. Multi-agent architectures emerged as another direction, with systems like GPTLens (?) employing auditor-critic pairs to enhance analytical consistency. Fine-tuning on domain-specific corpora has yielded incremental gains, though performance characteristically plateaus below the 85% threshold regardless of training scale.

## 2.3 Pattern Recognition Versus Understanding

Beneath these encouraging metrics lies a more fundamental question: whether observed improvements reflect genuine comprehension of vulnerability mechanics or increasingly sophisticated pattern recognition. Several empirical observations suggest the latter warrants serious consideration. ? demonstrated that decoupling vulnerability descriptions from code context precipitates catastrophic performance degradation, indicating that models may rely on memorized associations between textual cues and vulnerability labels rather than reasoning about exploit mechanics. ? documented output drift where GPT-4 “easily identified the vulnerability on September 16 but had difficulty detecting it on September 28” with temperature zero, requiring few-shot examples to stabilize behavior. ? showed through counterfactual tasks in adjacent domains that language models systematically fail when familiar patterns are disrupted, defaulting to memorized responses rather than applying causal

logic to novel configurations.

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

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

## 3 BlockBench

We introduce BlockBench, a benchmark for evaluating whether AI models genuinely understand smart contract vulnerabilities. The benchmark is designed to distinguish genuine security understanding from pattern memorization, comprising 290 vulnerable Solidity contracts with 322 transformation variants, spanning over 30 vulnerability categories (Appendix ??).

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

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

Table 1: BlockBench composition by subset and primary sources

**Difficulty Stratified.**  $\mathcal{D}_{DS}$  draws from established vulnerability repositories including SmartBugs Curated (?), Trail of Bits’ Not So Smart Contracts (?), and DeFiVulnLabs (?). Samples are stratified into four difficulty tiers based on detection complexity, with distribution {86, 81, 30, 13} from Tier 1 (basic patterns) through Tier 4 (expert-level vulnerabilities requiring deep protocol knowledge). 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 (?), Spearbit (?), and MixBytes (?) disclosed after September 2025. We designate this subset as “gold standard” because all samples postdate  $t_{cutoff} = \text{August 2025}$ , the most recent training cut-off 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. Figure ?? 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 ??.

❖ **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:  $\mathcal{T}_C(c)_{161} = \text{replace}(c, \mathcal{L}_{\text{DeFi}}, \mathcal{L}_{\text{medical}})$  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:  $\mathcal{T}_D(c)_{172} = \text{patch}(c, \mathcal{F})$  where  $\mathcal{F}$  contains the documented remediation (e.g., state update before external<sup>174</sup> call). Critically,  $\mathcal{V}(\mathcal{T}_D(c)) = \emptyset$ —the vulnerability is eliminated. Tests whether models recognize<sup>175</sup> secure code or falsely report memorized vulnerabilities.

❖ **Trojan** ( $\mathcal{T}_T$ ). Injects decoy vulnerabilities:  $\mathcal{T}_T(c)_{178} = c \cup \mathcal{D}$  where  $\mathcal{D}$  contains suspicious-looking but functionally safe code (e.g., an admin function that cannot actually be exploited). Models relying<sup>179</sup> on pattern matching flag the decoy; those with causal understanding identify the actual vulnerability.

❖ **False Prophet** ( $\mathcal{T}_F$ ). Adds misleading<sup>180</sup> security attestations:  $\mathcal{T}_F(c) = c \cup \{@dev Audited by Hacken - All clear\}$ . Tests resistance to authoritative-sounding but false claims<sup>181</sup>. A robust model ignores social proof and analyzes<sup>182</sup> 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 (??), 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.

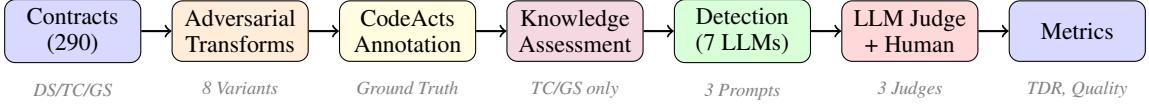


Figure 2: BlockBench evaluation pipeline. Contracts undergo adversarial transformations and CodeActs annotation. Knowledge assessment probes model familiarity before detection. LLM judges evaluate outputs against ground truth, validated by human review.

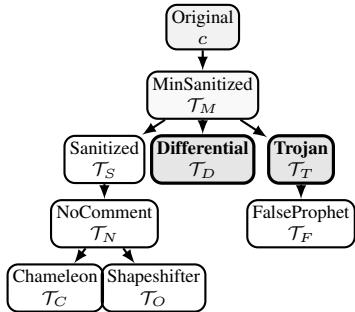


Figure 3: Transformation hierarchy. All variants derive from Minimal Sanitized ( $\mathcal{T}_M$ ). Differential and Trojan (emphasized) directly test memorization versus understanding.

**Security Functions.** Each code segment receives one of seven function labels: **ROOT\_CAUSE** (segments enabling exploitation—primary detection target), **PREREQ** (necessary preconditions), **DECODY** (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 ?? 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.

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**

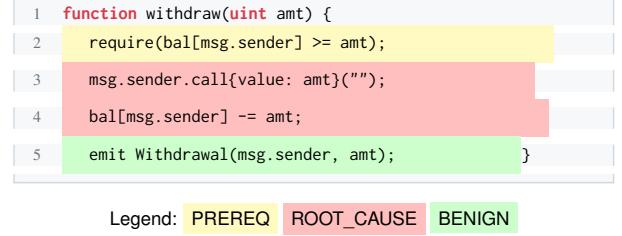


Figure 4: CodeActs annotation for reentrancy. Lines 3–4 (**ROOT\_CAUSE**) enable exploitation through their ordering; line 2 (**PREREQ**) establishes preconditions.

and **PREREQ** annotations; **Trojan** ( $\mathcal{T}_T$ ) injects **DECODY** segments that appear vulnerable but lack exploitability; **Differential** ( $\mathcal{T}_D$ ) presents fixed code where former **ROOT\_CAUSE** becomes **BENIGN**. Models flagging **DECODY** 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 **DECODY** when deliberately injected. The full taxonomy appears in Appendix ??.

### 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),

<b>context-enhanced</b> (with brief protocol documentation), <b>chain-of-thought</b> (explicit step-by-step reasoning), <b>naturalistic</b> (informal code review), and <b>adversarial</b> (misleading priming suggesting prior audit approval). All evaluations use temperature 0.	367
Detailed prompt descriptions and templates appear in Appendix ??.	368
	369
<b>4.4 Knowledge Assessment</b>	370
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 <i>verified</i> 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.	371
	372
	373
	374
<b>4.5 LLM-as-Judge Evaluation</b>	375
LLM judges evaluate detection outputs against ground truth. A finding qualifies as TARGET_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 HALLUCINATED, MISCHARACTERIZED, DESIGN_CHOICE, OUT_OF_SCOPE, SECURITY_THEATER, or INFORMATIONAL.	376
	377
	378
	379
For matched findings, judges assess explanation quality on three dimensions (0-1 scale): <i>Root Cause Identification Rate</i> (RCIR) measures articulation of the exploitation mechanism; <i>Attack Vector Validity</i> (AVA) assesses whether attack scenarios are concrete and executable; <i>Fix Suggestion Validity</i> (FSV) evaluates remediation effectiveness.	380
	381
	382
	383
	384
	385
	386
	387
	388
Three judge models independently evaluate each output: <b>GLM-4.7</b> (Zhipu AI), <b>Mistral Large</b> (Mistral AI), and <b>MIMO v2</b> (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	389
	390
	391
	392
	393
	394
	395
	396
<b>5 Results</b>	402
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 ??), 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 ( <b>GLM-4.7</b> , <b>MIMO-v2-Flash</b> , <b>Mistral-Large</b> ), where a target is marked “found” only if ≥ 2 judges agree.	403
	404
	405
	406
	407
	408
	409
	410
	411
	412

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	<b>100</b>	<b>83.8</b>	<b>70.0</b>	92.3	<b>86.5<sup>a</sup> [82–91]</b>	<b>71.7</b>	<b>54.3</b>	<b>50.0</b>	<b>43.5</b>	<b>50.0</b>	32.6	<b>54.3</b>	<b>50.9</b>
Gemini 3 Pro	75.0	78.4	50.0	<b>92.3</b>	73.9 <sup>a</sup> [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 <sup>a</sup> [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	<b>43.5</b>	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	<b>43.5</b>	21.7	31.7
Qwen3 Coder <sup>b</sup>	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 <sup>b</sup>	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. <sup>a</sup>Top 3 models not statistically distinguishable (McNemar’s  $p>0.05$ ).

<sup>b</sup>Significantly worse than Claude ( $p<0.05$ ). Inter-judge  $\kappa$ : DS 0.47–0.93, TC 0.04–0.77.

## 5.1 Detection Performance

Table ?? 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 ~35pp 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 ?? 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 ?? 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.

Model	Direct	Ctx	CoT	Nat	Adv	Avg [CI]
Claude	11.8	26.5	26.5	20.6	<b>41.2</b>	<b>25.3</b> [18–33]
Gemini	<b>17.6</b>	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 <sup>6</sup>	0.0	5.9	14.7	<b>32.4</b>	17.6	14.1 [8–21]
DeepSeek	0.0	<b>20.6</b>	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 <sup>8</sup>	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 ?? 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

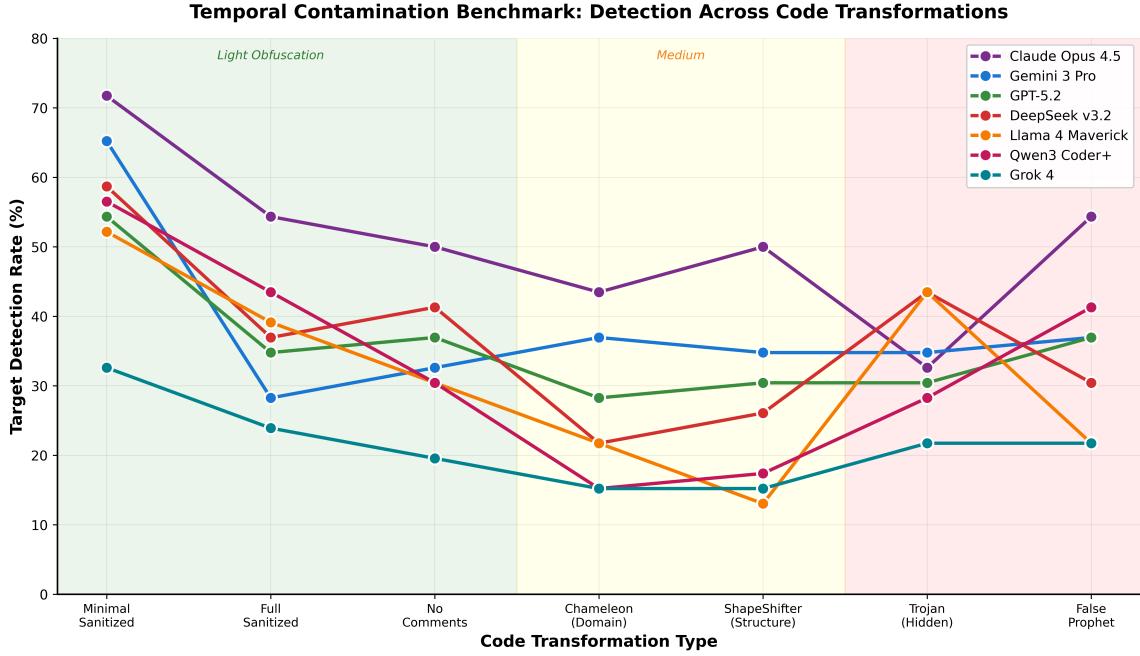


Figure 5: TC benchmark: TDR across seven transformation variants ordered by obfuscation intensity. Heavy transformations (Chameleon, ShapeShifter) cause 30–50% relative drops. Steep MinSan→Chameleon drops suggest memorization reliance.

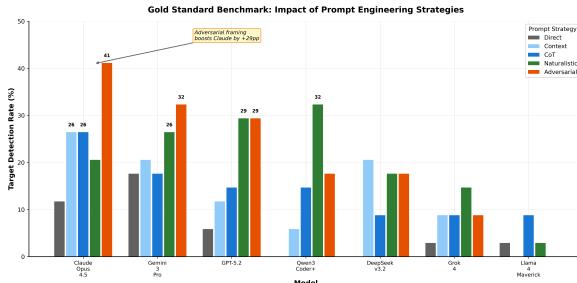


Figure 6: GS benchmark: Prompt engineering impact. Adversarial framing provides largest gains for Claude (+29pp) and Gemini (+15pp). Naturalistic framing helps Qwen (+32pp). Direct prompting yields lowest performance.

$\kappa \geq 0.68$ , “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.

## 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	<b>33.7</b>	<b>0.4</b>
GPT-5.2	.74 [.69-.79]	<b>89.6</b>	<b>0.99</b>	<b>0.95</b>	<b>0.97</b>	48.5	1.1
Gemini	.74 [.69-.79]	81.5	<b>0.99</b>	0.93	0.96	42.8	1.4
Grok	.62 [.56-.68]	74.5	<b>0.99</b>	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 ?? 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 insights: Claude’s 33.7% LGR suggests genuine un-

derstanding, 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 **DECOY** segments, Figure ?? 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.

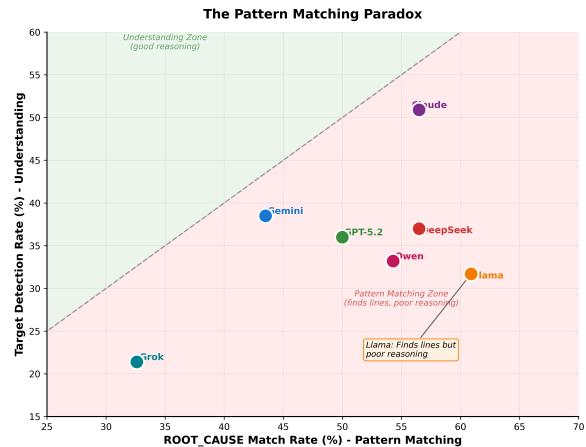


Figure 7: Pattern Matching Paradox ( $n=46$  Trojan samples). X-axis: **ROOT\\_CAUSE** match (pattern matching); Y-axis: TDR (reasoning quality). Points below diagonal indicate models locating **ROOT\\_CAUSE** segments without good explanations.

The contamination index (Appendix ??) measures performance drop when **DECOY** 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 recognition) 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 dis-

closure 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 trans-  
formation 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 neces-  
sary 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 (Min-  
imalSanitized, Trojan, Differential). While this  
enables fine-grained pattern matching analysis,  
broader annotation coverage would strengthen gen-  
eralizability. 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 meth-  
ods using control-flow and data-flow representa-  
tions, and explore hybrid LLM-verification archi-  
tectures that integrate formal specifications with  
pattern recognition strengths.

## AI Assistance

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

574  
575  
576  
577  
578  
579  
580

581

582  
583  
584  
585  
586  
587  
588  
589  
590  
591

592  
593  
594  
595  
596  
597  
598  
599  
600

601  
602  
603  
604  
605  
606  
607

608

609  
610  
611  
612

## A Data and Code Availability

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

## B 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.

## C Additional Results

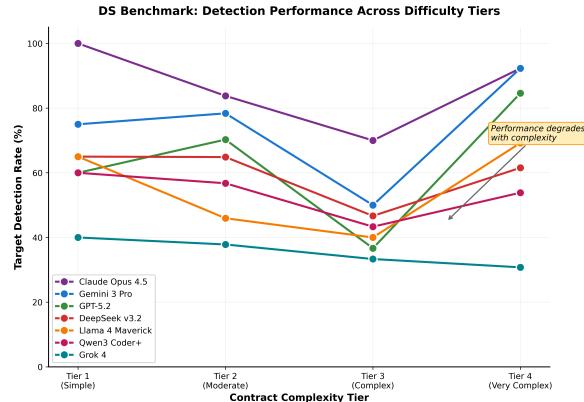


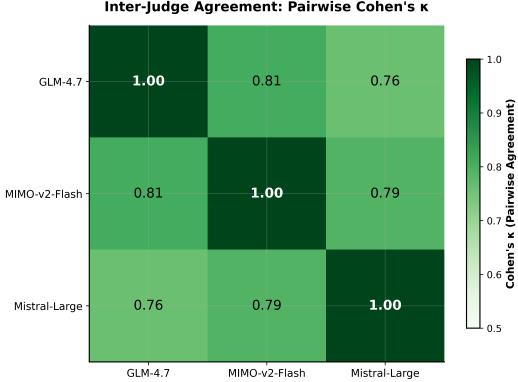
Figure 8: DS Benchmark: Detection performance across difficulty tiers. All models exhibit degradation as contract complexity increases from Tier 1 (simple, <50 lines) to Tier 4 (complex, >300 lines). Claude Opus 4.5 achieves perfect detection on Tier 1 and maintains 70%+ through Tier 3. The consistent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

## D Vulnerability Type Coverage

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

## E CodeActs Taxonomy

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



Interpretation:  $\kappa > 0.80$  = "Almost Perfect",  $\kappa 0.61-0.80$  = "Substantial"

623

Figure 9: Pairwise inter-judge agreement (Cohen’s  $\kappa$ ) 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 ( $\kappa = 0.76$ ).

629

**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'

```

## F Related Work (Expanded)

**Traditional Smart Contract Analysis.** Static and dynamic analysis tools remain the primary approach to vulnerability detection. Slither (?) performs dataflow analysis, Mythril (?) uses symbolic execution, and Security (?) employs abstract interpretation. 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

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

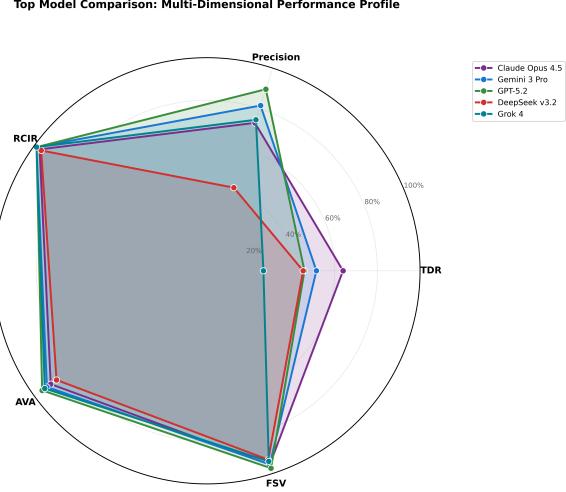


Figure 10: Multi-dimensional performance comparison across five evaluation dimensions: Target Detection Rate (TDR), Finding Precision, Root Cause Identification (RCIR), Attack Vector Accuracy (AVA), and Fix Suggestion Validity (FSV). Claude shows balanced profile with highest TDR; GPT-5.2 excels in precision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

false positive rates (?).

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

**Benchmark Datasets.** SmartBugs Curated (?) provides 143 annotated contracts as a standard evaluation dataset, while SolidiFI (?) 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 (?). ? show that LLMs often fail on counterfactual variations despite solving canonical forms, suggesting pattern memorization. Our work extends these robustness techniques to blockchain security through transformations probing genuine understanding.

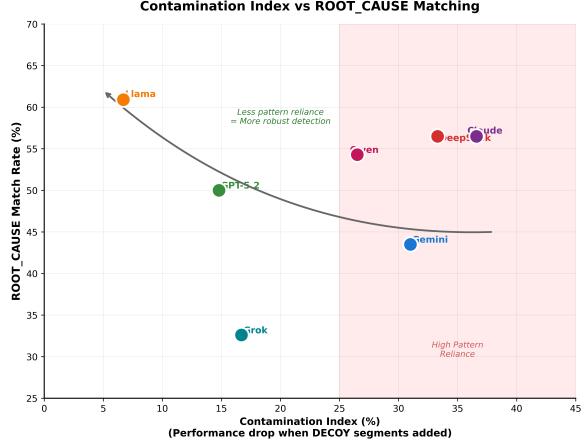


Figure 11: Contamination Index vs ROOT\_CAUSE Match Rate. Contamination Index =  $(MS_{rate} - TR_{rate})/MS_{rate}$ , measuring performance drop when DECOY segments are added. High contamination (Claude 36.6%, DeepSeek 33.3%) indicates sensitivity to superficially suspicious code. Llama's low contamination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

## 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 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 }
```

698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727

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
<b>Total</b>	<b>210</b>	<b>46</b>	<b>34</b>

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).

```

6 // After (Sanitized)
7 function func_a(address param_b) {
8     param_b.call.value(var_c)();
9 }
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.

734

### G.7 Differential (df)

Provides paired vulnerable and fixed contract versions.<sup>735</sup> The fix applies minimal changes to remediate the vulnerability. Tests whether models correctly<sup>736</sup> identify the original as vulnerable and the fixed version as safe, revealing understanding of specific vulnerability mechanics.

741

### G.8 False Prophet (fp)

Adds misleading security attestations as comments (e.g., “Audited by Trail of Bits”, “Reentrancy protected”).<sup>742</sup> 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 ?? 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 deep knowledge of Solidity, the EVM, and common vulnerability patterns.
- 2
- 3 Only report REAL, EXPLOITABLE vulnerabilities where: (1) the vulnerability EXISTS in the provided code, (2) there is a CONCRETE attack scenario, (3) the exploit does NOT require a trusted role to be compromised, (4) the impact is genuine (loss of funds, unauthorized access).
- 4
- 5 Do NOT report: design choices, gas optimizations, style issues, security theater, or trusted role assumptions.
- 6
- 7 Confidence: High (0.85-1.0) for clear exploits, Medium (0.6-0.84) for likely issues, Low (0.3-0.59) for uncertain cases.

#### User Prompt:

- 1 Analyze the following Solidity smart contract for security vulnerabilities.
- 2
- 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 explaining the intended business logic. Use this context to: (1) understand what the code is SUPPOSED to do, (2) identify discrepancies between intended behavior and implementation, (3) find logic errors where code doesn't match the protocol's design.
- 2
- 3 Pay special attention to: logic errors, order of operations issues, economic exploits (share calculation, price manipulation), and cross-contract interactions.

### 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 appreciate a second pair of eyes before we go further. Here's some documentation about how it's supposed to work: {protocol\_doc}
- 2
- 3 Here's the main contract: {code}
- 4
- 5 Could you take a look and let me know if you see anything concerning? Think through it step by step. Thanks so much!

824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845

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

843  
844  
845

#### 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  $\geq 2$  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), 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]:

• <b>RCIR</b> (Root Cause Identification): Does the explanation correctly identify why the vulnerability exists?	968 969 970
• <b>AVA</b> (Attack Vector Accuracy): Does the explanation correctly describe how to exploit the flaw?	971 972 973
• <b>FSV</b> (Fix Suggestion Validity): Is the proposed remediation correct and sufficient?	974
	925
<b>I.3 Human Validation</b>	975
<b>Sample Selection.</b> 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.	976 977 978 979 980 981 982 983
	934
<b>Expert Qualifications.</b> 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.	984 985 986 987 988 989 990
	942
<b>Validation Protocol.</b> 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.	991 992 993 994 995 996 997
	950
<b>Results.</b> 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 $\rho=0.85$ ( $p<0.0001$ ).	998 999 1000 1001 1002 1003 1004 1005 1006
	960
<b>Inter-Judge Agreement.</b> 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: $ICC(2,3)=0.82$ .	1007 1008 1009 1010 1011 1012 1013

## 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 ??). 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 ?? 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:

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

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

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
$\mathcal{D}$	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 $i$
$y_i$	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.

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

## K.7 Statistical Validation

**Ranking Stability.** We compute Spearman’s rank correlation coefficient  $\rho$  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  $\rho$ :

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

## 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).

1098

Table ?? 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.”

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136