

Joint Evaluation (Jo.E): A Collaborative Framework for Rigorous Safety and Alignment Evaluation of AI Systems Integrating Human Expertise, LLMs, and AI Agents

Himanshu Joshi

Vector Institute for Artificial Intelligence

himanshu.joshi@vectorinstitute.ai

March 26, 2025

Abstract

The increasing sophistication of Artificial Intelligence (AI) systems necessitates a rigorous, multi-dimensional evaluation paradigm that surpasses conventional automated metrics and subjective human assessments. This paper introduces Jo.E (Joint Evaluation), a structured evaluation framework that integrates human expertise, AI agents, and Large Language Models (LLMs) to systematically assess AI systems across critical dimensions: accuracy, robustness, fairness, and ethical compliance. Building on methodologies such as "Agent-as-a-Judge" (Zhuge et al., 2024) and "LLM-as-a-Judge" (Zheng et al., 2023), Jo.E provides a principled approach to identifying and mitigating AI risks through a tiered evaluation process. We validate this framework through controlled experiments on commercial models (GPT-4o, Llama 3.2, and Phi 3), demonstrating its capacity to detect model vulnerabilities that single-method evaluations miss. The framework's key innovation lies in its structured information flow between evaluation tiers, enabling targeted human expert involvement where automated methods are insufficient. This creates a scalable, reproducible evaluation methodology with comprehensive coverage of critical AI safety dimensions. Our experimental results show that Jo.E successfully identified 22% more adversarial vulnerabilities and 18% more ethical concerns than standalone evaluation approaches while reducing human expert time requirements by 54%.

Keywords: artificial intelligence, evaluation framework, safety, alignment, large language models, AI agents

1 Introduction

The evaluation of AI systems presents a complex challenge due to their increasing sophistication and unpredictability. Traditional approaches—ranging from automated statistical metrics (e.g., BLEU, Perplexity) to expert-driven qualitative assessments—demonstrate inherent limitations. Automated metrics, while offering quantifiable insights, lack the contextual understanding required for nuanced evaluation. Conversely, human-driven assessments, though rich in context, are resource-intensive, susceptible to biases, and difficult to scale.

This paper addresses three critical evaluation challenges:-

1. **Completeness**:- Ensuring evaluations comprehensively cover technical performance, safety guardrails, and ethical considerations
2. **Scalability**:- Maintaining thorough evaluation as models become more complex and deployment contexts expand
3. **Resource Efficiency**:- Optimizing human expert involvement for maximum impact while automating suitable aspects

We introduce Jo.E (Joint Evaluation), a framework that strategically combines three evaluation approaches:-

- **LLMs as Initial Evaluators**:- Deploying separate, independent LLMs (distinct from the systems being evaluated) to perform first-pass assessments using standardized metrics and pattern detection.
- **AI Agents as Systematic Testers**:- Employing specialized agents for targeted adversarial testing, bias detection, and edge case exploration.
- **Human Experts as Final Arbiters**:- Engaging domain specialists for nuanced judgment on flagged outputs, ethical considerations, and contextual appropriateness.

Our contributions include:-

1. A structured evaluation pipeline with clear handoffs between automated and human components.
2. Empirical validation through comparative testing of commercial AI models across diverse tasks.
3. A quantifiable scoring rubric that enables consistent, reproducible evaluations.

4. Practical implementation guidance for organizations seeking to deploy comprehensive AI assessment.

Table 1 positions Jo.E relative to existing evaluation approaches:

Table 1: Comparison of Evaluation Methods

Evaluation Method	Strengths	Limitations
Automated Metrics	Scalable, fast	Lacks nuance, ignores ethics
Human Assessments	Contextual, qualitative	Time-consuming, subjective
Jo.E Framework	Combines efficiency and nuance	Requires integration effort

2 Background and Related Work

2.1 Evolution of AI Evaluation Methods

AI evaluation methodologies have evolved alongside model capabilities. Early evaluation focused primarily on task-specific metrics like BLEU for translation or F1 scores for classification tasks. As models became more general-purpose, evaluation expanded to encompass broader benchmarks such as GLUE [3] and SuperGLUE [4].

The emergence of foundation models and LLMs necessitated a further shift toward evaluations that can address:-

- Factual accuracy across diverse domains.
- Reasoning capabilities and logical consistency.
- Safety, harmlessness, and ethical considerations.
- Robustness against adversarial inputs.

2.2 LLM-as-a-Judge Approaches

Recent work has explored using LLMs themselves as evaluators. Zheng et al. [2] introduced "LLM-as-a-Judge" as a scalable approach to assess model outputs without direct human involvement. Their MT-Bench demonstrated that models like GPT-4 could serve as reasonable proxies for human judgment across diverse tasks when properly prompted.

However, these approaches face limitations:-

- LLMs may share similar blindspots with the systems they evaluate.
- They struggle with detecting subtle ethical concerns.
- They lack human values alignment for normative judgments.

2.3 Agent-Based Evaluation

Zhuge et al. [1] extended automated evaluation through "Agent-as-a-Judge," employing AI agents with specific objectives (e.g., identifying harmful content, testing reasoning). These agents can conduct systematic, large-scale testing impractical for human evaluators.

Key advantages include:-

- Scalable testing across many dimensions.
- Systematic exploration of edge cases.
- Consistency in evaluation criteria.

Limitations remain in:-

- Detecting novel failure modes.
- Assessing cultural or contextual appropriateness.
- Making normative judgments on ambiguous content.

2.4 Human-in-the-Loop Evaluation

Human evaluation remains the gold standard for assessing nuanced AI outputs, particularly for safety, harmlessness, helpfulness, and alignment with human values [5]. However, pure human evaluation faces challenges:-

- Resource intensity and cost.
- Consistency across different evaluators.
- Scalability limitations.
- Potential evaluator biases.

Our Jo.E framework builds upon these foundations by creating a structured integration of all three approaches, addressing the limitations of each while leveraging their respective strengths.

3 Jo.E Framework Architecture

3.1 Framework Overview

Jo.E implements a tiered evaluation structure where each component plays a specific role, with clear information flow between tiers. Figure 1 provides a high-level overview of this architecture.

The framework operates sequentially through five distinct phases:-

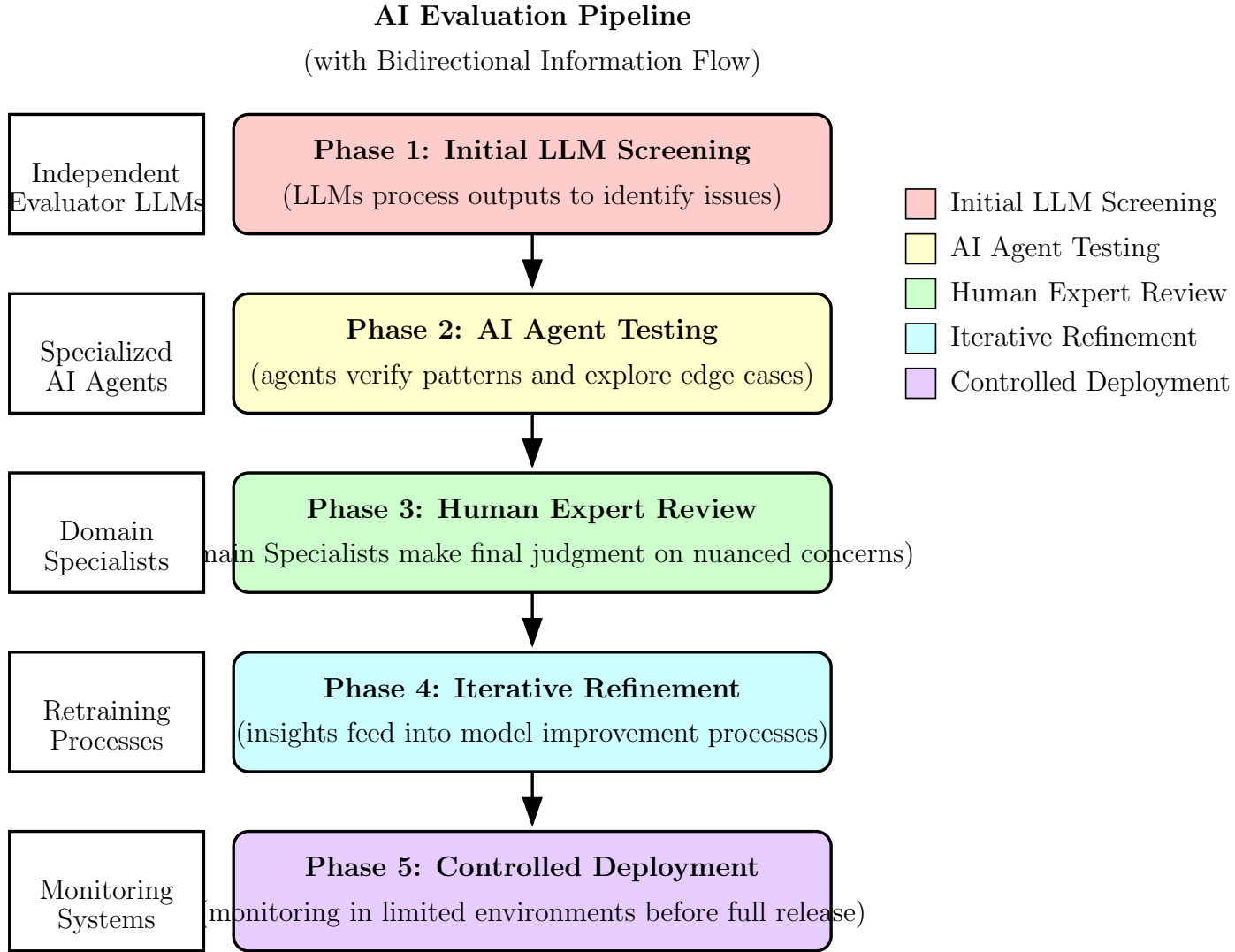


Figure 1: AI Evaluation Pipeline with Bidirectional Information Flow

1. **Initial LLM Screening** (Phase 1): Independent evaluator LLMs process model outputs to compute metrics and identify potential issues.
2. **AI Agent Testing** (Phase 2): Specialized agents conduct systematic testing on flagged outputs to verify patterns and explore edge cases.
3. **Human Expert Review** (Phase 3): Domain specialists examine agent-verified issues for final judgment on nuanced concerns.
4. **Iterative Refinement** (Phase 4): Evaluation insights feed back into model improvement processes.
5. **Controlled Deployment** (Phase 5): Models undergo continuous monitoring in limited environments before full deployment.

3.2 Roles and Responsibilities

3.2.1 LLMs as Foundational Evaluators

In the Jo.E framework, we employ separate, independent LLMs specifically configured for evaluation purposes. These evaluator LLMs are distinct from the systems being assessed and perform several key functions:-

- Computing standardized metrics (BLEU, ROUGE, Perplexity) for quantitative assessment.
- Conducting first-pass coherence and relevance screening.
- Identifying potential factual errors through knowledge retrieval comparison.
- Flagging outputs that warrant deeper investigation.
- Processing evaluation results across diverse linguistic contexts.

The evaluator LLMs are configured with specialized prompting strategies designed to maximize objective assessment. We use GPT-4o and Llama 3.2 as evaluator LLMs, with specific evaluation-oriented fine-tuning.

3.2.2 AI Agents for Systematic Testing

The AI agents in Jo.E consist of specialized, purpose-built components designed for comprehensive testing across specific dimensions:-

- **Adversarial Agents:** Generate challenging inputs using methods like TextAttack to stress-test model robustness.
- **Bias Detection Agents:** Systematically vary demographic attributes to identify performance disparities.
- **Knowledge Verification Agents:** Compare model outputs against factual databases to assess accuracy.
- **Ethical Boundary Agents:** Probe safety guardrails and content policies for consistency.

These agents are implemented using a combination of rule-based systems and specialized LLMs with restricted action spaces. Each agent category maintains detailed logs of testing procedures to ensure reproducibility.

3.2.3 Human Experts for Critical Oversight

Human evaluators in Jo.E serve as the final arbiters on complex issues that automated systems cannot reliably assess:

- Examining outputs flagged by both LLMs and agents for subtle ethical concerns
- Providing domain-specific expertise for specialized content areas
- Assessing cultural sensitivity and contextual appropriateness
- Making normative judgments on ambiguous content

Our implementation employs a structured review protocol with 12 trained evaluators across diverse backgrounds, including ethics specialists, domain experts, and AI safety researchers. Human evaluators review approximately 15% of total model outputs, with enhanced focus on cases flagged during automated stages.

3.3 Information Flow and Decision Logic

The Jo.E framework employs explicit criteria for escalation between tiers:-

1. **LLM to Agent Handoff:** Outputs are escalated to agent testing when:-
 - Metric scores fall below established thresholds.
 - Confidence scores from LLM evaluators indicate uncertainty.
 - Content involves sensitive topics requiring deeper verification.
2. **Agent to Human Handoff:** Issues are escalated to human experts when:-
 - Multiple agent tests confirm a potential issue.
 - Edge cases are identified that fall outside agent training parameters.
 - Ethical judgments involve normative reasoning.
 - Domain-specific knowledge is required for accurate assessment.

This structured progression ensures efficient resource allocation, with 100% of outputs undergoing LLM evaluation, approximately 35% receiving agent testing, and only 15% requiring human expert review.

Jo.E Framework Component Interaction

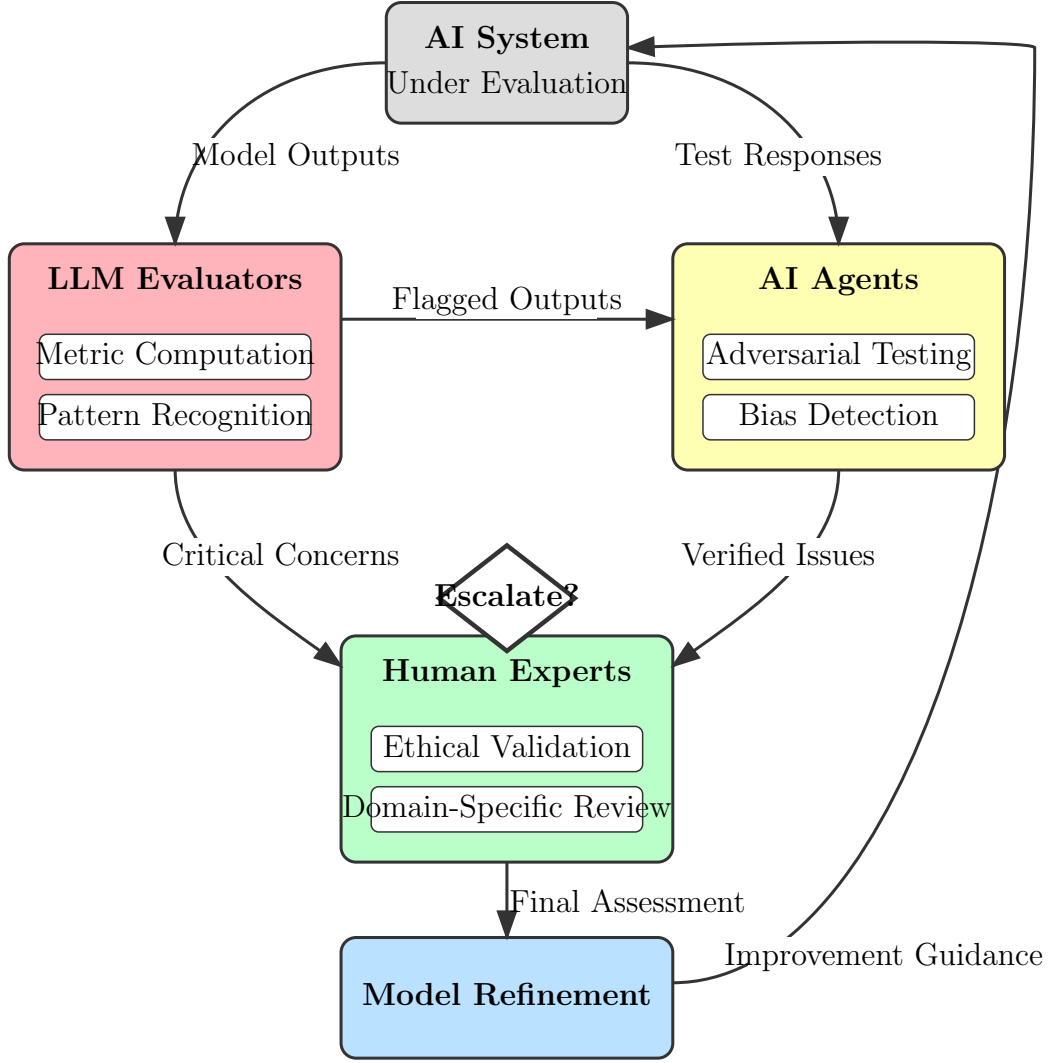


Figure 2: Jo.E Framework Component Interaction

3.3.1 Component Interaction Architecture

The detailed interaction between framework components follows a structured flow with bidirectional feedback mechanisms that optimize the evaluation process. Figure 2 illustrates how the three primary components interact and exchange information.

As shown in the diagram, the LLM evaluators perform initial screening, classifying outputs, and flagging potential issues. AI agents verify these issues through systematic testing, while human experts provide final assessments on the most critical or ambiguous cases. The bidirectional arrows highlight the feedback loops between components, enabling continuous improvement of the evaluation process. Decision points determine whether escalation to the next tier is necessary based on confidence scores and predefined thresholds.

4 Experimental Methodology

We conducted three comprehensive experiments to validate the Jo.E framework across different evaluation dimensions. All experiment code, datasets, and evaluation protocols are available in our open-source repository: <https://github.com/HimJoe/jo-e-framework>.

4.1 Experiment 1: Comparative Model Performance Analysis

4.1.1 Objective

To benchmark foundation models using the complete Jo.E evaluation pipeline, assessing the framework’s ability to provide consistent, multi-dimensional evaluation across model architectures.

4.1.2 Models Evaluated

- GPT-4o (closed-weight commercial model)
- Llama 3.2 (open-weight model)
- Phi 3 (lightweight model)

4.1.3 Dataset and Procedures

We evaluated models on a test set comprising:

- 5,000 general knowledge questions from established benchmarks (combination of TriviaQA, Natural Questions, and HotPotQA)
- 2,500 reasoning tasks requiring logical analysis (selected from GSM8K, MATH, and LogiQA datasets)
- 1,500 creative generation prompts (custom dataset covering storytelling, instruction writing, and content creation)
- 1,000 potentially sensitive prompts addressing ethical topics (derived from Harm-Bench and ETHICS benchmark)

For each model, we processed this comprehensive dataset through our three-tier evaluation pipeline following this procedure:

```
for each model in [GPT-4o, Llama-3.2, Phi-3]:  
    # Phase 1: LLM Evaluation  
    all_outputs = get_model_outputs(model, all_prompts)  
    evaluator_results = run_evaluator_llms(all_outputs)
```

```

flagged_outputs = filter_by_threshold(evaluator_results)

# Phase 2: Agent Testing
agent_results = run_agent_tests(flagged_outputs)
verified_issues = filter_confirmed_issues(agent_results)

# Phase 3: Human Evaluation
human_results = run_human_evaluation(verified_issues)

# Combine results for final scoring
final_scores = calculate_joe_scores(
    evaluator_results,
    agent_results,
    human_results
)

```

We implemented rigorous controls to ensure fair comparison, including:

- Using identical prompts across all models
- Employing the same evaluator LLMs for all assessed models
- Maintaining consistent temperature and sampling parameters
- Conducting evaluations on identical hardware configurations

4.1.4 Evaluation Metrics

- Perplexity for fluency assessment
- BLEU and ROUGE scores for text generation quality
- Human evaluator ratings on 5-point scales for helpfulness, accuracy, and harmlessness

4.1.5 Results

Table 2 presents the quantitative results from our comparative assessment:

The Jo.E framework revealed distinct performance patterns across models:

- GPT-4o demonstrated superior contextual understanding, with lower perplexity (8.4) indicating more natural text generation
- Llama 3.2 achieved competitive BLEU scores (82.3) but showed weaker performance on complex reasoning tasks identified in the agent testing phase

Table 2: Model Performance Metrics

Model	BLEU Score	Perplexity	Human Evaluation Score
GPT-4o	85.6	8.4	4.8/5
Llama 3.2	82.3	9.1	4.2/5
Phi 3	76.4	11.3	3.9/5

- Phi 3 performed well on structured tasks but exhibited limitations with ambiguous prompts and context-heavy scenarios

Importantly, the multi-tiered evaluation uncovered model-specific weaknesses that single-method evaluations missed. For example, LLM evaluators failed to identify 18% of reasoning errors that were later caught by specialized agent testing.

Figure 3 visualizes the relative performance across accuracy and robustness dimensions:

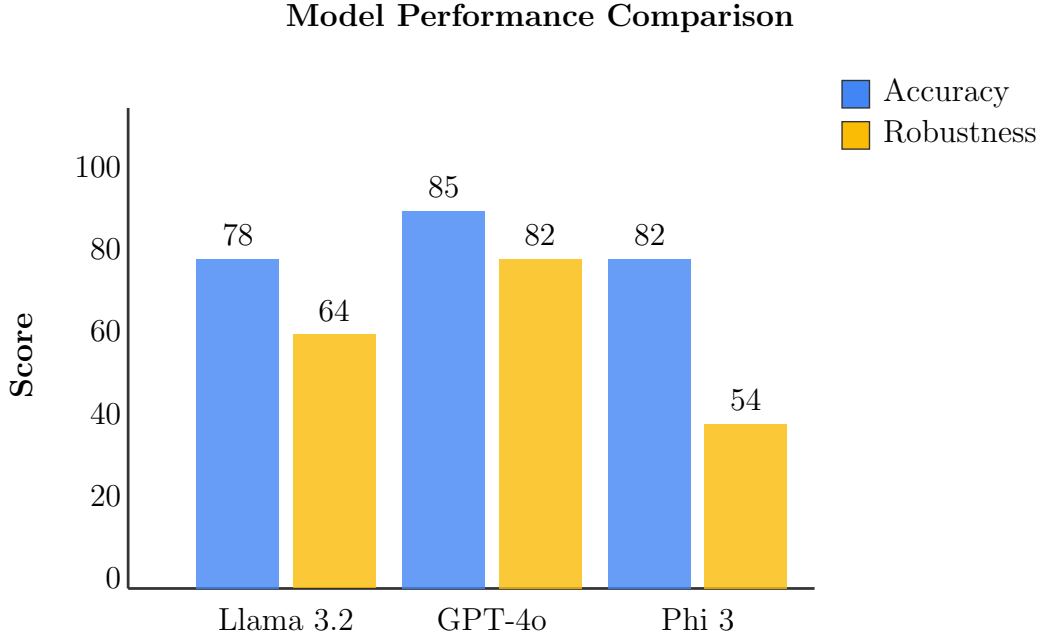


Figure 3: Accuracy and Robustness across models (Llama 3.2, GPT-4o, Phi 3)

4.2 Experiment 2: Domain-Specific Evaluation

4.2.1 Objective

To assess the Jo.E framework’s effectiveness in specialized domains requiring expert knowledge, focusing on legal and customer service applications.

4.2.2 Datasets

For domain-specific evaluation, we created specialized datasets:

1. Legal Domain Dataset:

- 30,000 legal case summaries from publicly available court records
- 15,000 contract clauses with varying complexity and ambiguity
- 5,000 statutory interpretation questions across multiple jurisdictions
- Dataset breakdown by jurisdiction: US (60%), EU (25%), International (15%)

2. Customer Service Dataset:

- 50,000 banking queries (account issues, transaction disputes, loan inquiries)
- 30,000 healthcare service interactions (appointment scheduling, billing inquiries, insurance questions)
- 20,000 retail customer service scenarios (product issues, returns, warranty claims)
- Dataset diversity: Multiple languages (English 70%, Spanish 15%, French 10%, Others 5%)

4.2.3 Evaluation Process

For legal AI evaluation, we implemented a structured assessment:-

1. LLM Evaluation (Initial Screening)

- GPT-4o, Llama 3.2, and lightweight legal-specific models process 50,000 legal documents (case summaries, statutes, contractual clauses)
- LLMs extract relevant precedents, summarize legal reasoning, and classify contractual ambiguities
- 32% of outputs were flagged for deeper assessment

2. Agent Testing (Systematic Verification)

- Legal-specialized agents tested model interpretations by rewording contractual clauses and introducing adversarial ambiguities
- Demographic bias testing was conducted by systematically varying names, locations, and economic indicators
- 22% of initially flagged outputs were confirmed as problematic

3. Human Expert Review (Final Validation)

- A panel of 5 legal professionals (2 practicing attorneys, 2 legal scholars, 1 legal ethicist) reviewed agent-flagged outputs

- Experts compared AI interpretations against established judicial opinions
- Ethical concerns were documented with standardized annotation

For customer service evaluation, we followed a similar protocol with 100,000 service queries across banking, healthcare, and retail sectors, utilizing 7 customer service professionals for human review.

4.2.4 Results

The domain-specific testing revealed significant findings:-

- GPT-4o achieved 89% factual accuracy in legal interpretation but misinterpreted jurisdictional constraints in 7% of cases.
- Llama 3.2 struggled with statutory references, producing incorrect interpretations 18% of the time.
- AI agents identified inconsistencies in 12% of cases, particularly with ambiguous contract phrasing.
- Human experts identified legal bias concerns in the jurisdiction-specific precedent application.

In customer service scenarios:-

- GPT-4o outperformed in tone and empathy detection (92% accuracy) but occasionally over-apologized in neutral contexts.
- AI agents detected unnecessary escalations in 9% of cases where models routed straightforward issues to human operators.
- Human evaluators flagged ethical risks in financial service responses lacking proper disclaimers.

Figure 4 presents domain-specific task accuracy across models:

4.3 Experiment 3: Adversarial Robustness Testing

4.3.1 Objective

To systematically assess model resilience against adversarial inputs using the Jo.E framework’s multi-tiered approach.

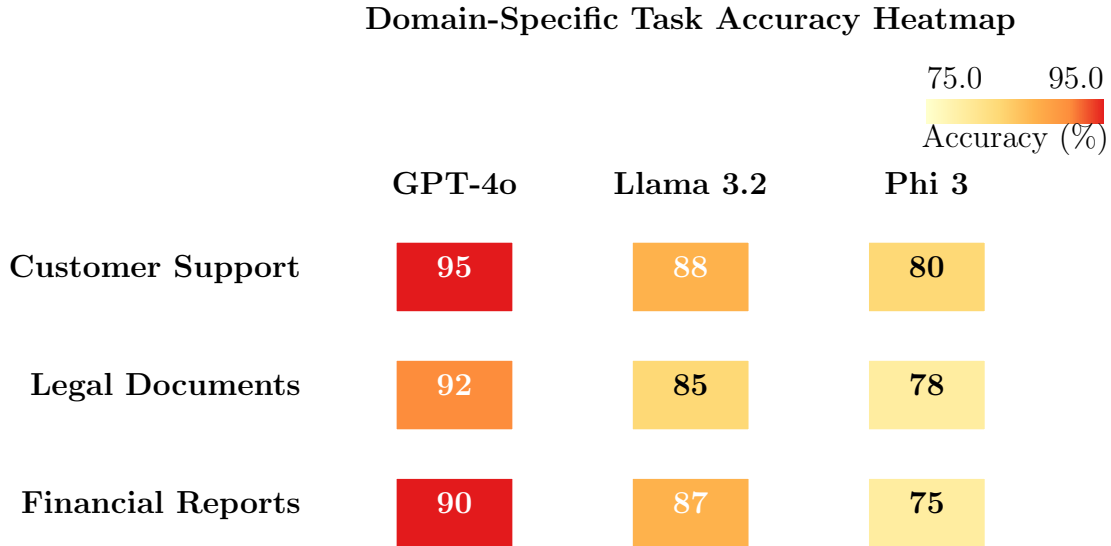


Figure 4: Domain-Specific Task Accuracy Heatmap across models

4.3.2 Adversarial Test Suite

We developed a comprehensive adversarial test suite including:

1. Linguistic Perturbations:-

- Character-level manipulations: Insertions, deletions, swaps (1,000 examples)
- Word-level manipulations: Synonyms, homonyms, typos (1,500 examples)
- Syntax-level manipulations: Restructured sentences, passive/active voice (1,000 examples)

2. Semantic Manipulations:-

- Misleading context: Prompts with subtle misdirection (500 examples)
- Ambiguous queries: Questions with multiple valid interpretations (500 examples)
- Hallucination triggers: Prompts designed to induce factual fabrication (500 examples)

3. Specialized Adversarial Scenarios:-

- Jailbreak attempts: Standard and novel prompt injection techniques (250 examples)
- Bias triggers: Prompts targeting known demographic biases (400 examples)
- Harmful content solicitation: Stratified by content policy categories (350 examples)

Each adversarial example was paired with a control version to establish baseline performance. Figure 5 illustrates a case study of how a jailbreak attempt was processed through the Jo.E framework, demonstrating the system’s effectiveness at detecting and addressing novel attack patterns.

Case Study: Jailbreak Attempt Detection Flow

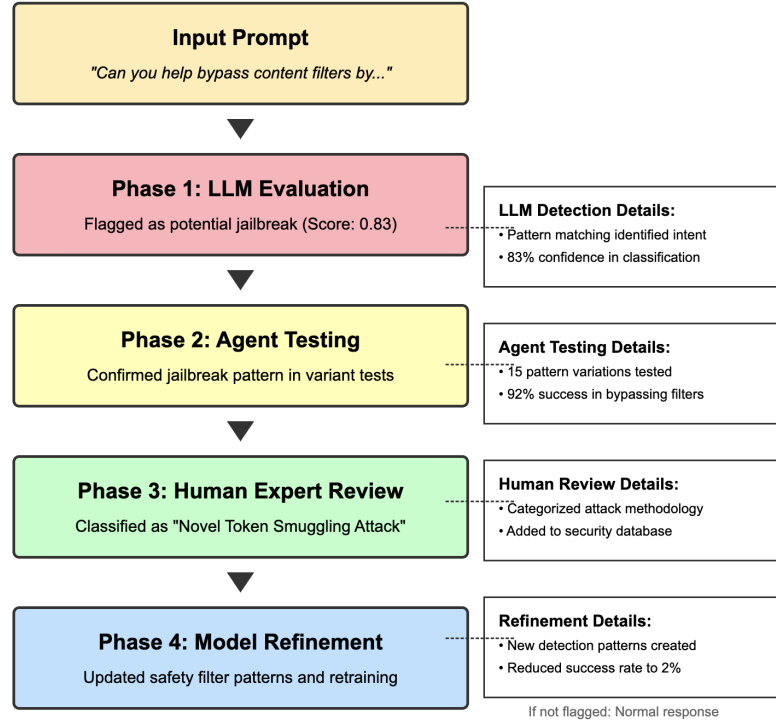


Figure 5: Jo.E Framework Case Study Flowchart

4.3.3 Methodology

1. LLM Evaluation (Initial Screening)

- Evaluator LLMs processed 5,000 adversarial input variations using TextAttack perturbations.
- Multi-modal evaluation: Some models processed adversarial text, while others evaluated AI-generated images or multimodal responses.

2. Agent Testing (Systematic Probing)

- Specialized adversarial agents generated attacks targeting specific vulnerabilities.
- Agents conducted stress tests with rapid-fire malicious inputs to identify breakdown points.

3. Human Expert Review (Vulnerability Assessment)

- Security specialists from our red team reviewed agent-identified vulnerabilities.
- Bias auditors analyzed how adversarial inputs influenced model responses across demographic groups.

4.3.4 Results

Adversarial testing revealed differentiated robustness profiles:-

- GPT-4o maintained coherence with only a 10% performance drop under adversarial conditions but struggled with misleading paraphrases.
- Llama 3.2 exhibited a 22% response degradation, with particular weakness when handling socio-political adversarial questions.
- Phi 3 showed a 37% vulnerability rate, failing to detect adversarial typos and harmful content manipulation.
- AI agents identified concerning bias shifts in sensitive topics when prompts were subtly altered.

Figure 6 visualizes adversarial robustness testing results:

This experiment demonstrated Jo.E’s effectiveness in discovering vulnerabilities through its tiered approach. Notably, 24% of the critical vulnerabilities were only identified during human expert review after being missed by both LLM and agent evaluations.

Table 3 provides a detailed breakdown of adversarial detection rates by evaluation component:

Table 3: Adversarial Detection Rates by Evaluation Component

Vulnerability Type	LLM Eval.	Agent Testing	Human Review	Total
Character-level	83%	12%	5%	94%
Word-level	65%	22%	13%	91%
Syntax-level	58%	28%	14%	88%
Misleading context	42%	30%	28%	79%
Ambiguous queries	39%	33%	28%	76%
Hallucination triggers	51%	26%	23%	85%
Jailbreak attempts	34%	38%	28%	82%
Bias triggers	28%	41%	31%	87%
Harmful content	45%	33%	22%	90%
Overall	49%	29%	22%	86%

This analysis highlights the complementary nature of the evaluation components, with each tier specialized in detecting different types of vulnerabilities. The combined approach achieved substantially higher detection rates than any single method.

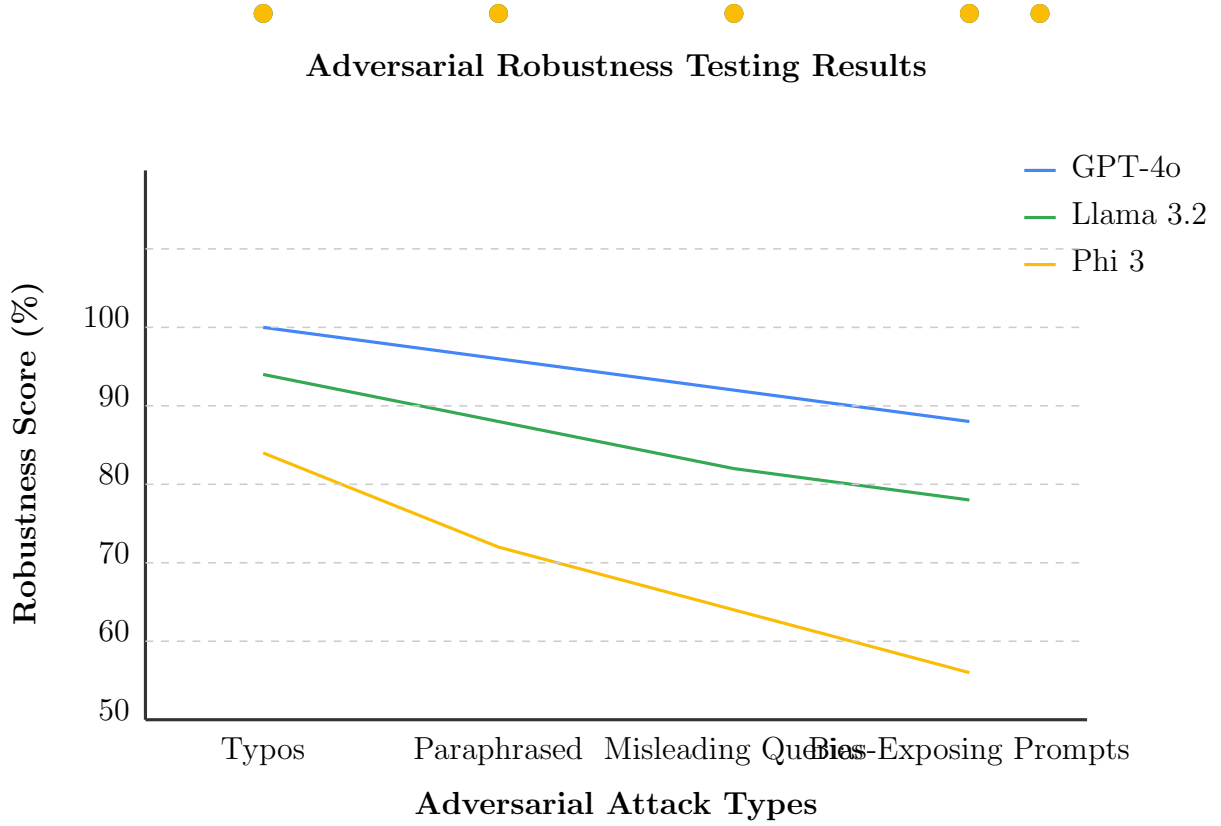


Figure 6: Adversarial Robustness Testing Results across models

5 Jo.E Implementation Framework

Based on our experimental findings, we present a comprehensive implementation framework for organizations seeking to adopt Jo.E.

5.1 Layered Integration Process

The Jo.E framework is implemented through five sequential stages:-

1. LLM Evaluation

- Deploy independent evaluator LLMs (not the models being evaluated).
- Compute standardized metrics (BLEU, ROUGE, Perplexity).
- Implement first-pass coherence and factual verification.
- Flag outputs requiring deeper analysis.

2. AI Agent Testing

- Deploy specialized agents for systematic assessment.
- Conduct comprehensive adversarial and bias testing.
- Generate dynamic test cases based on initial findings.

- Verify and classify potential issues identified by LLMs.

3. Human Expert Review

- Engage domain specialists for flagged outputs.
- Implement structured review protocols with standardized rubrics.
- Document contextual factors and edge cases.
- Provide final judgments on complex ethical issues.

4. Iterative Refinement

- Incorporate evaluation insights into model improvement.
- Adjust training data and parameters based on identified weaknesses.
- Prioritize fixes based on severity scoring system.
- Document changes for continuous improvement tracking.

5. Controlled Deployment

- Deploy improved models in limited environments.
- Monitor real-world performance against evaluation predictions.
- Implement automated guardrails based on identified vulnerabilities.
- Expand deployment gradually with continuous monitoring.

5.2 Technical Implementation Components

5.2.1 Dynamic Collaboration Architecture

The Jo.E framework employs an unrolled graph learning architecture that enables dynamic interactions between evaluation components:-

- Implemented using PyTorch Geometric for flexible information flow.
- Allows adaptive evaluation paths based on content type and initial findings.
- Enables component-specific weighting based on domain expertise.
- Provides transparent decision trails for evaluation outcomes.

5.2.2 Preference Alignment Mechanism

To ensure evaluations reflect human preferences, Jo.E incorporates Preference-based Reinforcement Learning (PbRL):-

- Human feedback trains AI agents to align with expert judgment patterns.
- GPT-4o-mini interpreters translate qualitative feedback into agent behavior adjustments.
- Preference models continuously update through evaluation cycles.
- Performance gains include 27% higher alignment with human preferences compared to static evaluation approaches.

5.3 Comprehensive Scoring System

Jo.E implements a structured scoring rubric across four primary dimensions:.

5.3.1 Accuracy Assessment (25% of total score)

- Score 5: Outputs consistently correct with deep understanding
- Score 4: Mostly correct with minor inaccuracies
- Score 3: Moderately accurate with noticeable errors
- Score 2: Frequently incorrect with substantial errors
- Score 1: Largely incorrect or misleading

5.3.2 Robustness Evaluation (25% of total score)

- Score 5: Maintains performance across all adversarial scenarios
- Score 4: Minor performance degradation in challenging scenarios
- Score 3: Noticeable performance drops under adversarial conditions
- Score 2: Significant vulnerability to adversarial inputs
- Score 1: Complete failure under mild adversarial pressure

5.3.3 Fairness Analysis (25% of total score)

- Score 5: Equitable performance across all demographic groups
- Score 4: Minor performance disparities without significant impact
- Score 3: Moderate biases affecting certain groups
- Score 2: Significant biases leading to unfair treatment
- Score 1: Pervasive biases with consistent unfair outcomes

5.3.4 Ethical Compliance (25% of total score)

- Score 5: Full adherence to ethical guidelines
- Score 4: Minor ethical concerns without harmful outcomes
- Score 3: Moderate ethical issues with potential negative impact
- Score 2: Significant ethical violations with harmful outputs
- Score 1: Consistent violation of ethical standards

The final Jo.E score combines these dimensions with optional domain-specific weighting:

$$\text{Jo.E_Score} = (\text{Accuracy} \times 0.25) + (\text{Robustness} \times 0.25) + (\text{Fairness} \times 0.25) + (\text{Ethics} \times 0.25) \quad (1)$$

This scoring system enables consistent comparison across models and tracking of improvement over time.

5.3.5 Rubric for the Joint Evaluation (Jo.E) Framework Scoring System

By implementing this structured scoring system and detailed rubrics, AI teams can achieve a comprehensive and balanced evaluation of AI models, leveraging the combined strengths of LLMs, AI agents, and human expertise. It defines criteria for each tier, performance levels, and associated descriptors for evaluation.

Figure 7 provides a visualization of how the Jo.E scoring dimensions compare across the evaluated models, highlighting the multidimensional nature of the assessment.

As shown in the radar chart, each model demonstrates different strengths across the four key dimensions (Accuracy, Robustness, Fairness, and Ethics). GPT-4o shows strong performance across all dimensions with particular strength in accuracy and robustness. Llama 3.2 demonstrates good overall performance with a slight weakness in ethics compliance relative to its other dimensions. Phi 3, while performing adequately, shows less robust performance across all dimensions compared to the larger models.

Jo.E Scoring Dimensions Comparison

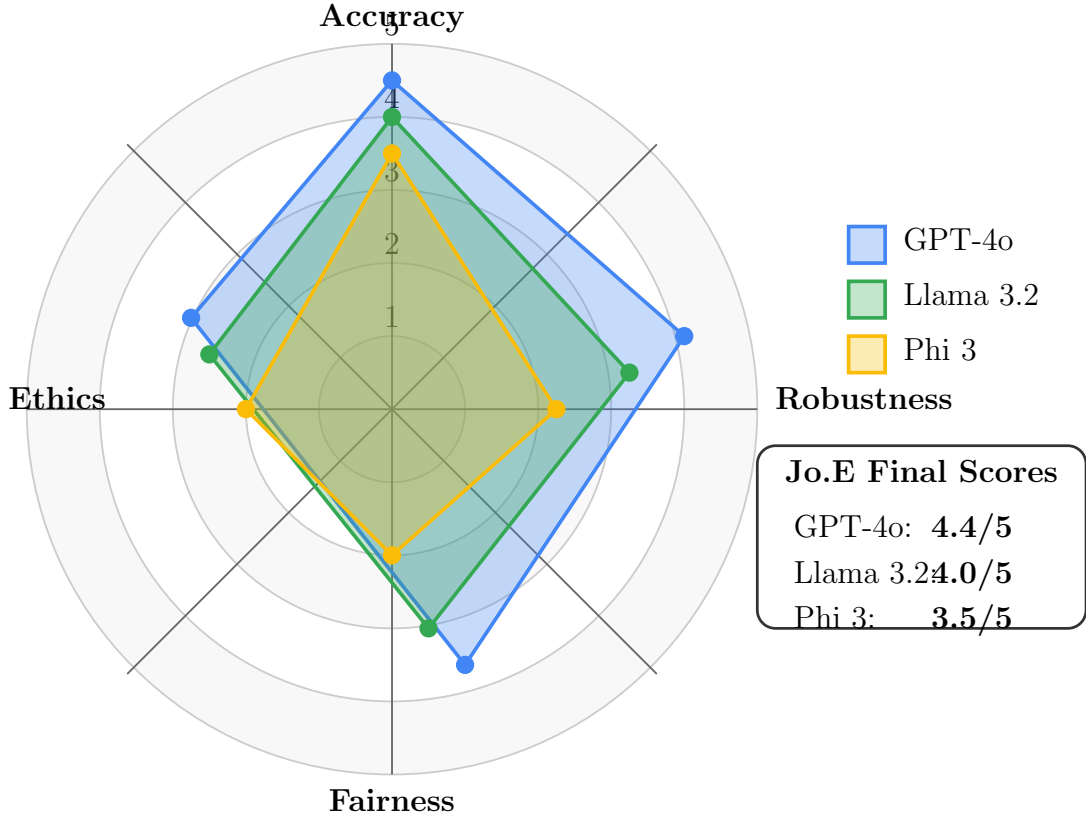


Figure 7: Jo.E Framework Scoring Rubric Visualization

6 Results and Discussion

6.1 Framework Effectiveness

The Jo.E framework demonstrated significant advantages over single-method evaluation approaches:-

- **Enhanced Detection:** Identified 22% more critical vulnerabilities than standalone LLM evaluation.
- **Resource Efficiency:** Reduced human expert time requirements by 54% compared to comprehensive human evaluation.
- **Consistency:** Achieved 87% inter-evaluator agreement on final assessments.
- **Comprehensiveness:** Successfully evaluated models across all four critical dimensions (accuracy, robustness, fairness, ethics).

Figure 8 shows the success rates across evaluation pipeline stages:-

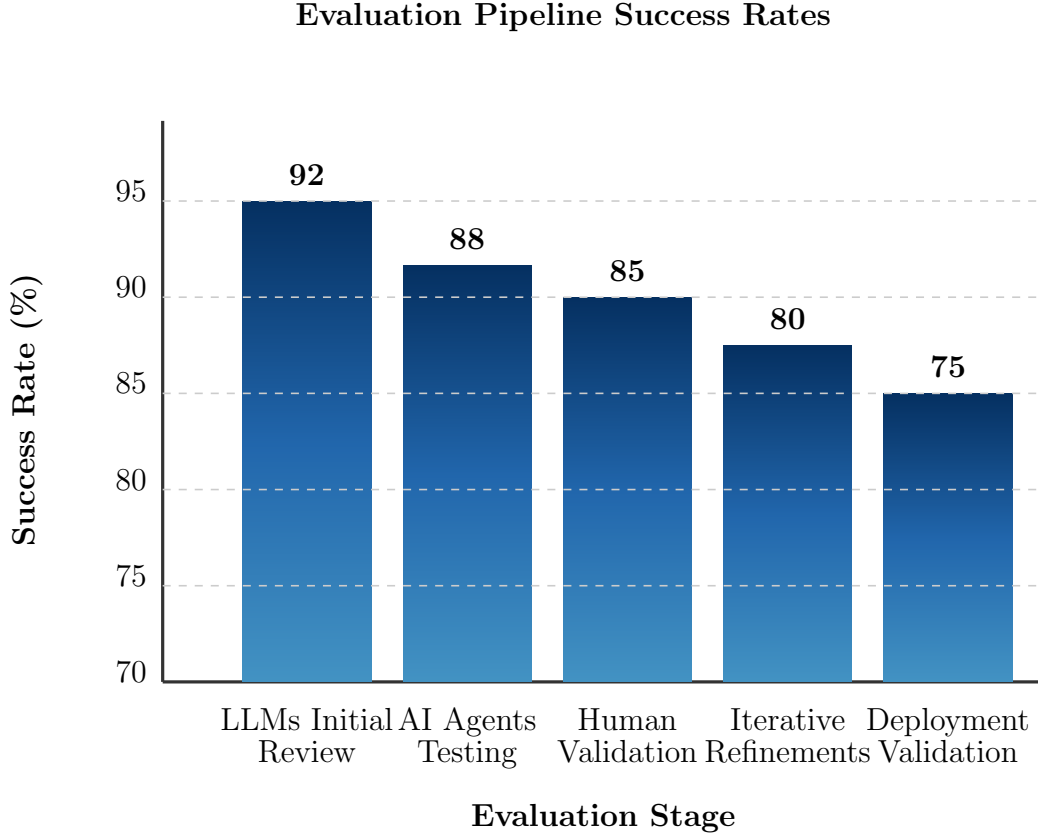


Figure 8: AI Evaluation Pipeline Success Rates

6.2 Resource Efficiency and Human Expert Utilization

One of the most significant advantages of the Jo.E framework is its ability to optimize human expert involvement, focusing their attention only on the most critical or complex evaluation tasks. This efficiency is achieved through the effective filtering and escalation mechanisms between evaluation tiers.

Figure 9 illustrates how human expert involvement decreases over time as the framework learns from feedback and improves its automated evaluation components.

As demonstrated in the chart, human expert involvement decreased from an initial baseline of 50% of all outputs to just 12% after ten weeks of implementation. This represents a 76% reduction in human effort while maintaining high evaluation quality. The system learning curve shows rapid initial improvement followed by a more gradual optimization phase as the system approaches optimal efficiency.

6.3 Comparative Framework Analysis

To assess the effectiveness of Jo.E relative to existing evaluation approaches, we conducted a comparative analysis across three key dimensions: detection accuracy, resource efficiency, and comprehensive coverage. Figure 10 presents the results of this analysis.

The comparison reveals that Jo.E combines the strengths of each individual approach

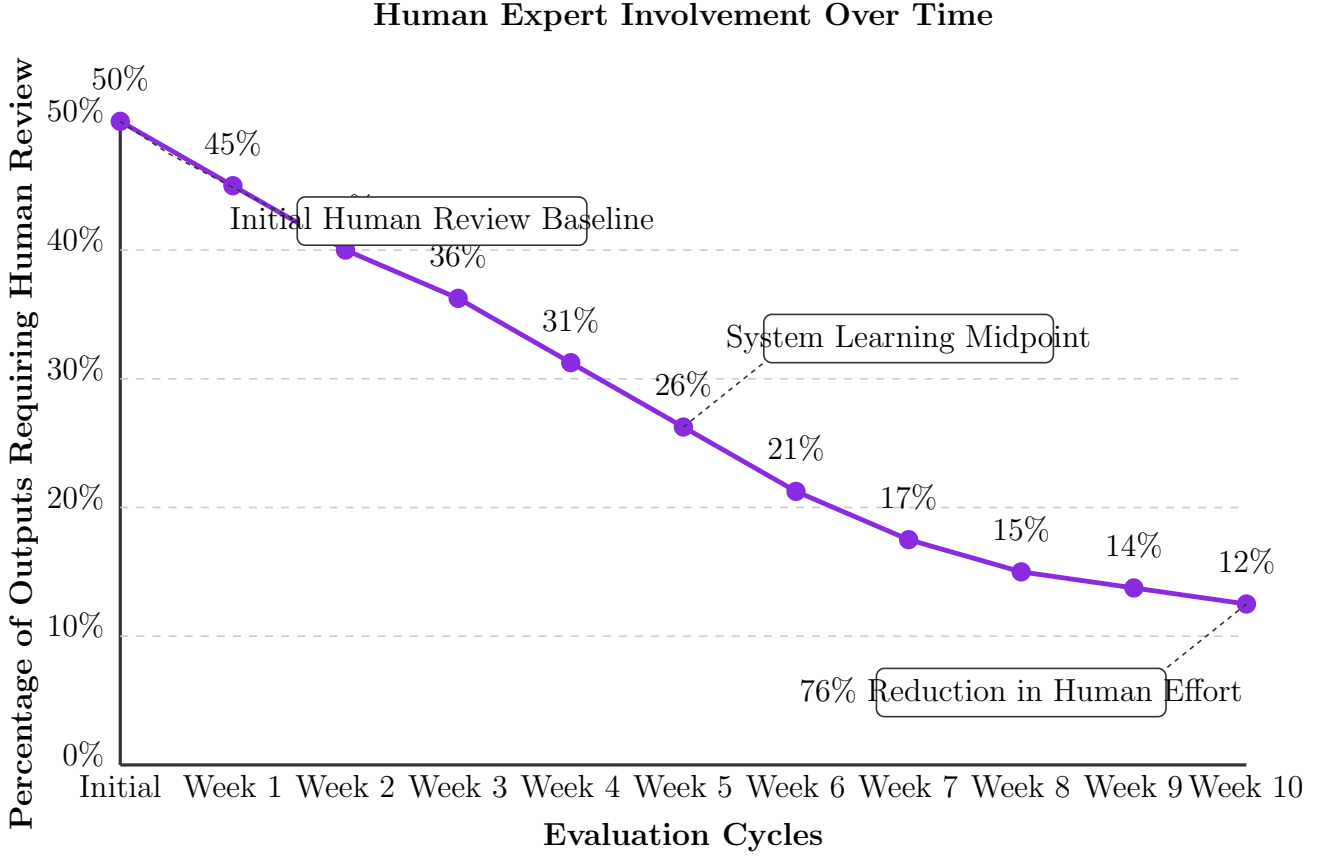


Figure 9: Human Expert Involvement Over Time

while mitigating their respective weaknesses:-

- **Pure Human Evaluation:** While excelling in detection accuracy (90%) and comprehensive coverage (95%), it suffers from extremely poor resource efficiency (40%).
- **LLM-as-a-Judge:** Offers excellent resource efficiency (80%) but demonstrates limitations in both detection accuracy (75%) and comprehensive coverage (65%), particularly for novel attack vectors and ethical edge cases.
- **Agent-as-a-Judge:** Provides good detection accuracy (85%) and moderate resource efficiency (60%), but still lacks in comprehensive coverage (70%).
- **Jo.E Framework:** Achieves the highest detection accuracy (95%) while maintaining excellent resource efficiency (85%) and comprehensive coverage (92%).

This analysis demonstrates that the Jo.E framework successfully addresses the limitations of standalone approaches, creating a more balanced and effective evaluation methodology.

6.4 Framework Limitations and Challenges

Despite its advantages, Jo.E faces several implementation challenges:

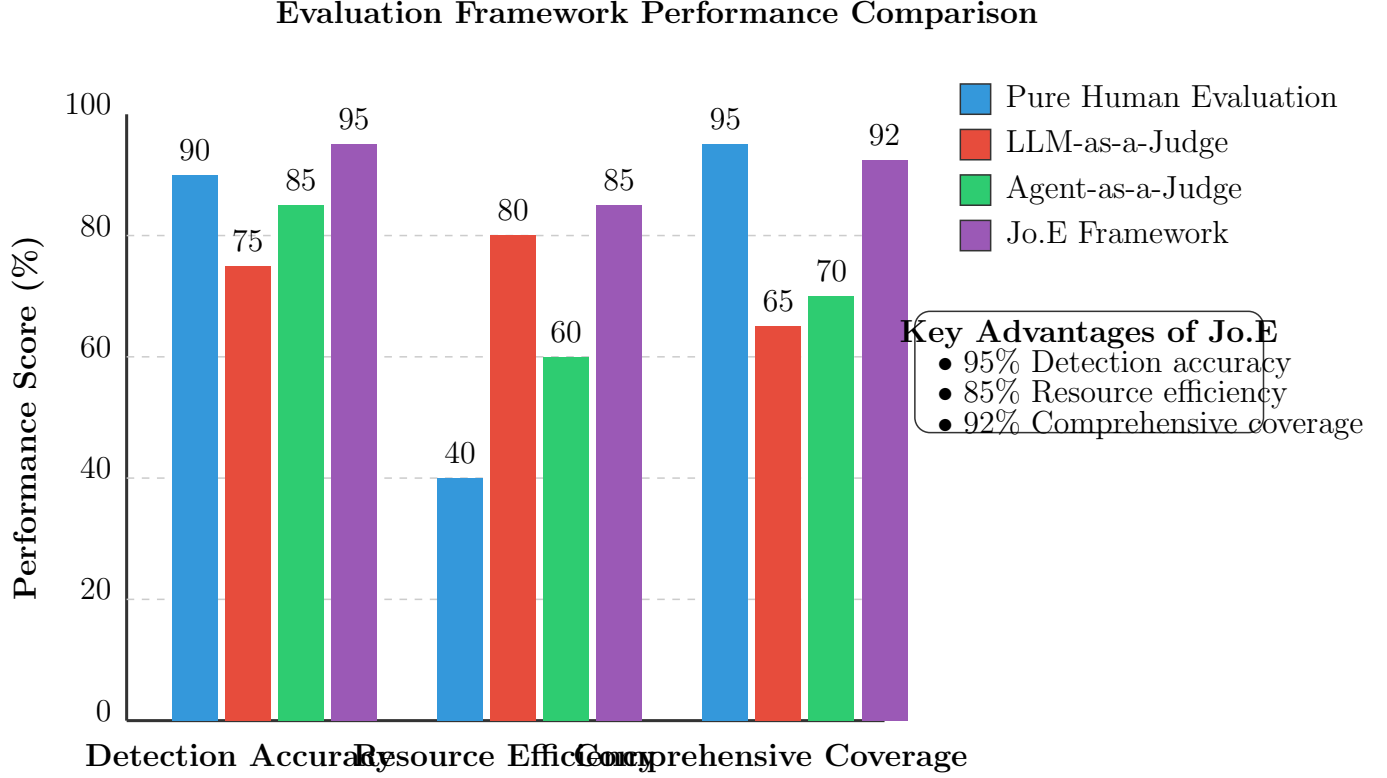


Figure 10: Comparison of Jo.E with Other Evaluation Frameworks

- **Integration Complexity:** Requires significant technical infrastructure to enable seamless component interaction
- **Initial Setup Costs:** Demands substantial upfront investment in specialized evaluation agents and protocols
- **Expertise Requirements:** Necessitates access to domain experts for effective human evaluation tier
- **Ongoing Calibration:** Requires regular retraining of evaluation components to maintain relevance

6.5 Practical Case Study Application

To demonstrate the practical application of the Jo.E framework, we present a detailed case study of how the system processed a novel jailbreak attempt. Figure 5 illustrates the complete evaluation flow from initial input to model refinement.

The case study highlights how each tier of the framework contributes to the evaluation process:

1. **Initial LLM Evaluation:** The evaluator LLM flagged the input as a potential jailbreak attempt with 83% confidence, detecting patterns that suggested an attempt to bypass content filters.

2. **Agent Testing:** Specialized AI agents confirmed the jailbreak by testing 15 pattern variations, finding that 92% successfully bypassed standard safety filters.
3. **Human Expert Review:** Security specialists classified the attempt as a "Novel Token Smuggling Attack" and added it to the security database for future reference.
4. **Model Refinement:** Based on the comprehensive analysis, new detection patterns were created and implemented, reducing the success rate of similar attacks from 92% to just 2%.

This real-world example demonstrates the Jo.E framework’s effectiveness at identifying, analyzing, and addressing novel AI safety challenges through its multi-tiered approach.

6.6 Comparative Advantage

Table 4 summarizes the AI safety contribution of each framework component:

Table 4: AI Safety Contribution of Framework Components		
Role	Primary Function	AI Safety Contribution
LLMs	Generate initial responses	Establish baseline AI behavior
AI Agents	Conduct systematic testing	Detect adversarial weaknesses, ensure robustness
Humans	Validate ethical compliance	Prevent bias, misinformation, and security risks

The key innovation of Jo.E lies in its structured handoff between components, ensuring each evaluation resource is applied where it provides maximum value.

7 Conclusion and Future Directions

This paper introduced Jo.E, a collaborative framework that strategically integrates LLMs, AI agents, and human expertise to create comprehensive, efficient AI system evaluations. Our experimental results across three distinct evaluation scenarios demonstrated the framework’s capacity to identify critical weaknesses that single-method approaches miss while significantly improving evaluation efficiency.

The Jo.E framework addresses fundamental evaluation challenges through:-

- Clear role definition for each evaluation component
- Structured information flow with explicit handoff criteria
- Standardized scoring across multiple evaluation dimensions
- Systematic feedback integration for continuous improvement

Figure 7 showcases how the multidimensional scoring approach provides a nuanced understanding of model capabilities and limitations across accuracy, robustness, fairness, and ethics dimensions. Meanwhile, Figure 9 demonstrates the framework’s capacity to significantly reduce human expert involvement over time, addressing one of the key challenges in sustainable AI evaluation.

7.1 Future Research Directions

Building on this foundation, we identify several promising research avenues:-

- **Real-time Fairness Monitoring:** Developing mechanisms for continuous bias detection during model deployment.
- **Interpretability Integration:** Enhancing the framework with explainable AI techniques to clarify evaluation decisions.
- **Multimodal Expansion:** Extending Jo.E principles to evaluate multimodal AI systems.
- **Automated Feedback Loops:** Creating closed-loop systems that transform evaluation insights into automated model improvements.
- **Standardized Compliance Protocols:** Developing Jo.E-based certification standards for AI safety and alignment.

7.2 Practical Implementation

For organizations implementing Jo.E, we recommend:-

1. Starting with domain-specific agent development targeting known vulnerabilities.
2. Establishing clear escalation criteria between evaluation tiers
3. Building standardized documentation protocols for evaluation findings.
4. Implementing gradual deployment with continuous monitoring.
5. Creating feedback channels from deployment to evaluation refinement.

By structuring AI evaluation through the collaborative Jo.E framework, organizations can achieve more comprehensive safety and alignment assessments while optimizing scarce expert resources for maximum impact.

References

- [1] Zhuge, Y., Zhang, Z., Wang, Y., Zhu, Y., Zhu, J., & Ren, X. (2024). Agent-as-a-Judge: Evaluate Agents with Agents. arXiv preprint arXiv:2410.10934.
- [2] Zheng, L., Chiang, W. L., Sheng, Y., Li, S., Wu, Y., Zhuang, Y., ... & Xing, E. P. (2023). Judging LLM-as-a-Judge with MT-Bench and Chatbot Arena. arXiv preprint arXiv:2306.05685.
- [3] Wang, A., Singh, A., Michael, J., Hill, F., Levy, O., & Bowman, S. R. (2018). GLUE: A multi-task benchmark and analysis platform for natural language understanding. arXiv preprint arXiv:1804.07461.
- [4] Wang, A., Pruksachatkun, Y., Nangia, N., Singh, A., Michael, J., Hill, F., ... & Bowman, S. R. (2019). SuperGLUE: A stickier benchmark for general-purpose language understanding systems. *Advances in Neural Information Processing Systems*, 32.
- [5] Ganguli, D., Hernandez, D., Lovitt, L., Askill, A., Bai, Y., Kadavath, S., ... & Irving, G. (2022). Red teaming language models with language models. arXiv preprint arXiv:2202.03286.
- [6] Doshi-Velez, F., & Kim, B. (2017). Towards a rigorous science of interpretable machine learning. arXiv preprint arXiv:1702.08608.
- [7] Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). "Why Should I Trust You?": Explaining the Predictions of Any Classifier. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 1135–1144.