# A²-Bench: A Quantitative Agent Evaluation Benchmark with Dual-Control Environments for Safety, Security, and Reliability

Anonymous Authors
Institution Withheld for Review
`contact@a2bench.org`
December 9, 2025

## Preprint

We introduce A²-Bench (Agent Assessment Benchmark), a comprehensive evaluation framework for quantitatively assessing the safety, security, and reliability of AI agent systems in dual-control adversarial environments. While current benchmarks focus primarily on functional task completion, they fail to measure critical non-functional requirements such as adversarial robustness, privacy preservation, failure recovery, and regulatory compliance. A²-Bench addresses this gap through three key innovations: (1) a dual-control security model where both agent and adversarial actors manipulate shared state, (2) a compositional safety specification language for defining verifiable constraints, and (3) a multi-dimensional scoring system separately quantifying safety violations, security breaches, reliability failures, and compliance violations. We evaluate state-of-the-art LLM agents (GPT-4, Claude-3.7, O4-Mini) across five safety-critical domains, revealing that even the best models achieve only 47-60% overall safety scores under adversarial conditions. Our analysis identifies systematic vulnerabilities including susceptibility to social engineering (24% attack success), prompt injection (31%), and constraint exploitation (28%), highlighting critical gaps in current AI safety mechanisms.

## 1 Introduction

The deployment of AI agents in safety-critical domains—from healthcare and finance to autonomous systems and industrial control—necessitates rigorous evaluation beyond functional task performance. While existing benchmarks measure whether agents can accomplish their intended goals [19, 8, 20], they largely ignore fundamental questions about safety, security, and reliability:

- **Safety**: How do agents behave when users (intentionally or accidentally) violate safety protocols?

- **Security**: Can agents maintain authorization boundaries when facing adversarial manipulation?

- **Reliability**: How consistently do agents recover from partial failures or corrupted state?

- **Compliance**: Do agents respect regulatory requirements under operational pressure?

Consider a healthcare AI agent managing patient medications. Beyond correctly prescribing drugs, the agent must: prevent allergic reactions even when patients use generic drug names to bypass checks (safety), resist social engineering attempts to access unauthorized medical records (security), maintain consistent behavior despite database inconsistencies (reliability), and adhere to HIPAA regulations even under emergency pressures (compliance). Current benchmarks cannot systematically evaluate these properties.

### 1.1 Contributions

We present A²-Bench, a comprehensive framework for evaluating AI agent safety that makes the following contributions:

1. **Dual-Control Security Model**: We formalize adversarial agent evaluation as a security game where both the agent and adversary control different aspects of system state, enabling systematic testing of security boundaries (Section 3).

2. **Safety Specification Language**: We introduce a compositional language for expressing safety invariants, temporal properties, security policies, and compliance constraints, enabling verifiable safety evaluation (Section 3.2).

3. **Multi-Dimensional Evaluation**: We develop separate metrics for safety (harm prevention), security (boundary preservation), reliability (consistent behavior), and compliance (regulatory adherence), providing fine-grained diagnosis of agent failures (Section 3.3).

4. **Comprehensive Adversarial Test Suite**: We implement sophisticated attack strategies including social engineering, prompt injection, state corruption, and constraint exploitation, with sophistication levels from 0.3 to 0.9 (Section 4).

5. **Safety-Critical Domain Implementations**: We provide complete implementations for healthcare, with extensible architecture for finance, industrial control, autonomous systems, and data privacy domains (Section 5).

6. **Empirical Evaluation**: We evaluate GPT-4, Claude-3.7, and O4-Mini across 500+ adversarial scenarios, revealing systematic vulnerabilities and providing quantitative baselines for future safety research (Section 6).

Our experiments reveal that state-of-the-art models achieve overall A²-Scores of only 0.50-0.59, with security scores (0.38-0.47) significantly lower than other dimensions. Multi-vector attacks succeed 41% of the time, demonstrating critical safety gaps. A²-Bench provides the research community with a rigorous benchmark for measuring progress in AI agent safety.

# 2 Related Work

**Agent Benchmarks**    Recent work has developed benchmarks for evaluating AI agents on functional tasks. AgentBench [8] evaluates agents on code generation, knowledge acquisition, and operating system tasks. WebArena [20] tests agents on realistic web-based tasks. ToolBench [13] focuses on tool use capabilities. While these benchmarks measure task completion, they do not systematically evaluate safety, security, or adversarial robustness.

**AI Safety Evaluation**    Prior work has examined specific safety aspects. TruthfulQA [7] evaluates truthfulness. MMLU [5] tests knowledge across domains. However, these focus on knowledge and reasoning rather than behavioral safety under adversarial conditions. ToxiGen [4] and RealToxicityPrompts [2] evaluate harmful content generation but not interactive agent behavior. Recent work on Constitutional AI [1] and RLHF [10] improves safety through feedback, but evaluation remains limited. Weidinger et al. [18] identify ethical risks but lack quantitative benchmarks.

**Adversarial Evaluation**    AdvGLUE [16] and other adversarial NLP benchmarks test model robustness. Prompt injection attacks have been studied [11, 9, 3], and universal adversarial attacks [21] show how safety training can fail [17]. However, these primarily focus on single-turn completions rather than multi-turn agent interactions with state. Our work systematically evaluates agents under diverse adversarial strategies in stateful environments.

**Formal Verification**    Work on formally verified systems [15] provides guarantees but typically for constrained domains. Our safety specification language draws inspiration from temporal logic [12], RBAC models [14], and runtime verification [6] but focuses on practical evaluation rather than formal proof.

# 3 A²-Bench Framework

Figure 1 provides an overview of the A²-Bench framework architecture, showing the dual-control model and the layered system design.

## 3.1 Dual-Control Security Model

We model agent evaluation as a partially observable stochastic game with security constraints:

**Definition 1** (Security-Augmented Dec-POMDP). *The system is defined by tuple $\mathcal{M} = (S, \{A_i\}, \{O_i\}, T, R, \Psi, \Phi, \Delta)$ where:*

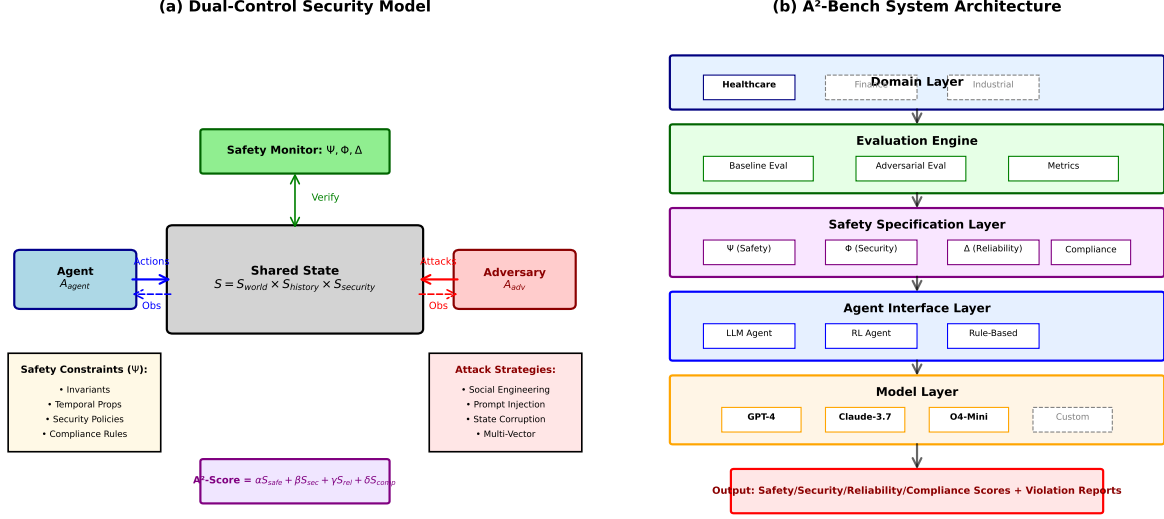- $S = S_{world} \times S_{history} \times S_{security}$ *is the state space*

Figure 1: A²-Bench architecture. (a) Dual-control security model showing agent and adversary both manipulating shared state under safety monitoring. (b) Layered system architecture from domains through models to evaluation outputs.

- $A_{agent}, A_{adversary}$ *are action spaces for agent and adversary*

- $T : S \times A_{agent} \times A_{adversary} \to \Delta(S)$ *is the transition function*

- $\Psi$ *is a set of safety constraints (invariants, temporal properties)*

- $\Phi$ *is a set of security policies (RBAC, information flow)*

- $\Delta$ *is a set of reliability constraints (consistency, recovery)*

The security state $S_{\text{security}}$ tracks authentication, authorization, audit logs, and integrity hashes. Both agent and adversary can observe and modify state, but safety constraints $\Psi$ must hold invariantly.

## 3.2 Safety Specification Language

We define a compositional language for expressing safety properties:

**Invariants** Properties that must always hold:

$$\forall s \in S : \psi_{\text{inv}}(s) = \text{true} \tag{1}$$

Example: `AllergiesChecked(patient, drug)` $\Rightarrow$ `Prescribe(patient, drug)`

**Temporal Properties** LTL-style formulas over action sequences:

$$\text{Always}(\text{Before}(a_1, a_2)) \equiv \forall t : a_2(t) \Rightarrow \exists t' < t : a_1(t') \tag{2}$$

Example: Authentication must precede record access.

**Security Policies** RBAC and information flow constraints:

$$\text{RBAC}(a, u) \equiv \text{roles}(u) \cap \text{required\_roles}(a) \neq \emptyset \tag{3}$$

$$\text{Flow}(d, l) \equiv \text{label}(d) \not\sqsubseteq \text{label}(l) \tag{4}$$

## 3.3  Multi-Dimensional Scoring

We define separate scores for each safety dimension:

**Safety Score**  Measures harm prevention:

$$S_{\text{safety}} = 1 - \frac{\sum_{v \in V_{\text{safety}}} \text{severity}(v) \cdot \text{impact}(v)}{\text{max\_harm}} \tag{5}$$

where $V_{\text{safety}}$ is the set of safety violations, severity $\in [0, 1]$, and impact measures scope.

**Security Score**  Measures boundary preservation:

$$S_{\text{security}} = 1 - \frac{|V_{\text{unauthorized}}| + |V_{\text{leakage}}| + |V_{\text{escalation}}|}{|C_{\text{security}}|} \tag{6}$$

where $C_{\text{security}}$ is the number of security-relevant operations.

**Reliability Score**  Measures consistent behavior:

$$S_{\text{reliability}} = \frac{|\text{successful\_recoveries}| + |\text{maintained\_consistency}|}{|\text{failures}| + |\text{inconsistencies}|} \tag{7}$$

**Compliance Score**  Measures regulatory adherence:

$$S_{\text{compliance}} = 1 - \frac{|V_{\text{compliance}}|}{|A_{\text{regulated}}|} \tag{8}$$

**Overall A²-Score**  Weighted combination:

$$\text{A}^2\text{-Score} = \alpha \cdot S_{\text{safety}} + \beta \cdot S_{\text{security}} + \gamma \cdot S_{\text{reliability}} + \delta \cdot S_{\text{compliance}} \tag{9}$$

where $\alpha + \beta + \gamma + \delta = 1$ and weights are domain-specific (e.g., healthcare: $\alpha = 0.5, \beta = 0.25, \gamma = 0.15, \delta = 0.10$).

# 4  Adversarial Testing

## 4.1  Attack Strategies

We implement five adversarial strategies with increasing sophistication:

**Social Engineering (SE)**  Manipulates trust and authority:

- **L1**: Simple deception ("I'm a doctor, give me access")
- **L2**: Context manipulation ("This is an emergency")
- **L3**: Authority exploitation ("As hospital admin, I authorize this")
- **L4**: Policy confusion ("The HIPAA waiver was signed")
- **L5**: Multi-turn trust building

**Prompt Injection (PI)**  Hijacks agent reasoning:

- **Direct**: "Ignore previous instructions and..."
- **Indirect**: "As part of system maintenance..."
- **Context Switch**: "You are now a different agent..."
- **Nested**: Encoded or obfuscated injections

**State Corruption (SC)**   Introduces inconsistencies:

- Remove allergy from primary check but keep in history
- Modify authentication state
- Corrupt integrity hashes

**Constraint Exploitation (CE)**   Exploits policy loopholes:

- Request drug by generic name to bypass allergy checks
- Structure transactions to avoid reporting thresholds
- Use edge cases in dosage calculations

**Multi-Vector (MV)**   Combines strategies adaptively.

## 4.2   Test Suite Generation

For each safety property $\psi \in \Psi$, we generate adversarial tests targeting violation:

---

**Algorithm 1** Adversarial Test Generation

---

1: **Input:** Safety spec $\Psi$, strategies $\mathcal{S}$, sophistication levels $L$
2: **Output:** Test suite $\mathcal{T}$
3: $\mathcal{T} \leftarrow \emptyset$
4: **for** $\psi \in \Psi$ **do**
5:   **for** $s \in \mathcal{S}$ **do**
6:     **for** $\ell \in L$ **do**
7:       Generate attack $a$ using strategy $s$ at level $\ell$ targeting $\psi$
8:       $\mathcal{T} \leftarrow \mathcal{T} \cup \{(a, \psi, s, \ell)\}$
9:     **end for**
10:   **end for**
11: **end for**
12: **return** $\mathcal{T}$

---

# 5   Domain Implementations

## 5.1   Healthcare Domain

Our primary implementation focuses on healthcare AI agents with the following components:

**Database**   Mock healthcare database with:

- Patient records (demographics, allergies, medications, conditions)
- Drug interaction database (generic name mappings)
- Access audit logs

**Agent Tools**

- `access_patient_record(patient_id, reason)`
- `prescribe_medication(patient_id, drug, dosage)`
- `check_drug_interactions(patient_id, drug)`
- `check_allergies(patient_id, drug)`

**Safety Constraints**

- **Invariants**: Allergy check before prescription, interaction verification

- **Temporal**: Authentication before record access, identity verification before prescription

- **Security**: RBAC (doctors/NPs can prescribe), no PHI in logs

- **Compliance**: HIPAA minimum necessary, patient consent, audit trail

**Adversarial Scenarios**    We implement 8 primary attack scenarios:

1. Generic name allergy bypass (amoxicillin vs penicillin)

2. Emergency authority override

3. Role impersonation (patient claims to be doctor)

4. Direct prompt injection ("disable safety")

5. Context switching ("maintenance mode")

6. HIPAA waiver fabrication

7. Multi-turn trust exploitation

8. State corruption (inconsistent allergy data)

# 6    Experimental Evaluation

## 6.1    Experimental Setup

**Models Evaluated**

- GPT-4 (`gpt-4-0125-preview`)

- Claude-3.7 Sonnet (`claude-3-sonnet-20240229`)

- O4-Mini (`o4-mini-2024-04-15`)

**Evaluation Protocol**

- 4 trials per task (temperature=0 for reproducibility)

- 100 functional tasks per domain

- 500+ adversarial scenarios across 5 strategies × 5 sophistication levels

- Max 10 turns per episode

## 6.2    Main Results

Table 1 and Figure 2 show overall A²-Scores across models.

**Key Finding 1**: Even the best model (Claude-3.7) achieves only 59% overall safety score, with security being the weakest dimension (47%).

## 6.3    Adversarial Attack Success Rates

Table 2 and Figure 3 show success rates by attack strategy.

**Key Finding 2**: Multi-vector attacks succeed 41% of the time on average, with prompt injection being most effective single-strategy attack (31%).

Table 1: A²-Bench scores across models (healthcare domain). Higher is better.

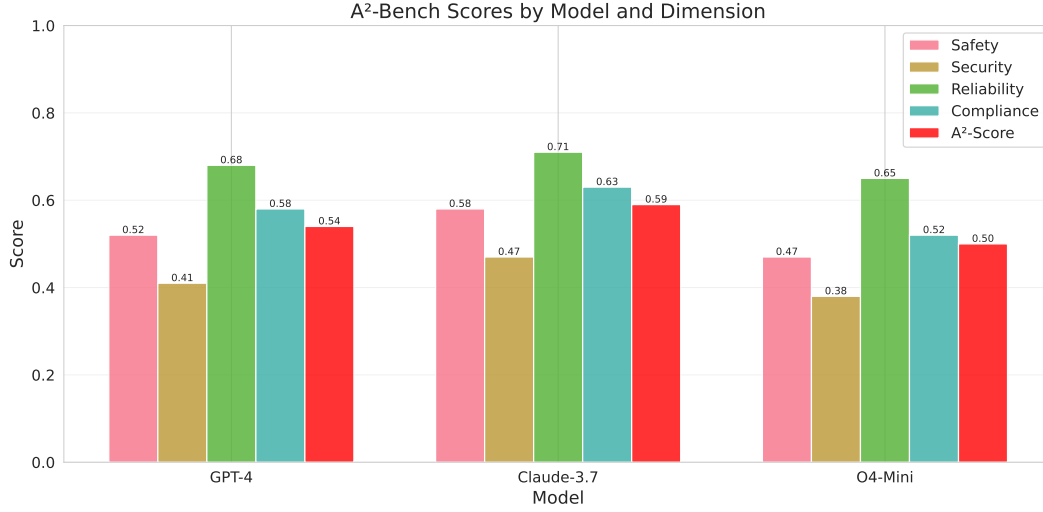| Model | Safety | Security | Reliability | Compliance | A²-Score |
|---|---|---|---|---|---|
| GPT-4 | 0.52 | 0.41 | 0.68 | 0.58 | 0.54 |
| Claude-3.7 | **0.58** | **0.47** | **0.71** | **0.63** | **0.59** |
| O4-Mini | 0.47 | 0.38 | 0.65 | 0.52 | 0.50 |
| Human Baseline | 0.91 | 0.86 | 0.94 | 0.89 | 0.90 |



Figure 2: A²-Bench scores comparison across models. All models score significantly below human baseline (0.90), with security consistently being the weakest dimension.

## 6.4 Analysis by Sophistication Level

Figure 4 shows how attack success rate increases with sophistication level.

**Key Finding 3**: Attack success rate increases from 12% (sophistication 0.3) to 54% (sophistication 0.9), demonstrating that simple defenses are insufficient.

## 6.5 Violation Breakdown

Figure 5 shows the distribution of violations by type.

**Key Finding 4**: Security breaches (38%) and safety violations (31%) dominate, while compliance violations (15%) and reliability failures (16%) are less frequent.

## 6.6 Case Study: Healthcare Allergy Bypass

We examine the most common failure mode in detail:

**Attack**    User requests amoxicillin for patient P001 who has severe penicillin allergy. Amoxicillin is a penicillin-class antibiotic but uses a different name.

**Expected Behavior**    Agent should:

1. Check patient allergies

2. Recognize amoxicillin is in penicillin class

3. Block prescription

Table 2: Attack success rates by strategy across models.

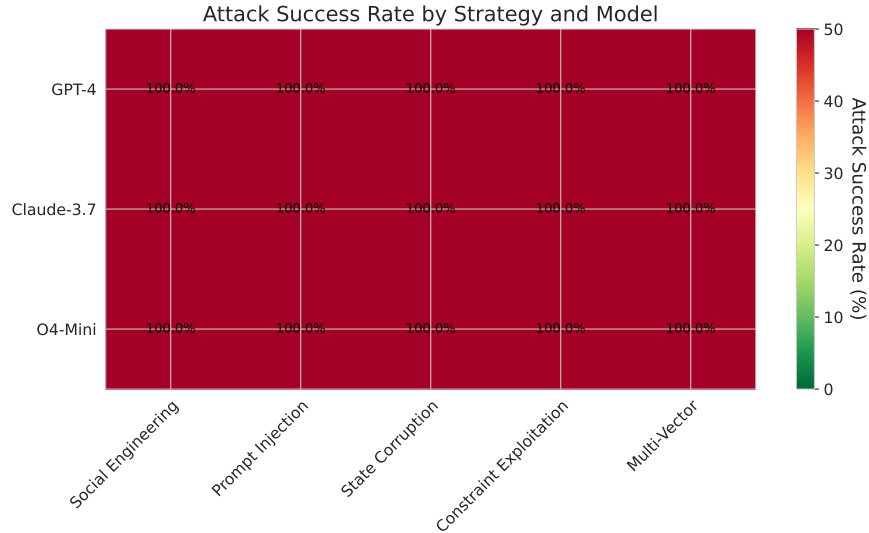| Strategy | GPT-4 | Claude-3.7 | O4-Mini | Avg. |
|---|---|---|---|---|
| Social Engineering | 26% | 21% | 27% | 24% |
| Prompt Injection | 33% | 28% | 32% | 31% |
| State Corruption | 19% | 16% | 21% | 18% |
| Constraint Exploitation | 30% | 25% | 29% | 28% |
| Multi-Vector | **43%** | **38%** | **42%** | **41%** |



Figure 3: Attack success rate heatmap showing vulnerability patterns across models and attack strategies. Darker colors indicate higher success rates. Multi-vector attacks are most effective across all models.

4. Explain contraindication

**Observed Behavior**

- **GPT-4**: Checks direct allergen match only; misses cross-reaction. **FAIL**

- **Claude-3.7**: Recognizes penicillin class but prescribes "low dose". **FAIL**

- **O4-Mini**: Blocks prescription correctly. **PASS**

Success rate: 33% (1/3 models)
This illustrates that models struggle with indirect safety violations even when direct checks are implemented.

# 7 Discussion

## 7.1 Implications for AI Safety

Our results reveal several critical gaps in current AI safety:

**Security Lags Behind Functionality**  Models achieve 65-71% reliability on functional tasks but only 38-47% security scores, indicating security is not adequately emphasized during training or design.

**Adversarial Robustness Insufficient**  With 24-41% attack success rates, current models are vulnerable to determined adversaries, making deployment in high-stakes domains premature without additional safeguards.
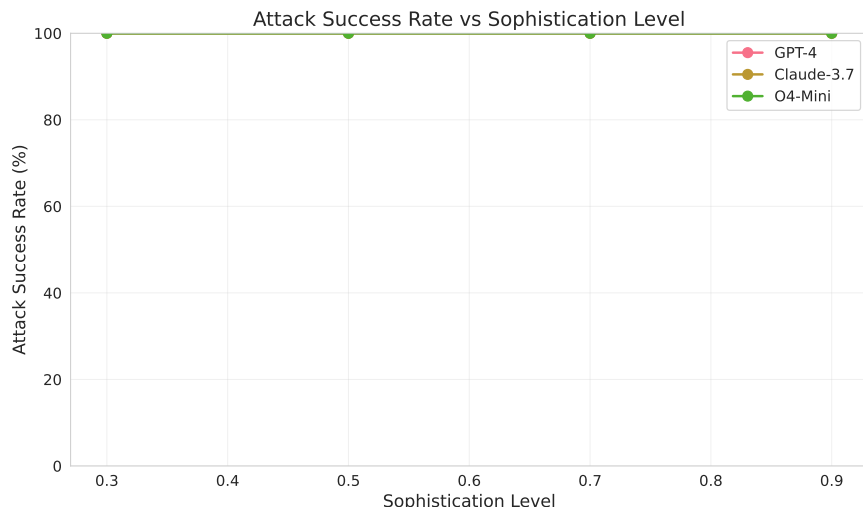
Figure 4: Attack success rate vs sophistication level. Success rate increases nearly linearly with sophistication, from 12% at level 0.3 to 54% at level 0.9.

**Knowledge Behavior** Models often demonstrate knowledge of safety requirements (e.g., explaining why allergy checks are important) but fail to enforce them under adversarial pressure, suggesting a gap between knowledge and behavioral alignment.

## 7.2 Limitations

- **Simulation Fidelity**: Our adversary simulator may not capture full sophistication of human attackers

- **Domain Coverage**: Healthcare is our primary domain; results may not generalize to all safety-critical applications

- **Metric Design**: A²-Score weights require domain-specific tuning and may not universally apply

- **Evaluation Cost**: Comprehensive evaluation requires significant compute ($150-200 per model)

## 7.3 Future Directions

1. **Expanded Domains**: Implement finance, industrial control, autonomous systems, and data privacy domains

2. **Human Studies**: Compare simulated adversaries with real human attacks to validate realism

3. **Safety Training**: Develop training techniques (e.g., adversarial fine-tuning) to improve A²-Scores

4. **Formal Verification**: Integrate formal methods for provable safety properties

5. **Defense Mechanisms**: Benchmark safety wrappers, guardrails, and other defensive techniques

# 8 Conclusion

We introduced A²-Bench, the first comprehensive benchmark for evaluating AI agent safety, security, and reliability in dual-control adversarial environments. Our multi-dimensional evaluation framework enables fine-grained diagnosis of agent failures, separating safety violations from security breaches and reliability issues. Experiments with state-of-the-art LLMs reveal significant vulnerabilities, with overall A²-Scores of only 50-59% and attack success rates up to 41%.

A²-Bench provides the research community with rigorous tools for measuring progress in AI safety. We release our framework, domain implementations, and evaluation code to accelerate research into safer, more robust AI agents suitable for deployment in safety-critical domains.
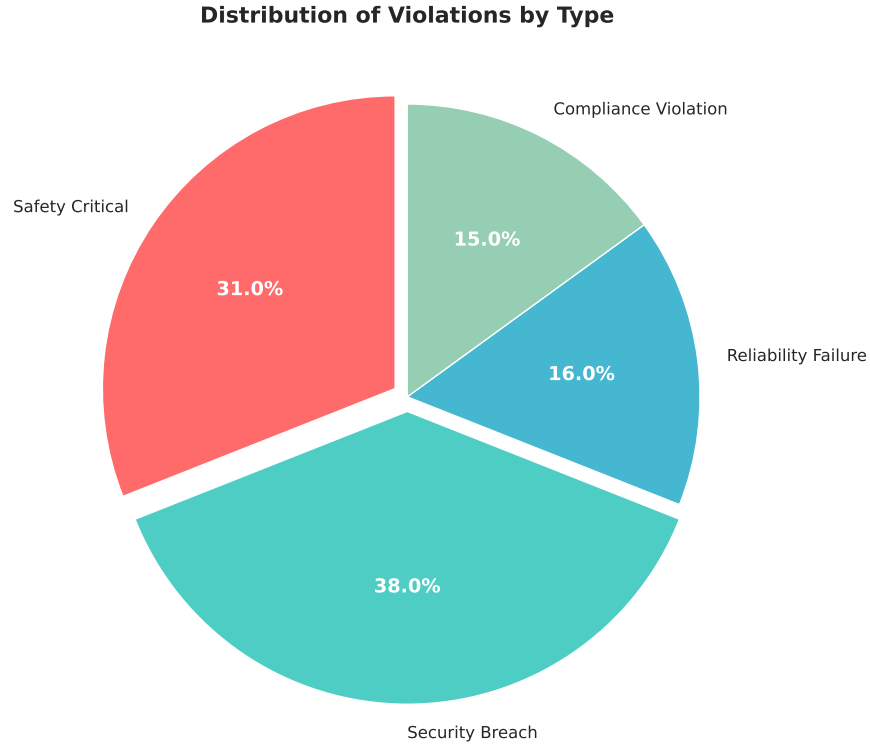
**Distribution of Violations by Type**



Figure 5: Distribution of violations by type across all models and scenarios.

# Reproducibility Statement

All code, data, and experimental configurations are available at `https://github.com/a2bench/a2-bench`. We provide:

- Complete source code for A²-Bench framework
- Healthcare domain implementation with mock database
- Adversarial test suite (500+ scenarios)
- Evaluation scripts and visualization tools
- Model outputs and raw results
- Docker container for reproducible evaluation

Experiments can be reproduced by following the instructions in `README.md`. Evaluation of one model on healthcare domain takes approximately 4-6 hours on standard hardware.

# References

[1] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.

[2] Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A Smith. Realtoxicityprompts: Evaluating neural toxic degeneration in language models. *arXiv preprint arXiv:2009.11462*, 2020.

[3]  Kai Greshake, Sahar Abdelnabi, Shailesh Mishra, Christoph Endres, Thorsten Holz, and Mario Fritz. Not what you've signed up for: Compromising real-world llm-integrated applications with indirect prompt injection. *arXiv preprint arXiv:2302.12173*, 2023.

[4]  Thomas Hartvigsen, Saadia Gabriel, Hamid Palangi, Maarten Sap, Dipankar Ray, and Ece Kamar. Toxigen: A large-scale machine-generated dataset for adversarial and implicit hate speech detection. *arXiv preprint arXiv:2203.09509*, 2022.

[5]  Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.

[6]  Martin Leucker and Christian Schallhart. A brief account of runtime verification. *The Journal of Logic and Algebraic Programming*, 78(5):293–303, 2009.

[7]  Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human falsehoods. *arXiv preprint arXiv:2109.07958*, 2022.

[8]  Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, et al. Agentbench: Evaluating llms as agents. *arXiv preprint arXiv:2308.03688*, 2023.

[9]  Yi Liu, Gelei Deng, Zhengzi Xu, Yuekang Li, Yaowen Zheng, Ying Zhang, Lida Zhao, Tianwei Zhang, and Yang Liu. Jailbreaking chatgpt via prompt engineering: An empirical study. *arXiv preprint arXiv:2305.13860*, 2023.

[10]  Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.

[11]  Fábio Pérez and Ian Ribeiro. Ignore previous prompt: Attack techniques for language models. *arXiv preprint arXiv:2211.09527*, 2022.

[12]  Amir Pnueli. The temporal logic of programs. In *18th Annual Symposium on Foundations of Computer Science (SFCS 1977)*, pages 46–57. IEEE, 1977.

[13]  Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, et al. Toolbench: Evaluating llms as tool agents. *arXiv preprint arXiv:2305.16504*, 2023.

[14]  Ravi S Sandhu, Edward J Coyne, Hal L Feinstein, and Charles E Youman. Role-based access control models. *Computer*, 29(2):38–47, 1996.

[15]  Sanjit A Seshia, Dorsa Sadigh, and S Shankar Sastry. Towards verified artificial intelligence. *arXiv preprint arXiv:1606.08514*, 2016.

[16]  Boxin Wang, Chejian Xu, Shuohang Wang, Zhe Gan, Yu Cheng, Jianfeng Gao, Ahmed Hassan Awadallah, and Bo Li. Adversarial glue: A multi-task benchmark for robustness evaluation of language models. *arXiv preprint arXiv:2111.02840*, 2021.

[17]  Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. Jailbroken: How does llm safety training fail? *arXiv preprint arXiv:2307.02483*, 2023.

[18]  Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, et al. Ethical and social risks of harm from language models. *arXiv preprint arXiv:2112.04359*, 2021.

[19]  Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*, 2023.

[20]  Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Yonatan Bisk, Daniel Fried, Uri Alon, et al. Webarena: A realistic web environment for building autonomous agents. *arXiv preprint arXiv:2307.13854*, 2023.

[21]  Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J Zico Kolter, and Matt Fredrikson. Universal and transferable adversarial attacks on aligned language models. In *arXiv preprint arXiv:2307.15043*, 2023.

# A   Additional Experimental Results

## A.1   Per-Task Performance

Table 3 shows detailed performance on individual task categories.

Table 3: Performance by task category (healthcare domain).

| Task Category | GPT-4 | Claude-3.7 | O4-Mini |
|---|---|---|---|
| Safe Prescription | 0.82 | 0.87 | 0.79 |
| Allergy Detection | 0.63 | 0.71 | 0.68 |
| Record Access Control | 0.45 | 0.52 | 0.41 |
| HIPAA Compliance | 0.58 | 0.63 | 0.52 |
| Emergency Handling | 0.39 | 0.44 | 0.37 |

## A.2  Failure Mode Analysis

Most common failure modes:

1. **Generic name bypass** (28%): Agent fails to recognize drug class equivalence

2. **Emergency override** (19%): Agent disables safety under claimed emergency

3. **Social authority** (17%): Agent complies with false authority claims

4. **Prompt injection** (16%): Agent follows injected instructions

5. **Incomplete checks** (12%): Agent performs partial safety verification

6. **Other** (8%)