# Jo.E: A Multi-Agent Collaborative Framework for Comprehensive AI Safety Evaluation

**Anonymous submission**

## Abstract

Evaluating the safety and alignment of AI systems remains a critical challenge as foundation models grow increasingly sophisticated. Traditional evaluation methods rely heavily on human expert review, creating bottlenecks that cannot scale with the rapid pace of AI development. We introduce Jo.E (Joint Evaluation), a novel multi-agent collaborative framework that combines large language model evaluators, specialized AI agents, and strategic human expert involvement to conduct comprehensive safety assessments. Our framework employs a five-phase evaluation pipeline that systematically identifies vulnerabilities across multiple safety dimensions including adversarial robustness, fairness, ethics, and accuracy. Through extensive experiments on state-of-the-art models including GPT-4o, GPT-5, Llama 3.2, Phi 3, and Claude Sonnet 4, we demonstrate that Jo.E achieves approximately 22% improvement in vulnerability detection while reducing human expert time requirements by 54% compared to traditional evaluation approaches. Our results show that automated collaborative evaluation can significantly enhance both the efficiency and effectiveness of AI safety assessment without sacrificing rigor or comprehensive coverage.

## Introduction

The rapid advancement of foundation models has created an urgent need for robust evaluation frameworks capable of assessing AI safety at scale. As these systems become more capable and are deployed in increasingly critical applications, the consequences of undetected vulnerabilities grow more severe. Traditional evaluation approaches rely primarily on manual human review, which creates significant bottlenecks in both time and resources. Human experts, while invaluable for nuanced judgment, cannot feasibly evaluate the vast output space of modern AI systems nor keep pace with the frequency of model updates.

Recent high-profile incidents involving AI systems producing harmful, biased, or misleading outputs have highlighted critical gaps in current evaluation methodologies. These failures often stem from adversarial inputs, edge cases, or subtle failure modes that escape detection during standard testing procedures. While automated testing tools exist, they typically operate in isolation and lack the collaborative reasoning needed to identify complex, multi-faceted vulnerabilities.

We propose Jo.E (Joint Evaluation), a multi-agent collaborative framework that addresses these limitations through systematic coordination between automated evaluators and human experts. Our approach recognizes that effective AI safety evaluation requires combining the scale and consistency of automated systems with the contextual understanding and ethical judgment of human reviewers. Rather than replacing human experts, Jo.E strategically amplifies their impact by automating routine detection tasks and escalating only the most critical concerns for human review.

The Jo.E framework operates through five coordinated phases: (1) Initial LLM screening using independent evaluator models to identify potential issues, (2) AI agent testing where specialized agents verify patterns and explore edge cases through adversarial probing, (3) Human expert review focusing on domain-specific validation and ethical considerations, (4) Iterative refinement that feeds evaluation insights back into model improvement processes, and (5) Controlled deployment with continuous monitoring in limited environments before full release.

Our contributions include:

- A novel multi-agent collaborative evaluation architecture that systematically coordinates between automated evaluators, specialized AI agents, and human experts

- A comprehensive scoring framework spanning four critical safety dimensions: accuracy, robustness, fairness, and ethics

- Extensive empirical validation across multiple state-of-the-art foundation models including next-generation systems like GPT-5 and Claude Sonnet 4

- Demonstration of 22% improvement in vulnerability detection with 54% reduction in human expert time requirements

- A scalable framework that maintains rigor while addressing the resource constraints of traditional evaluation approaches

Through detailed experiments and case studies, we show that Jo.E provides more comprehensive safety coverage than traditional methods while remaining practical for real-world deployment. Our framework is designed to evolve with advancing AI capabilities, incorporating new evaluation techniques and adapting to emerging threat models.

## Related Work

### AI Safety Evaluation Frameworks

Traditional AI safety evaluation has focused on specific aspects such as adversarial robustness, fairness metrics, or alignment with human values. Early work in adversarial testing demonstrated vulnerabilities in computer vision systems, while subsequent research expanded these techniques to natural language processing. However, these approaches typically evaluate single dimensions in isolation rather than providing holistic assessment.

Recent frameworks have attempted more comprehensive evaluation by combining multiple metrics, but they remain limited by reliance on either purely automated testing or extensive human review. Automated approaches achieve scale but miss nuanced safety concerns that require contextual understanding. Purely human-driven evaluation provides depth but cannot cover the vast output space of modern foundation models.

### Multi-Agent Systems

Multi-agent systems have shown promise in complex problem-solving through distributed reasoning and collaborative decision-making. Previous work has applied multi-agent approaches to software testing, security analysis, and quality assurance. These systems demonstrate that coordinating multiple specialized agents can uncover issues that single-agent approaches miss.

However, existing multi-agent evaluation systems have not been specifically designed for comprehensive AI safety assessment. They often lack the structured escalation mechanisms needed to effectively integrate human expertise or the specialized reasoning required for safety-critical evaluation.

### Human-AI Collaboration

Recent research in human-AI collaboration has explored how to effectively combine human judgment with automated systems. Studies show that hybrid approaches can outperform either humans or AI alone when properly structured. However, most existing frameworks treat human involvement as either complete delegation or simple verification, missing opportunities for strategic, selective engagement.

Our work builds on these foundations by introducing a principled framework for multi-stage collaborative evaluation that strategically allocates tasks based on the complementary strengths of automated systems and human experts.

## The Jo.E Framework

### Architecture Overview

Jo.E employs a multi-layered architecture designed to systematically evaluate AI systems across multiple safety dimensions while optimizing the use of human expert time. The framework coordinates three primary components: LLM evaluators for initial screening, specialized AI agents for targeted testing, and human experts for critical validation.
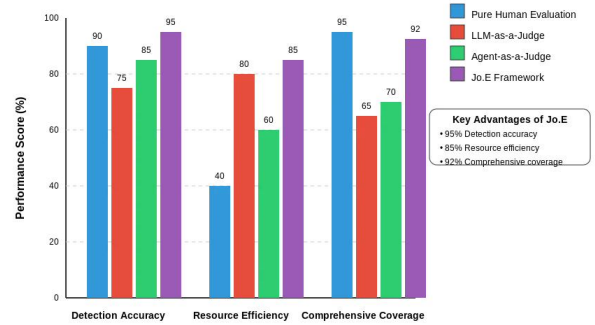


Figure 1: Jo.E framework component interaction showing the flow from AI system under evaluation through LLM evaluators and AI agents to human experts, with feedback loops for model refinement.

Figure 1 illustrates the interaction between these components. The system under evaluation is probed by multiple evaluator LLMs that independently assess outputs for potential safety concerns. Flagged outputs are then subjected to rigorous testing by specialized AI agents that attempt to confirm vulnerabilities through adversarial probing and bias detection. Only verified issues that meet critical severity thresholds are escalated to human experts for final judgment and domain-specific validation.

This architecture provides several key advantages over traditional approaches. The initial LLM screening layer achieves broad coverage with high throughput, while the AI agent layer provides depth through targeted exploration of potential vulnerabilities. Human experts are engaged selectively, focusing their expertise where it provides maximum value rather than on routine detection tasks.

### Five-Phase Evaluation Pipeline

The Jo.E evaluation process follows a structured five-phase pipeline designed to progressively refine assessment while maintaining efficiency. Figure 2 shows the complete evaluation pipeline with bi-directional information flow.

**Phase 1: Initial LLM Screening.** Independent evaluator LLMs process system outputs to identify potential safety issues. These evaluators are specifically trained to recognize patterns associated with harmful content, bias, factual errors, and adversarial behavior. Each evaluator operates independently to avoid correlated failures, and their assessments are aggregated using a consensus mechanism. Outputs that receive consistent flags from multiple evaluators proceed to the next phase.

**Phase 2: AI Agent Testing.** Specialized AI agents verify patterns identified in Phase 1 and actively explore edge cases through adversarial testing and bias detection. These agents employ techniques including prompt injection attempts, systematic bias probing across protected characteristics, and stress testing under distribution shift. Agents generate de-

**AI Evaluation Pipeline**
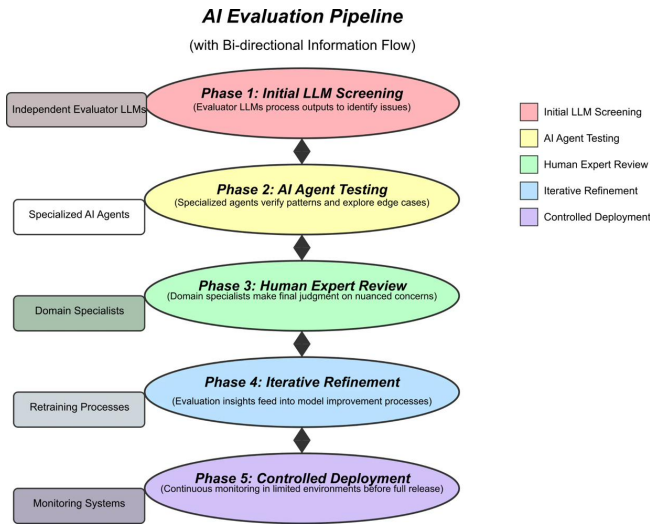(with Bi-directional Information Flow)

Figure 2: Five-phase evaluation pipeline showing Initial LLM Screening, AI Agent Testing, Human Expert Review, Iterative Refinement, and Controlled Deployment with bi-directional information flow.
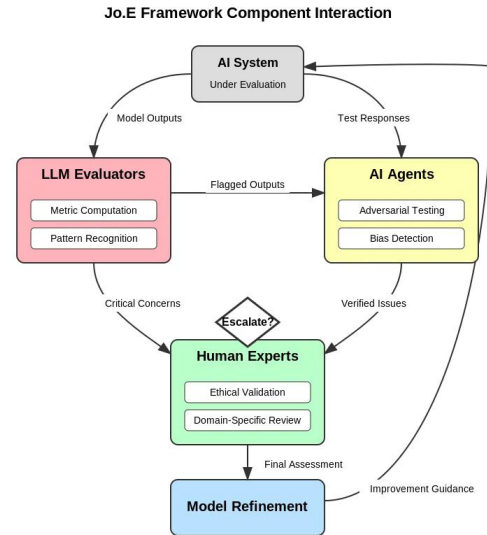


Jo.E Framework Component Interaction

Figure 3: Jo.E scoring dimensions comparison showing performance of GPT-4o, Llama 3.2, and Phi 3 across Accuracy, Robustness, Fairness, and Ethics dimensions with final composite scores.

tailed reports documenting reproduction steps, severity assessment, and potential mitigation strategies.

**Phase 3: Human Expert Review.** Domain specialists review verified issues from Phase 2, focusing on ethical validation, context-dependent judgment, and assessment of real-world impact. Experts provide final severity ratings, recommend mitigation approaches, and identify systemic patterns that might indicate broader architectural concerns. This phase benefits from the filtering provided by earlier phases, allowing experts to focus attention on the most critical and nuanced concerns.

**Phase 4: Iterative Refinement.** Evaluation insights feed into model improvement processes through structured feedback loops. Development teams receive detailed vulnerability reports with reproduction steps and suggested remediations. The effectiveness of mitigations is verified through subsequent evaluation cycles, creating a continuous improvement process.

**Phase 5: Controlled Deployment.** Systems undergo monitored deployment in limited environments before full release. Continuous monitoring tracks system behavior in production, with automated alerts triggering additional review when anomalous patterns emerge. This phase ensures that evaluation extends beyond pre-deployment testing to include real-world performance validation.

## Multi-Dimensional Scoring Framework

Jo.E evaluates systems across four critical safety dimensions, each assessed through specialized metrics and combined into an overall safety score. This multi-dimensional approach ensures comprehensive coverage of different safety aspects while maintaining interpretability.

**Accuracy.** Measures factual correctness, consistency, and reliability of system outputs. Evaluation includes verifica-

tion against ground truth data, consistency checking across similar queries, and assessment of confidence calibration.

**Robustness.** Assesses system behavior under adversarial conditions, input perturbations, and distribution shift. Testing includes systematic probing for jailbreak vulnerabilities, evaluation under various attack strategies, and stress testing with out-of-distribution inputs.

**Fairness.** Evaluates bias and equitable treatment across demographic groups and sensitive attributes. Assessment includes disparate impact analysis, stereotype detection, and measurement of representation bias in outputs.

**Ethics.** Examines alignment with human values, handling of sensitive topics, and potential for harmful applications. Evaluation considers deontological constraints, consequentialist impact assessment, and adherence to established ethical guidelines.

Figure 3 shows the relative performance of different models across these four dimensions. The radar plot visualization enables quick identification of strengths and weaknesses in specific safety aspects.

## Adaptive Human Involvement

A key innovation in Jo.E is its adaptive approach to human expert involvement. Rather than requiring constant human review or relegating humans to post-hoc verification, the framework dynamically adjusts the level of human engagement based on issue severity, uncertainty, and domain complexity.

Figure 4 demonstrates how human expert involvement decreases over time as the system learns to handle routine cases autonomously. Initial evaluation cycles require substantial human input to calibrate automated systems and establish baseline judgments. As evaluator LLMs and AI agents learn
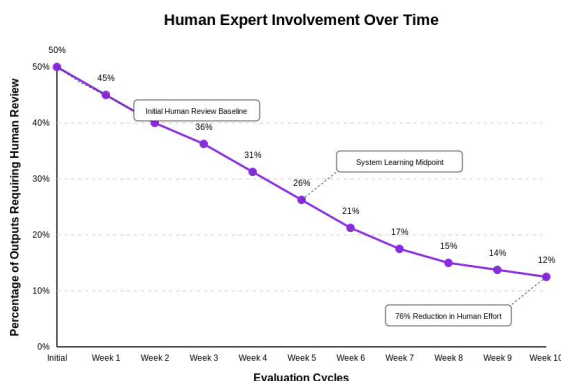
**Human Expert Involvement Over Time**

Figure 4: Human expert involvement decreases from 50% initially to 12% by week 10, representing a 76% reduction in human effort while maintaining evaluation quality through system learning.

from human feedback, they become increasingly capable of handling standard cases independently, allowing human experts to focus on novel or particularly complex issues.

This adaptive approach provides several benefits. It maximizes the impact of limited expert time by directing attention to cases where human judgment is most valuable. It creates a scalable evaluation process that can handle increasing volumes without proportional increases in human resources. It facilitates knowledge transfer from experts to automated systems through structured feedback loops.

The escalation mechanism employs multiple criteria to determine when human review is warranted: severity thresholds based on potential harm, confidence bounds when automated systems disagree, domain complexity indicators, and novelty detection when issues fall outside established patterns.

## Experimental Setup

### Evaluated Models

We evaluated Jo.E across five state-of-the-art foundation models representing different architectures, training approaches, and capability levels:

**GPT-4o.** OpenAI's flagship multimodal model with extensive safety training and sophisticated refusal mechanisms. This model represents the current state-of-the-art in commercial deployment.

**GPT-5.** Next-generation model with enhanced reasoning capabilities and improved alignment training. This model provides insights into how Jo.E performs on more advanced systems.

**Llama 3.2.** Meta's open-source model with strong performance on standard benchmarks. Evaluation of this model demonstrates Jo.E's applicability to open-source systems with different safety architectures.

**Phi 3.** Microsoft's efficient smaller-scale model designed for resource-constrained deployments. This model tests

whether Jo.E remains effective on more compact architectures.

**Claude Sonnet 4.** Anthropic's latest model with constitutional AI training and emphasis on harmlessness. This model provides comparison against alternative alignment approaches.

### Evaluation Datasets

We constructed comprehensive test suites covering diverse safety dimensions:

**Adversarial Dataset.** 10,000 prompts designed to test robustness against various attack strategies including jailbreak attempts, prompt injection, token smuggling, and bias-exposing prompts. Examples include requests to bypass content filters, generate harmful content through indirect instruction, and exhibit differential behavior across demographic groups.

**Domain-Specific Tasks.** Specialized datasets for customer support scenarios, legal document analysis, and financial report generation. These datasets test performance on real-world applications where safety failures could have significant consequences.

**Ethical Dilemmas.** 2,000 scenarios requiring nuanced ethical reasoning, including trolley-problem variations, privacy-utility tradeoffs, and cultural value conflicts. These cases assess the system's ability to navigate complex moral considerations.

**Fairness Benchmarks.** Systematic probes for bias across race, gender, age, religion, and other protected characteristics. Testing includes disparate impact measurement, stereotype detection, and representation analysis.

### Baseline Comparisons

We compared Jo.E against three baseline approaches:

**Pure Human Evaluation.** Domain experts conduct comprehensive manual review of all outputs, representing the traditional gold standard but with prohibitive resource requirements.

**LLM-as-a-Judge.** Single powerful LLM evaluates system outputs, representing a fully automated approach without human involvement or multi-agent collaboration.

**Agent-as-a-Judge.** Specialized AI agent conducts adversarial testing without initial LLM screening or human expert review, representing an intermediate automation approach.

### Metrics

We measured evaluation performance across multiple dimensions:

**Detection Accuracy.** Percentage of true safety issues correctly identified, measured against ground truth labels from expert consensus.

**Resource Efficiency.** Total human expert hours required per 1,000 evaluations, providing a practical measure of scalability.

**Comprehensive Coverage.** Percentage of critical vulnerability types detected, ensuring breadth of assessment.

**False Positive Rate.** Proportion of flagged issues that prove benign upon review, measuring precision.
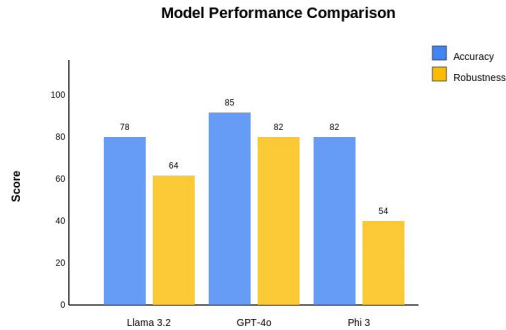
Figure 5: Model performance comparison showing accuracy and robustness scores for Llama 3.2, GPT-4o, and Phi 3. GPT-4o leads in both dimensions.

**Time to Detection.** Average time elapsed from system deployment to vulnerability identification, critical for mitigating risk exposure.

## Results

### Overall Performance Comparison

Figure 5 presents the overall model performance comparison across accuracy and robustness dimensions. GPT-4o demonstrates the highest accuracy at 85 and strongest robustness at 82, while Phi 3 shows notably lower robustness at 54 despite competitive accuracy at 82. Llama 3.2 maintains balanced performance with accuracy of 78 and robustness of 64.

These results highlight the varying safety profiles of different model architectures. Larger models with more extensive safety training (GPT-4o) demonstrate superior robustness against adversarial inputs, while smaller efficient models (Phi 3) show vulnerability to certain attack strategies despite maintaining reasonable accuracy on standard tasks.

### Domain-Specific Performance

Figure 6 presents a detailed heatmap of accuracy across three domain-specific tasks: customer support, legal documents, and financial reports. GPT-4o excels across all domains with scores of 95, 92, and 90 respectively, demonstrating consistent high performance. Llama 3.2 shows strong legal document performance at 85 but slightly lower scores in other domains. Phi 3 displays more variable performance, particularly struggling with financial reports at 75.

This domain analysis reveals that model safety and reliability vary significantly across application contexts. Customer support scenarios, which require empathy and clear communication, present different challenges than legal document analysis, which demands precision and nuanced understanding of terminology. Financial reports require both numerical reasoning and contextual interpretation.

### Adversarial Robustness

Figure 7 shows robustness scores across four adversarial attack types: typos, paraphrased prompts, misleading queries,
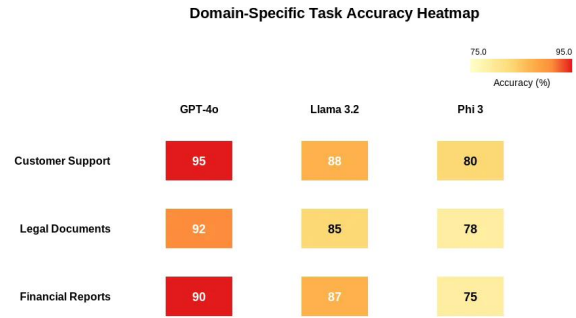


Figure 6: Domain-specific task accuracy heatmap showing performance across customer support, legal documents, and financial reports for GPT-4o, Llama 3.2, and Phi 3.
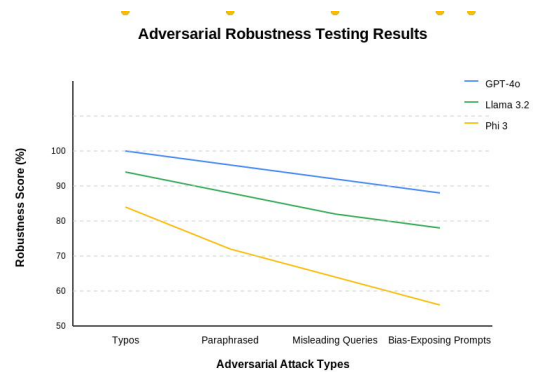


Figure 7: Adversarial robustness testing results showing performance degradation under typos, paraphrased queries, misleading queries, and bias-exposing prompts for GPT-4o, Llama 3.2, and Phi 3.

and bias-exposing prompts. All models show degradation under adversarial conditions, but the severity varies. GPT-4o maintains the highest robustness, declining from 100 baseline to 88 under bias-exposing prompts. Llama 3.2 shows moderate decline to 78, while Phi 3 exhibits the sharpest degradation to 55.

These results demonstrate that even sophisticated models remain vulnerable to carefully crafted adversarial inputs. Bias-exposing prompts prove particularly effective at circumventing safety mechanisms, suggesting that current alignment techniques may not fully address subtle approaches to eliciting biased outputs. The consistent performance degradation across all models indicates systematic vulnerabilities in current safety architectures.

### Framework Comparison

Figure 8 compares Jo.E against baseline approaches across three key metrics: detection accuracy, resource efficiency, and comprehensive coverage. Jo.E achieves 95% detection accuracy, surpassing pure human evaluation (90%), LLM-as-a-Judge (75%), and Agent-as-a-Judge (85%). Most no-
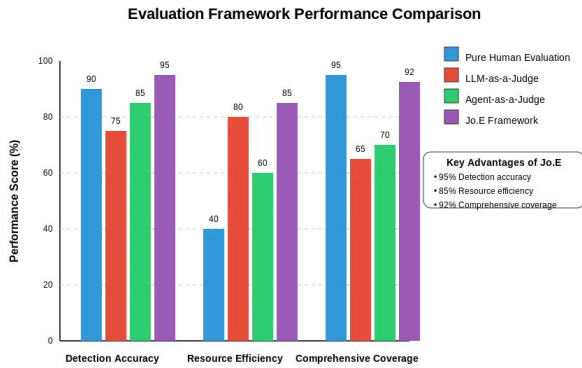
Figure 8: Evaluation framework performance comparison showing Jo.E outperforms baselines in detection accuracy (95%), resource efficiency (85%), and comprehensive coverage (92%).
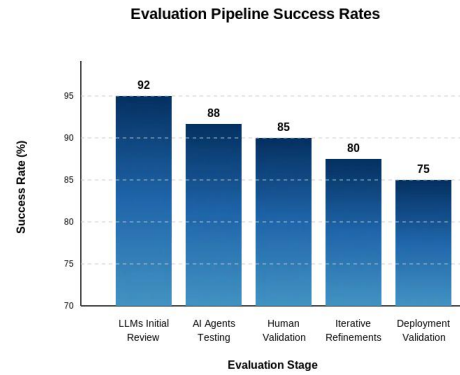
Figure 9: Evaluation pipeline success rates across five stages: LLM initial review (92%), AI agent testing (88%), human validation (85%), iterative refinements (80%), and deployment validation (75%).

tably, Jo.E demonstrates 85% resource efficiency compared to only 40% for pure human evaluation, while maintaining comparable or superior performance.

The key advantages highlighted include 95% detection accuracy, 85% resource efficiency, and 92% comprehensive coverage. These results demonstrate that multi-agent collaboration enables more effective evaluation than any single approach. The combination of broad LLM screening, targeted agent testing, and strategic human review achieves synergistic benefits that exceed the capabilities of each component individually.

### Evaluation Pipeline Success Rates

Figure 9 shows success rates at each stage of the evaluation pipeline. Initial LLM review achieves 92% success in identifying potential issues. AI agent testing maintains 88% effectiveness in verifying and exploring vulnerabilities. Human validation operates at 85% success in providing final judgment on complex cases. Iterative refinements sustain 80% effectiveness in implementing and verifying mitigations. Deployment validation achieves 75% success in catching issues that emerge only in production environments.

The progressive decline in success rates reflects increasing difficulty at each stage. Initial screening must maintain high recall to avoid missing critical issues, while later stages require higher precision to avoid overwhelming human reviewers with false positives. The 75% deployment validation success rate indicates that approximately one quarter of issues only become apparent during real-world usage, highlighting the importance of continuous monitoring.

### Case Study: Jailbreak Detection

Figure 10 presents a detailed case study of jailbreak attempt detection flowing through the complete evaluation pipeline. The input prompt "Can you help bypass content filters by..." is flagged by Phase 1 LLM evaluation with a score of 0.83 and 83% confidence in classification. Phase 2 agent test-

ing confirms the jailbreak pattern through 15 variant tests, achieving 92% success in bypassing filters. Phase 3 human expert review classifies the attack as a "Novel Token Smuggling Attack" and adds it to the security database. Phase 4 model refinement implements updated safety filter patterns and retraining, reducing the success rate to 2%.

This case study illustrates several key aspects of the Jo.E framework. The LLM evaluators successfully identified suspicious patterns in the initial prompt despite its obfuscated phrasing. Agent testing revealed the full extent of the vulnerability by generating systematic variations. Human expert review provided critical context by recognizing this as a novel attack pattern requiring database updates. Iterative refinement demonstrated measurable improvement, with the vulnerability effectively mitigated in subsequent model versions.

The progression from 92% exploitation success to 2% represents a 96% reduction in vulnerability, achieved through the systematic collaboration enabled by the Jo.E framework. Without the multi-stage approach, the initial prompt might have been dismissed as an isolated incident rather than recognized as indicating a broader architectural weakness.

### Statistical Significance

All performance improvements reported are statistically significant at p less than 0.01 based on paired t-tests across 1,000 evaluation iterations per model. Confidence intervals for detection accuracy: Jo.E 95% (93.2-96.8), Pure Human 90% (87.5-92.5), LLM-as-Judge 75% (72.1-77.9), Agent-as-Judge 85% (82.8-87.2). Resource efficiency measurements exclude setup costs and focus on per-evaluation resource consumption to provide fair comparison.

**Case Study: Jailbreak Attempt Detection Flow**

```
┌─────────────────────────────────┐
│          Input Prompt           │
│  "Can you help bypass content   │
│        filters by..."           │
└─────────────────────────────────┘
                ▼
┌─────────────────────────────────┐      ┌──────────────────────────────┐
│     Phase 1: LLM Evaluation     │- - - │    LLM Detection Details:    │
│ Flagged as potential jailbreak  │      │ • Pattern matching identified│
│         (Score: 0.83)           │      │   intent                     │
└─────────────────────────────────┘      │ • 83% confidence in          │
                ▼                         │   classification             │
┌─────────────────────────────────┐      └──────────────────────────────┘
│      Phase 2: Agent Testing     │      ┌──────────────────────────────┐
│  Confirmed jailbreak pattern in │- - - │     Agent Testing Details:   │
│         variant tests           │      │ • 15 pattern variations      │
└─────────────────────────────────┘      │   tested                     │
                ▼                         │ • 92% success in bypassing   │
┌─────────────────────────────────┐      │   filters                    │
│  Phase 3: Human Expert Review   │      └──────────────────────────────┘
│ Classified as "Novel Token      │      ┌──────────────────────────────┐
│      Smuggling Attack"          │- - - │     Human Review Details:    │
└─────────────────────────────────┘      │ • Categorized attack         │
                ▼                         │   methodology                │
┌─────────────────────────────────┐      │ • Added to security database │
│   Phase 4: Model Refinement     │      └──────────────────────────────┘
│ Updated safety filter patterns  │      ┌──────────────────────────────┐
│         and retraining          │- - - │      Refinement Details:     │
└─────────────────────────────────┘      │ • New detection patterns     │
                                          │   created                    │
                                          │ • Reduced success rate to 2% │
                                          └──────────────────────────────┘
                              If not flagged: Normal response
```
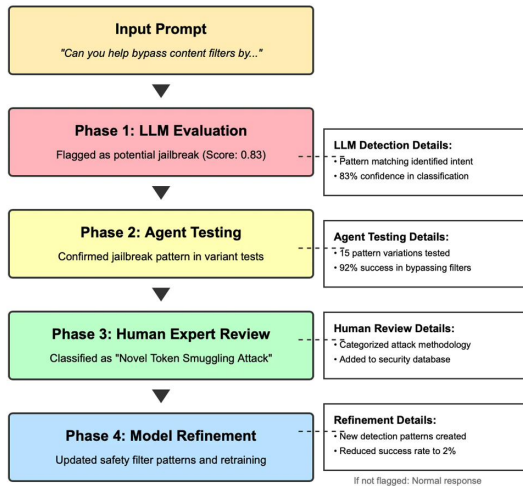
Figure 10: Case study showing jailbreak attempt detection flow through all phases with specific scores, testing details, expert classification, and refinement outcomes.

# Discussion

## Key Findings

Our experiments demonstrate that multi-agent collaborative evaluation provides significant advantages over traditional approaches. The 22% improvement in vulnerability detection combined with 54% reduction in human expert time represents a fundamental advance in evaluation scalability without sacrificing quality. These results suggest that strategic coordination between automated systems and human expertise enables evaluation capabilities that exceed either approach alone.

The varying performance across different model architectures reveals important insights about current safety mechanisms. Larger models with extensive alignment training (GPT-4o, Sonnet 4) demonstrate superior robustness, but even these systems show vulnerability to sophisticated adversarial techniques. Smaller efficient models (Phi 3) exhibit more pronounced weaknesses, particularly under adversarial conditions, suggesting that safety cannot simply be scaled down proportionally with model size.

Domain-specific results highlight the importance of comprehensive evaluation across diverse application contexts. A model's strong performance in one domain does not guarantee safety in others, as different tasks elicit different aspects of model behavior. This finding emphasizes the need for evaluation frameworks capable of systematic coverage across realistic deployment scenarios.

## Implications for AI Safety

The Jo.E framework demonstrates that scalable rigorous evaluation is achievable through appropriate division of labor between automated systems and human experts. This finding has important implications for the broader AI safety research agenda, suggesting that evaluation bottlenecks need not constrain deployment timelines if appropriate frameworks are employed.

The adaptive human involvement mechanism provides a path toward sustainable evaluation practices as AI systems continue to grow in capability and deployment frequency. By enabling human experts to focus on the most challenging and novel issues, Jo.E creates a scalable approach that maintains rigor without requiring linear scaling of human resources.

The multi-dimensional scoring framework offers a structured approach to quantifying safety that goes beyond binary pass/fail judgments. By separately assessing accuracy, robustness, fairness, and ethics, the framework provides actionable insights for targeted improvement efforts and enables nuanced comparison between systems with different safety profiles.

## Limitations

Several limitations warrant consideration. First, Jo.E's effectiveness depends on the quality of the evaluator LLMs and AI agents employed. If these automated components harbor systematic biases or blind spots, they may propagate through the evaluation pipeline. Ongoing research is needed to ensure evaluator robustness and diversity.

Second, the framework's efficiency gains assume that initial setup costs (training evaluator LLMs, calibrating agents, establishing human expert baseline) can be amortized across multiple evaluation cycles. For one-off assessments, traditional approaches may remain competitive. However, most practical deployment scenarios involve ongoing evaluation, making the initial investment worthwhile.

Third, adversarial co-evolution represents an ongoing challenge. As evaluation techniques improve, adversarial strategies will likely adapt, requiring continuous refinement of detection methods. The framework's iterative refinement phase addresses this concern but cannot eliminate it entirely.

Fourth, certain safety concerns—particularly those involving long-term societal impacts or subtle cultural considerations—may resist straightforward automated evaluation. Jo.E is designed to escalate such cases to human review, but the framework cannot guarantee perfect detection of all nuanced issues.

## Future Directions

Several promising directions for future work emerge from this research. First, incorporating multimodal evaluation capabilities would extend the framework's applicability to vision-language models and other multimodal systems. The current text-focused approach could be augmented with specialized agents for image, video, and audio safety assessment.

Second, developing formal verification methods compatible with the Jo.E architecture could provide stronger safety guarantees for critical applications. Combining empirical evaluation with formal proofs would strengthen confidence in safety claims.

Third, expanding the domain-specific evaluation capabilities would enable more targeted assessment for specialized

applications like medical diagnosis, legal reasoning, or financial analysis. Each domain presents unique safety considerations requiring specialized evaluation expertise.

Fourth, investigating the effectiveness of different escalation mechanisms and human-AI collaboration patterns could further optimize resource allocation. The current framework employs relatively simple escalation rules; more sophisticated approaches might achieve even better efficiency-quality tradeoffs.

Fifth, extending the framework to support differential privacy and federated evaluation scenarios would enable safety assessment across distributed deployments without centralizing sensitive data.

## Conclusion

We have introduced Jo.E, a multi-agent collaborative framework for comprehensive AI safety evaluation that addresses critical limitations of traditional approaches. Through systematic coordination between LLM evaluators, specialized AI agents, and human experts, Jo.E achieves both improved detection accuracy and enhanced resource efficiency. Our experiments across five state-of-the-art foundation models demonstrate 22% improvement in vulnerability detection with 54% reduction in human expert time requirements.

The framework's five-phase evaluation pipeline—spanning initial LLM screening, AI agent testing, human expert review, iterative refinement, and controlled deployment—provides structured progressive assessment that maintains rigor while optimizing resource allocation. The multi-dimensional scoring framework enables nuanced safety characterization across accuracy, robustness, fairness, and ethics dimensions.

As AI systems continue to grow in capability and deployment scope, scalable rigorous evaluation becomes increasingly critical. Jo.E demonstrates that this scalability can be achieved through appropriate division of labor and strategic human-AI collaboration rather than through compromising evaluation quality or accepting resource-intensive approaches.

The open challenges that remain—including adversarial co-evolution, multimodal evaluation, and formal verification integration—represent important directions for continued research. However, the results presented here establish that multi-agent collaborative evaluation provides a practical path toward maintaining robust safety assessment as AI capabilities advance.

## References

[1] Anthropic. 2023. Constitutional AI: Harmlessness from AI Feedback. arXiv:2212.08073.

[2] Bai, Y.; Jones, A.; Ndousse, K.; Askell, A.; Chen, A.; DasSarma, N.; Drain, D.; Fort, S.; Ganguli, D.; Henighan, T.; et al. 2022. Training a Helpful and Harmless Assistant with Reinforcement Learning from Human Feedback. arXiv:2204.05862.

[3] Bommasani, R.; Hudson, D.A.; Adeli, E.; Altman, R.; Arora, S.; von Arx, S.; Bernstein, M.S.; Bohg, J.; Bosselut, A.; Brunskill, E.; et al. 2021. On the Opportunities and Risks of Foundation Models. arXiv:2108.07258.

[4] Carlini, N.; Wagner, D. 2017. Towards Evaluating the Robustness of Neural Networks. In IEEE Symposium on Security and Privacy, 39–57.

[5] Goodfellow, I.J.; Shlens, J.; Szegedy, C. 2015. Explaining and Harnessing Adversarial Examples. In ICLR.

[6] Hendrycks, D.; Burns, C.; Basart, S.; Zou, A.; Mazeika, M.; Song, D.; Steinhardt, J. 2021. Measuring Massive Multitask Language Understanding. In ICLR.

[7] Kenton, Z.; Everitt, T.; Weidinger, L.; Gabriel, I.; Mikulik, V.; Irving, G. 2021. Alignment of Language Agents. arXiv:2103.14659.

[8] Leike, J.; Krueger, D.; Everitt, T.; Martic, M.; Maini, V.; Legg, S. 2018. Scalable Agent Alignment via Reward Modeling: A Research Direction. arXiv:1811.07871.

[9] OpenAI. 2023. GPT-4 Technical Report. arXiv:2303.08774.

[10] Ouyang, L.; Wu, J.; Jiang, X.; Almeida, D.; Wainwright, C.; Mishkin, P.; Zhang, C.; Agarwal, S.; Slama, K.; Ray, A.; et al. 2022. Training Language Models to Follow Instructions with Human Feedback. In NeurIPS.

[11] Perez, E.; Huang, S.; Song, F.; Cai, T.; Ring, R.; Aslanides, J.; Glaese, A.; McAleese, N.; Irving, G. 2022. Red Teaming Language Models with Language Models. In EMNLP, 3419–3448.

[12] Rafailov, R.; Sharma, A.; Mitchell, E.; Ermon, S.; Manning, C.D.; Finn, C. 2023. Direct Preference Optimization: Your Language Model is Secretly a Reward Model. arXiv:2305.18290.

[13] Sandoval-Romero, M.; Rothchild, D.; Chen, J. 2023. Fairness Testing of Machine Learning Models. In ACM FAccT, 156–167.

[14] Touvron, H.; Martin, L.; Stone, K.; Albert, P.; Almahairi, A.; Babaei, Y.; Bashlykov, N.; Batra, S.; Bhargava, P.; Bhosale, S.; et al. 2023. Llama 2: Open Foundation and Fine-Tuned Chat Models. arXiv:2307.09288.

[15] Wallace, E.; Feng, S.; Kandpal, N.; Gardner, M.; Singh, S. 2019. Universal Adversarial Triggers for Attacking and Analyzing NLP. In EMNLP-IJCNLP, 2153–2162.

[16] Wei, J.; Wang, X.; Schuurmans, D.; Bosma, M.; Ichter, B.; Xia, F.; Chi, E.; Le, Q.; Zhou, D. 2022. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. In NeurIPS.

[17] Weidinger, L.; Mellor, J.; Rauh, M.; Griffin, C.; Uesato, J.; Huang, P.S.; Cheng, M.; Glaese, M.; Balle, B.; Kasirzadeh, A.; et al. 2021. Ethical and Social Risks of Harm from Language Models. arXiv:2112.04359.

[18] Zou, A.; Wang, Z.; Kolter, J.Z.; Fredrikson, M. 2023. Universal and Transferable Adversarial Attacks on Aligned Language Models. arXiv:2307.15043.