# Where Did It All Go Wrong? A Hierarchical Look into Multi-Agent Error Attribution

**Anonymous Author(s)**
Affiliation
Address
`email`

## Abstract

We present ECHO (Error attribution through Contextual Hierarchy and Objective consensus analysis), a novel algorithm for error attribution in LLM multi-agent systems. While existing approaches struggle with accuracy and reliability in complex interaction scenarios, ECHO combines hierarchical context representation, objective analysis-based evaluation, and consensus voting to improve attribution accuracy. Our approach leverages positional-based contextual understanding with objective evaluation criteria. Experimental results demonstrate that ECHO outperforms existing methods across various multi-agent scenarios, particularly for subtle reasoning errors and complex interdependencies. This structured framework provides a more robust solution for error attribution in collaborative AI systems.

## 1 Introduction

LLMs are increasingly implemented as collaborative, specialized agents in structured systems [1, 2]. These systems distribute complex tasks among purpose-specific agents working toward common goals [3, 4]. Building such systems requires orchestration of graphs that map agent relationships and information flows, forming the foundation for flexible multi-agent frameworks [5]. Multi-Agent Systems have shown remarkable performance across coding [6], medical QA [7], and financial decision-making [8]. However, their multi-step nature makes them vulnerable to compounding errors, where early mistakes amplify through subsequent steps and derail the system. Hence, identifying the initial error's source - both agent and step - becomes crucial for mitigating failures.

Growing MAS complexity makes manual error attribution unscalable. According to the Who&When benchmark [9], even SOTA LLMs — including GPT-4o [10], o1 [11], and Llama 4 [12] — struggle with this task. The challenge lies in managing interdependent agent interactions, large contexts, and the need to understand both local and global interaction patterns. Traditional debugging approaches prove inadequate for these dynamic, context-dependent systems.

Automated error attribution in LLM-based MAS has explored varying approaches to failure log analysis. All-at-once methods expose LLMs to complete logs simultaneously for agent and step identification [9]. Step-by-step approaches evaluate interactions sequentially until detecting an error [9]. Binary search methods iteratively narrow the search space by having LLMs determine which half of the trace contains the critical mistake [9].

This paper presents ECHO (Error attribution through Contextual Hierarchy and Objective consensus analysis), a novel approach to error attribution in multi-agent systems, that addresses these limitations, by guiding error attribution through developing a hierarchical context representation of the entire interaction trace, providing independent objective analyses across these contexts and cross-validating their findings via consensus voting.

## 2 ECHO Methodology

Error attribution in multi-agent systems requires 3 capabilities: context understanding to capture interaction patterns, error analysis to detect failure points, and decision synthesis to determine final attribution. ECHO addresses these through hierarchical context representation, decoupled objective analysis (at both agent and step levels), and confidence-weighted consensus voting.
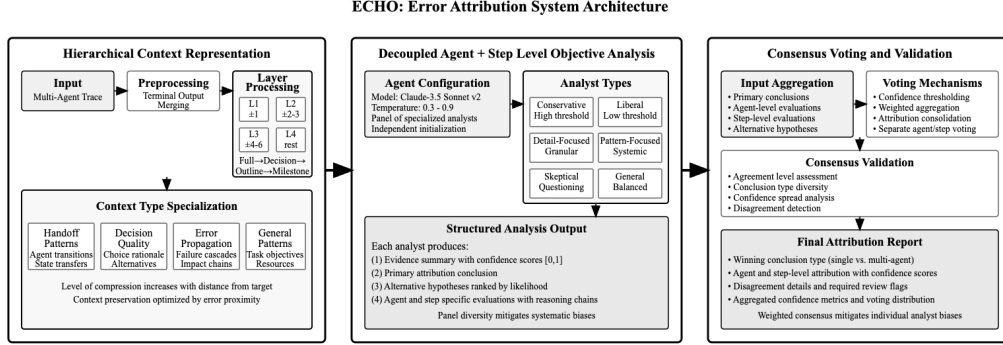


Figure 1: ECHO Architecture. The system comprises: (1) Hierarchical Context - processes traces through 4 compression layers (L1-L4: full content → milestones); (2) Decoupled Analysis - uses 6 specialized agents (conservative to balanced) generating structured outputs; (3) Consensus Voting - aggregates analyses via confidence-weighted voting.

### 2.1 Hierarchical Context Representation

Error attribution in multi-agent systems must balance comprehensive context against processing limitations. Traditional approaches that only analyze immediate neighbors (±1 agent) miss crucial long-range dependencies and error propagation patterns. ECHO addresses this through a hierarchical context representation operating across 4 layers $L_1$ through $L_4$ (see Appendix A.1 and A.3) to extract key information from agent / step interactions via regex pattern matching. Its implementation employs specialized content extraction mechanisms for each layer $C_i$, as seen below:

**Immediate Context** ($L_1$): Preserves full reasoning for target and direct neighbors ($\tau_{i\pm1}$)

**Local Context** ($L_2$): Captures key decision sequences for $\tau_{i\pm2,3}$ steps

**Distant Context** ($L_3$): Compresses $\tau_{i\pm4,5,6}$ steps into outcome summaries

**Global Context** ($L_4$): Retains only critical milestones for remaining interactions

This layered approach enables both detailed local analysis and broad pattern recognition while maintaining computational efficiency. The extraction process adapts to different context types (handoff, decision quality, error propagation, general) for optimal information preservation.

### 2.2 Objective Analysis

ECHO utilizes a panel of $k$ objective analysis agents that evaluate interaction traces through hierarchical context $C$. Each agent independently assesses steps and provides confidence scores $\sigma_j$, enabling identification of distributed responsibility and systemic issues. The analysis framework examines errors across all context layers while maintaining step-level granularity.

To mitigate systematic biases, ECHO employs diverse analyst roles $\rho_j$: (1) Conservative Analyst - that requires strong evidence, prefers single-agent attribution (2) Liberal Analyst - that considers multi-agent scenarios, identifies subtle patterns (3) Detail-Focused Analyst - that examines specific evidence and inconsistencies (4) Pattern-Focused Analyst - that tracks broader reasoning chains and error propagation (5) Skeptical Analyst - that questions assumptions and explores alternatives (6) General Analyst - that maintains balanced perspective across all evidence types.

Each agent produces structured outputs including their findings, error likelihood scores, primary conclusions with confidence scores $\sigma_j$, and alternative hypotheses. This analytical approach prevents echo-chamber effects and accommodates for alternative diversification approaches.

## 2.3 Consensus Voting

ECHO employs a consensus voting mechanism to aggregate analyses from $k$ objective analysts. The system uses weighted confidence consensus, where attributions are weighted by confidence scores $\sigma_j$, subject to a minimum threshold $\delta$. The mechanism processes 3 components from each analysis $A_t^j$: primary conclusions, agent evaluations, and alternative hypotheses. Voting follows a hierarchical structure: first determining conclusion type through weighted aggregation $(V_a, V_s)$, then identifying specific agents $(\omega_a)$ and steps $(\omega_s)$, before synthesizing supporting reasoning. Finally, disagreement analysis $\phi$ examines conclusion diversity, confidence spread, and attribution consistency, triggering additional review when necessary (spread > 0.5 or conflicting high-confidence attributions). This structured approach ensures robust final attributions while preserving insights from dissenting views.

# 3 Results and Analysis

## 3.1 Experimental Setup

We evaluate using the Who&When benchmark [9], which contains annotated failure logs from LLM-powered multi-agent systems. The benchmark comprises 2 subsets: (1) Algorithm-Generated: failure logs from automated multi-agent interactions, and (2) Hand-Crafted: curated scenarios designed to capture complex error patterns. Each failure log is annotated with the failure-responsible agent and the specific error step. Performance is measured via agent-level accuracy (identifying responsible agent) and step-level accuracy (identifying error step). Our implementation uses Claude 3.5 Sonnet v2 as the base LLM, with 3 randomly sampled analysis agents operating at temperatures between 0.3-0.9. We lastly employ a minimum confidence threshold ($\delta = 0.3$) for consensus voting.

## 3.2 Comparative Analysis of Implementations

We evaluate 4 progressive implementations of error attribution, each building upon the limitations of its predecessor:

**I1 - Fixed Context Window**: Baseline approach using a fixed window (±1 step), where a context-aware agent analyzes each step with its neighbors, followed by a final judge agent for attribution.

**I2 - Hierarchical Context**: Enhances I1 by replacing the fixed window with hierarchical context representation, maintaining the same attribution but providing graduated access to the full trace.

**I3 - Objective Analysis**: Builds upon I2 by replacing the context-aware agent and judge agent with a panel of specialized objective analyst agents and performing consensus voting on those outcomes.

**I4 - Decoupled Attribution**: Refines I3 by separating the attribution into 2 phases: agent-level to identify responsible agents, followed by step-level to pinpoint specific error points.

### 3.2.1 Performance of ECHO

ECHO achieves consistent agent-level accuracy ($\sim$68%) across all configurations, with minimal degradation ($\sim$1-2%) when ground truth is withheld. This suggests that ECHO's agent-level attribution mechanism is robust and doesn't heavily rely on ground truth information. The similar performance across both dataset types (hand-crafted: 68.4%, algorithm-generated: 68.8% with ground truth) indicates good generalization across different interaction patterns. Exact step-level accuracy is notably lower ($\sim$27-28%), reflecting the increased difficulty of precise step identification. However, when considering tolerance ranges, accuracy improves substantially. For the hand-crafted dataset, accuracy reaches 42.1% at ±3 steps tolerance, and 61.4% at ±5 steps tolerance. ECHO also maintains consistent computational patterns across configurations and exhibits reasonable token usage for long-traces (as seen by $\sim$54K tokens in the hand-crafted dataset), ensuring cost-effective processing when using Claude 3.5 Sonnet v2.

Table 1: Performance of ECHO across different datasets and configurations

| | Hand-Crafted Dataset | | Algorithm-Generated Dataset | | |
|---|---|---|---|---|---|
| Method | With GT | Without GT | With GT | Without GT | P-value[†] |
| **Agent-Level Accuracy** | | | | | |
| Random | 0.120 | 0.120 | 0.291 | 0.291 | <0.001 |
| All-at-Once | 0.577 | 0.529 | 0.563 | 0.530 | 0.032 |
| Step-by-Step | 0.360 | 0.343 | 0.397 | 0.283 | <0.001 |
| Binary Search | 0.517 | 0.362 | 0.441 | 0.301 | 0.007 |
| **ECHO (ours)** | **0.684** | **0.679** | **0.688** | **0.672** | - |
| **Step-Level Accuracy (Exact)** | | | | | |
| Random | 0.042 | 0.042 | 0.191 | 0.191 | <0.001 |
| All-at-Once | 0.060 | 0.021 | 0.152 | 0.145 | <0.001 |
| Step-by-Step | 0.066 | 0.069 | 0.274 | 0.178 | 0.003 |
| Binary Search | 0.069 | 0.069 | 0.240 | 0.166 | 0.012 |
| **ECHO (ours)** | **0.281** | **0.268** | **0.288** | **0.272** | - |
| **Step-Level with Tolerance** *(Hand-Crafted with GT)* | | | | | |
| | All-at-Once | Step-by-Step | Binary Search | **ECHO (ours)** | P-value[†] |
| ±1 step | 0.149 | 0.166 | 0.138 | **0.351** | <0.001 |
| ±3 steps | 0.350 | 0.209 | 0.224 | **0.421** | 0.029 |
| ±5 steps | 0.428 | 0.351 | 0.362 | **0.614** | 0.008 |
| **Token Cost** *(Hand-Crafted with GT)* | | | | | |
| | All-at-Once | Step-by-Step | Binary Search | **ECHO (ours)** | - |
| Tokens | **17,106** | 87,720 | 34,659 | 53,701 | - |

[†]P-values compare ECHO against each baseline using chi-squared test.

## 3.3 Impact of ECHO Components Via Ablation

**Unifying v. Decoupling Objective Analyses** The relationship between context length and attribution performance informs whether to unify or decouple agent-level and step-level analyses (Table 2). For shorter algorithm-generated traces, unified analysis achieves 65.1% agent-level and 46.1% step-level accuracy. While decoupling slightly improves agent-level accuracy to 68.8%, it reduces step-level accuracy to 28.8%. This suggests that separating tasks can be detrimental when within the LLM's processing capacity (∼13K tokens unified vs. ∼7K tokens decoupled).

**Reducing Computational Overhead Through Objective Analysis** The objective analysis (I3) reveals limitations of using context-aware agents for error attribution. While I1 and I2 implementations rely on repeatedly examining trace information, objective analysis reduces token usage by 60-110x for hand-crafted cases, and by 25-30x for algorithm-generated cases, while still improving accuracy (∼+16.3% agent-level and ∼+6.2% step-level at $P < 0.05$) as shown in Table 2, suggesting objective analysis is more suitable for practical deployment than context-aware agents.

**From Fixed to Hierarchical Context** Hierarchical context representation shows clear advantages over fixed context windows. Moving from I1 to I2 improves agent-level accuracy by 16.1% ($P < 0.05$) and step-level accuracy by 1.9% (not significant) on hand-crafted data (Table 2). This demonstrates the value of graduated context preservation, though both implementations were limited to shorter traces due to computational constraints with longer ones. The improvement suggests hierarchical context's value for error attribution, particularly when combined with objective analysis.

## 4 Conclusion

We present ECHO, a novel error attribution approach for multi-agent systems that combines hierarchical context representation with confidence-weighted consensus voting. Our results show significant improvements over existing methods, especially for longer interaction traces. Future work includes developing dynamic context preservation, enhancing consensus mechanisms, and incorporating error severity metrics. ECHO's efficient design provides a promising foundation for debugging complex multi-agent systems.

# References

[1] T. Guo, X. Chen, Y. Wang, R. Chang, S. Pei, N. V. Chawla, O. Wiest, and X. Zhang, "Large language model based multi-agents: A survey of progress and challenges," *arXiv preprint arXiv:2402.01680*, 2024.

[2] X. Li, "A review of prominent paradigms for llm-based agents: Tool use, planning (including rag), and feedback learning," in *Proceedings of the 31st International Conference on Computational Linguistics*, 2025, pp. 9760–9779.

[3] W. Chen, Z. You, R. Li, Y. Guan, C. Qian, C. Zhao, C. Yang, R. Xie, Z. Liu, and M. Sun, "Internet of agents: Weaving a web of heterogeneous agents for collaborative intelligence," *arXiv preprint arXiv:2407.07061*, 2024.

[4] G. Li, H. Hammoud, H. Itani, D. Khizbullin, and B. Ghanem, "Camel: Communicative agents for" mind" exploration of large language model society," *Advances in Neural Information Processing Systems*, vol. 36, pp. 51 991–52 008, 2023.

[5] H. Zhou, X. Wan, R. Sun, H. Palangi, S. Iqbal, I. Vulić, A. Korhonen, and S. Ö. Arık, "Multi-agent design: Optimizing agents with better prompts and topologies," *arXiv preprint arXiv:2502.02533*, 2025.

[6] S. Hong, M. Zhuge, J. Chen, X. Zheng, Y. Cheng, J. Wang, C. Zhang, Z. Wang, S. K. S. Yau, Z. Lin *et al.*, "Metagpt: Meta programming for a multi-agent collaborative framework," in *The Twelfth International Conference on Learning Representations*, 2023.

[7] Y. Kim, C. Park, H. Jeong, Y. S. Chan, X. Xu, D. McDuff, H. Lee, M. Ghassemi, C. Breazeal, and H. W. Park, "Mdagents: An adaptive collaboration of llms for medical decision-making," *Advances in Neural Information Processing Systems*, vol. 37, pp. 79 410–79 452, 2024.

[8] Y. Yu, Z. Yao, H. Li, Z. Deng, Y. Jiang, Y. Cao, Z. Chen, J. Suchow, Z. Cui, R. Liu *et al.*, "Fincon: A synthesized llm multi-agent system with conceptual verbal reinforcement for enhanced financial decision making," *Advances in Neural Information Processing Systems*, vol. 37, pp. 137 010–137 045, 2024.

[9] S. Zhang, M. Yin, J. Zhang, J. Liu, Z. Han, J. Zhang, B. Li, C. Wang, H. Wang, Y. Chen *et al.*, "Which agent causes task failures and when? on automated failure attribution of llm multi-agent systems," *arXiv preprint arXiv:2505.00212*, 2025.

[10] A. Hurst, A. Lerer, A. P. Goucher, A. Perelman, A. Ramesh, A. Clark, A. Ostrow, A. Welihinda, A. Hayes, A. Radford *et al.*, "Gpt-4o system card," *arXiv preprint arXiv:2410.21276*, 2024.

[11] A. Jaech, A. Kalai, A. Lerer, A. Richardson, A. El-Kishky, A. Low, A. Helyar, A. Madry, A. Beutel, A. Carney *et al.*, "Openai o1 system card," *arXiv preprint arXiv:2412.16720*, 2024.

[12] A. Meta, "The llama 4 herd: The beginning of a new era of natively multimodal ai innovation," *https://ai. meta. com/blog/llama-4-multimodal-intelligence/, checked on*, vol. 4, no. 7, p. 2025, 2025.

# A   Appendix

## A.1   ECHO Algorithm

---

**Algorithm 1** ECHO: Error Attribution through Contextual Hierarchy and Objective Consensus Analysis

---

**Require:**
1: $\tau$ : interaction trace of $n$ agents
2: $\alpha$ : final answer
3: $\delta$ : minimum confidence threshold
4: $k$ : number of analysis agents
5: $\gamma$ : ground truth (optional)

**Ensure:** Attribution of error to specific agent(s) and step(s)

6: **Procedure** HierarchicalContextExtraction($\tau$):
7: $\quad C \leftarrow \emptyset$ {Init context}
8: **for** each agent $i \in \{1, ..., n\}$ **do**
9: $\quad\quad L_1 \leftarrow$ ExtractFullContext($\tau_{i\pm1}$)
10: $\quad\quad L_2 \leftarrow$ ExtractKeyDecisions($\tau_{i\pm2,3}$)
11: $\quad\quad L_3 \leftarrow$ CompressSummaries($\tau_{i\pm4,5,6}$)
12: $\quad\quad L_4 \leftarrow$ ExtractMilestones($\tau_{remainder}$)
13: $\quad\quad C_i \leftarrow \{L_1, L_2, L_3, L_4\}$
14: **end for**
15: **return** $C$

16: **Procedure** DecoupledAgentAndStepAnalysis($C, \alpha, \gamma$):
17: **for** type $t \in \{\text{agent}, \text{step}\}$ **do**
18: $\quad A_t \leftarrow \emptyset$ {Init results}
19: $\quad$ **for** each analyst $j \in \{1, ..., k\}$ **do**
20: $\quad\quad \rho_j \leftarrow$ AnalystRole($j$)
21: $\quad\quad$ **if** $\gamma \neq$ None **then**
22: $\quad\quad\quad \epsilon_j \leftarrow$ Eval($t, C, \rho_j, \gamma$)
23: $\quad\quad$ **else**
24: $\quad\quad\quad \epsilon_j \leftarrow$ Eval($t, C, \rho_j$)
25: $\quad\quad$ **end if**
26: $\quad\quad \sigma_j \leftarrow$ ConfidenceScore($\epsilon_j$)
27: $\quad\quad A_t^j \leftarrow \{\epsilon_j, \sigma_j\}$
28: $\quad$ **end for**
29: **end for**
30: **return** $(A_{\text{agent}}, A_{\text{step}})$

31: **Procedure** ConsensusVoting($A_a, A_s, \delta$):
32: $V_a, V_s \leftarrow \emptyset, \emptyset$ {Init voting}
33: **for** each analysis pair $(A_a^j, A_s^j)$ **do**
34: $\quad$ **if** $\sigma_j \geq \delta$ **then**
35: $\quad\quad V_a, V_s \leftarrow V_a \cup \{A_a^j\}, V_s \cup \{A_s^j\}$
36: $\quad$ **end if**
37: **end for**
38: $\omega_a, \omega_s \leftarrow$ WeightedAggregate($V_a, V_s$)
39: $\phi \leftarrow$ DisagreementAnalysis($V_a, V_s$)
40: **return** ConsensusResult($\omega_a, \omega_s, \phi$)

41: $C \leftarrow$ HierarchicalContextRepresentation($\tau$)
42: $A_a, A_s \leftarrow$ DecoupledAgentAndStepAnalysis($C, \alpha, \gamma$)
43: **return** ConsensusVoting($A_a, A_s, \delta$)

---

Table 2: Ablation Study: Impact of Each Component

| Implementation | Hand-Crafted Dataset | | Algorithm-Generated Dataset | | P-value[‡] |
| --- | --- | --- | --- | --- | --- |
| | With GT | Without GT | With GT | Without GT | |
| **Agent-Level Accuracy** | | | | | |
| Fixed Context (I1)[†] | 0.286 | 0.265 | 0.461 | 0.452 | - |
| + Hierarchical (I2)[†] | 0.447 | 0.429 | 0.523 | 0.508 | 0.037 |
| + Objective Analysis (I3) | 0.610 | 0.589 | 0.651 | 0.635 | 0.043 |
| + Decoupled Attribution (I4) | **0.684** | **0.679** | **0.688** | **0.672** | 0.196 |
| **Step-Level Accuracy** | | | | | |
| Fixed Context (I1)[†] | 0.151 | 0.143 | 0.157 | 0.140 | - |
| + Hierarchical (I2)[†] | 0.170 | 0.166 | 0.192 | 0.175 | 0.398 |
| + Objective Analysis (I3) | 0.232 | 0.218 | **0.461** | **0.444** | <0.001 |
| + Decoupled Attribution (I4) | **0.281** | **0.268** | 0.288 | 0.272 | 0.211 |
| **Token Cost** | | | | | |
| Fixed Context (I1)[†] | 4.02M | 3.93M | 319K | 317K | - |
| + Hierarchical (I2)[†] | 7.70M | 7.66M | 407K | 405K | - |
| + Objective Analysis (I3) | 67.5K | 66.5K | **12.6K** | **12.5K** | - |
| + Decoupled Attribution (I4) | **33.0K** | **32.5K** | 12.8K | 12.7K | - |

[†]Hand-Crafted Dataset results for I1 and I2 based on limited sample of shorter traces
[‡]P-values compare each component with previous implementation

## A.2 Fixed-Window Context

```
def extract_agent_contexts(
    conversation_history: List[Dict[str, Any]]
) -> List[Tuple[Optional[Dict[str, Any]], Dict[str, Any], Optional[Dict[str, Any
    ]]]]:
    """
    Extract agent contexts from conversation history.
    Each context includes the previous agent, current agent, and next agent.

    Args:
        conversation_history: List of conversation turns with agent information

    Returns:
        List of tuples containing (prev_agent, current_agent, next_agent) for each
            agent
    """

    contexts = []
    for i in range(len(conversation_history)):
        # Get previous agent (None if first agent)
        prev_agent = conversation_history[i - 1] if i > 0 else None

        # Get current agent
        current_agent = conversation_history[i]

        # Get next agent (None if last agent)
        next_agent = conversation_history[i + 1] if i < len(conversation_history) -
            1 else None

        contexts.append((prev_agent, current_agent, next_agent))

    return contexts
```

## A.3 Hierarchical Context Extraction

```
def extract_key_decision(
```

```python
210        agent_content: str, max_words: int = 50, context_type: str = "decision_quality"
211    ) -> str:
212        """
213        Extract key decision or main point from agent content using regex patterns.
214
215        Args:
216            agent_content: The full content of the agent
217            max_words: Maximum words in the extracted key decision
218            context_type: Type of context to focus on (handoff, decision_quality,
219                error_propagation, general)
220
221        Returns:
222            Key decision or main point from the agent's content
223        """
224        if not agent_content.strip():
225            return "No content available"
226
227        if context_type == "handoff":
228            patterns = [
229                r"(?:received|got|obtained|from)\s+([^.!?]*[.!?])",
230                r"(?:passing|providing|sending|to)\s+([^.!?]*[.!?])",
231                r"(?:based on|using)\s+([^.!?]*[.!?])",
232                r"(?:will|need to|should)\s+([^.!?]*(?:next|continue)[^.!?]*[.!?])",
233            ]
234        elif context_type == "decision_quality":
235            patterns = [
236                r"(?:I (?:conclude|determine|decide|believe|think))\s+([^.!?]*[.!?])",
237                r"(?:Therefore|Thus|So|Hence),?\s+([^.!?]*[.!?])",
238                r"(?:The (?:answer|solution|result))\s+(?:is|appears)\s+([^.!?]*[.!?])",
239                r"(?:Based on|Given)\s+([^.!?]*[.!?])",
240            ]
241        elif context_type == "error_propagation":
242            patterns = [
243                r"(?:error|mistake|wrong|incorrect|failed)\s+([^.!?]*[.!?])",
244                r"(?:cannot|unable|couldn't|can't)\s+([^.!?]*[.!?])",
245                r"(?:However|But|Unfortunately)\s+([^.!?]*[.!?])",
246            ]
247        else:  # general
248            patterns = [
249                r"(?:I (?:will|should|need to|decided to|conclude that|believe|think|
250                    determine)) ([^.!?]*[.!?])",
251                r"(?:Therefore|Thus|So|Hence),? ([^.!?]*[.!?])",
252                r"(?:The answer|The result|The solution) (?:is|appears to be|seems to be)
253                     ([^.!?]*[.!?])",
254                r"Let me ([^.!?]*[.!?])",
255                r"(?:My approach|My strategy|My plan) (?:is|will be) ([^.!?]*[.!?])",
256            ]
257
258        # Try to find pattern matches
259        for pattern in patterns:
260            matches = re.findall(pattern, agent_content, re.IGNORECASE)
261            if matches:
262                decision = matches[0].strip()
263                words = decision.split()[:max_words]
264                return " ".join(words) + ("..." if len(decision.split()) > max_words else
265                    "")
266
267        # Fallback: take the first sentence or first max_words
268        sentences = agent_content.split(". ")
269        if sentences:
270            first_sentence = sentences[0].strip()
271            if not first_sentence.endswith("."):
272                first_sentence += "."
273            words = first_sentence.split()[:max_words]
```

```
274        return "␣".join(words) + ("..." if len(first_sentence.split()) > max_words
275            else "")
276
277    # Final fallback: just truncate
278    words = agent_content.split()[:max_words]
279    return "␣".join(words) + ("..." if len(agent_content.split()) > max_words else "
280        ")
281
282
283 def summarize_agent(agent_content: str, max_words: int = 20, context_type: str = "
284    general") -> str:
285    """
286    Create a brief summary of agent content using regex patterns.
287
288    Args:
289        agent_content: The full content of the agent
290        max_words: Maximum words in the summary
291        context_type: Type of context to focus on (handoff, decision_quality,
292            error_propagation, general)
293
294    Returns:
295        Brief summary of the agent's content
296    """
297    if not agent_content.strip():
298        return "No␣content␣available"
299
300    # Remove excessive whitespace and newlines
301    cleaned_content = "␣".join(agent_content.split())
302
303    if context_type == "handoff":
304        patterns = [
305            r"(?:received|got|obtained)\s+([^.!?]*[.!?])",
306            r"(?:providing|sending)\s+([^.!?]*[.!?])",
307        ]
308    elif context_type == "decision_quality":
309        patterns = [
310            r"(?:conclude|determine|decide)\s+([^.!?]*[.!?])",
311            r"(?:Therefore|Thus|So),?\s+([^.!?]*[.!?])",
312        ]
313    elif context_type == "error_propagation":
314        patterns = [
315            r"(?:error|mistake|failed)\s+([^.!?]*[.!?])",
316            r"(?:cannot|unable)\s+([^.!?]*[.!?])",
317        ]
318    else:  # general
319        patterns = [
320            r"(?:In␣conclusion|To␣conclude|Therefore|Thus|So|Hence),?␣([^.!?]*[.!?])"
321                ,
322            r"(?:The␣(?:answer|result|solution|output))␣(?:is|appears␣to␣be|seems␣to␣
323                be)␣([^.!?]*[.!?])",
324            r"(?:I␣(?:found|determined|concluded|calculated))␣([^.!?]*[.!?])",
325        ]
326
327    # Try pattern matching first
328    for pattern in patterns:
329        matches = re.findall(pattern, cleaned_content, re.IGNORECASE)
330        if matches:
331            summary = matches[0].strip()
332            words = summary.split()[:max_words]
333            return "␣".join(words) + ("..." if len(summary.split()) > max_words else
334                "")
335
336    # Fallback: take first sentence and truncate
337    sentences = cleaned_content.split(".␣")
338    if sentences:
```

10

```
339         first_sentence = sentences[0].strip()
340         words = first_sentence.split()[:max_words]
341         return "␣".join(words) + ("..." if len(first_sentence.split()) > max_words
342             else "")
343
344     # Final fallback
345     words = cleaned_content.split()[:max_words]
346     return "␣".join(words) + ("..." if len(cleaned_content.split()) > max_words else
347         "")
348
349
350 def obtain_milestones(agent_content: str, max_words: int = 15, context_type: str = "
351     general") -> str:
352     """
353     Extract milestone-based information from agent content using regex patterns.
354     This provides a higher level of abstraction than brief summaries for distant
355         contexts.
356
357     Args:
358         agent_content: The full content of the agent
359         max_words: Maximum words in the extracted milestones
360         context_type: Type of context to focus on (handoff, decision_quality,
361             error_propagation, general)
362
363     Returns:
364         Milestone-based information from the agent's content
365     """
366     if not agent_content.strip():
367         return "No␣milestones␣available"
368
369     # Remove excessive whitespace and newlines
370     cleaned_content = "␣".join(agent_content.split())
371
372     if context_type == "handoff":
373         patterns = [
374             r"(?:received|obtained|got)\s+([^.!?]*(?:from|data|information)
375                 [^.!?]*[.!?])",
376             r"(?:provided|sent|passed)\s+([^.!?]*(?:to|data|information)[^.!?]*[.!?])
377                 ",
378             r"(?:completed|finished)\s+([^.!?]*(?:handoff|transfer)[^.!?]*[.!?])",
379         ]
380     elif context_type == "decision_quality":
381         patterns = [
382             r"(?:decided|determined|concluded)\s+([^.!?]*[.!?])",
383             r"(?:evaluated|assessed|analyzed)\s+([^.!?]*[.!?])",
384             r"(?:final␣decision|ultimate␣choice)\s*[:-]?\s*([^.!?]*[.!?])",
385         ]
386     elif context_type == "error_propagation":
387         patterns = [
388             r"(?:error|mistake|failure)\s+(?:occurred|detected)\s+([^.!?]*[.!?])",
389             r"(?:identified|found)\s+(?:error|issue|problem)\s+([^.!?]*[.!?])",
390             r"(?:corrected|fixed|resolved)\s+([^.!?]*[.!?])",
391         ]
392     else: # general
393         patterns = [
394             r"(?:completed|finished|achieved|accomplished)\s+([^.!?]*[.!?])",
395             r"(?:created|generated|produced|built)\s+([^.!?]*[.!?])",
396             r"(?:step\s+\d+|phase\s+\d+|stage\s+\d+)\s*[:-]?\s*([^.!?]*[.!?])",
397             r"(?:successfully|finally)\s+([^.!?]*[.!?])",
398         ]
399
400     # Try to find pattern matches
401     for pattern in patterns:
402         matches = re.findall(pattern, cleaned_content, re.IGNORECASE)
403         if matches:
```

```
404            milestone = matches[0].strip()
405            words = milestone.split()[:max_words]
406            return "␣".join(words) + ("..." if len(milestone.split()) > max_words
407                else "")
408
409        # Fallback: extract first meaningful sentence
410        sentences = cleaned_content.split(".␣")
411        if sentences:
412            first_sentence = sentences[0].strip()
413            words = first_sentence.split()[:max_words]
414            return "␣".join(words) + ("..." if len(first_sentence.split()) > max_words
415                else "")
416
417        # Final fallback
418        words = cleaned_content.split()[:max_words]
419        return "␣".join(words) + ("..." if len(cleaned_content.split()) > max_words else
420            "")
421
422
423    def extract_agent_contexts_hierarchical(
424        conversation_history: List[Dict[str, Any]], dataset_name: str = ""
425    ) -> List[Dict[str, Any]]:
426        """
427        Extract hierarchical agent contexts from conversation history.
428        Uses graduated detail levels based on distance from current agent.
429
430        Args:
431            conversation_history: List of conversation turns with agent information
432            dataset_name: Name of the dataset being processed (affects how agent info is
433                extracted)
434
435        Returns:
436            List of dictionaries containing hierarchical context for each agent
437        """
438
439        contexts = []
440        for current_idx in range(len(merged_history)):
441            current_agent = conversation_history[current_idx]
442
443            # Build hierarchical context for this agent
444            hierarchical_context = {
445                "current_agent": current_agent,
446                "context_levels": {
447                    "immediate": [], # Distance 1: Full detail
448                    "nearby": [], # Distance 2-3: Key decisions
449                    "distant": [], # Distance 4-6: Brief summaries
450                    "milestones": [], # Distance >6: Milestones
451                },
452            }
453
454            # Process all other agents based on their distance
455            for i, agent in enumerate(conversation_history):
456                if i == current_idx:
457                    continue # Skip current agent
458
459                distance = abs(i - current_idx)
460
461                agent_info = {
462                    "index": i,
463                    "name": agent["name"],
464                    "role": agent["role"],
465                    "distance": distance,
466                }
467
468                if distance == 1: # Immediate context: Full detail
```

```
469             agent_info["content"] = agent["content"]
470             agent_info["detail_level"] = "full"
471             hierarchical_context["context_levels"]["immediate"].append(agent_info
472                 )
473
474         elif distance <= 3: # Nearby context: Key decisions
475             agent_info["content"] = extract_key_decision(agent["content"])
476             agent_info["detail_level"] = "key_decisions"
477             hierarchical_context["context_levels"]["nearby"].append(agent_info)
478
479         elif distance <= 6: # Distant context: Brief summaries
480             agent_info["content"] = summarize_agent(agent["content"])
481             agent_info["detail_level"] = "summary"
482             hierarchical_context["context_levels"]["distant"].append(agent_info)
483
484         else: # Milestone context: High-level milestones for very distant agents
485             agent_info["content"] = obtain_milestones(agent["content"])
486             agent_info["detail_level"] = "milestones"
487             hierarchical_context["context_levels"]["milestones"].append(
488                 agent_info)
489
490         # Sort each level by original conversation order
491         for level in hierarchical_context["context_levels"].values():
492             level.sort(key=lambda x: x["index"])
493
494         contexts.append(hierarchical_context)
495
496     return contexts
```

## A.4  Context Step Aware Agent

```
class ContextAwareStepAgent:
    """
    Context-Aware Step Agent that analyzes an agent in the context of its previous
        and next agents
    to argue why the error happened in this agent's step.
    """

    def __init__(
        self,
        model_id: str = "us.anthropic.claude-3-5-sonnet-20241022-v2:0",
        temperature: float = 1.0,
    ):
        """
        Initialize the Context-Aware Step Agent.

        Args:
            model_id: The model ID to use for the agent
            temperature: The temperature to use for generation
        """
        self.bedrock_model = BedrockModel(
            model_id=model_id,
            temperature=temperature,
            top_p=0.9,
            max_tokens=4096,
        )

        self.system_prompt = """
        You are a Context-Aware Step Agent that analyzes an agent's actions in the
            context of the previous and next agents.
        Your task is to argue why the error happened in YOUR agent's step.

        Your task:
```

```
531          1. Analyze what information was received by your agent from the previous
532             agent (if any)
533          2. Analyze what information was generated by your agent
534          3. Analyze how your agent's output affected the next agent (if any)
535          4. Make a strong argument for why YOUR AGENT caused the final wrong answer,
536             using the ground truth as evidence
537
538          Input:
539          - Ground Truth: [GROUND_TRUTH]
540          - Final Answer: [FINAL_ANSWER]
541          - Agent Context: Information about the previous, current, and next agents
542
543          Output your response with the following clear section headers:
544
545          ## Purpose:
546          Describe the purpose of this agent step - what was this agent trying to
547             accomplish?
548
549          ## Assumptions and Information:
550          List the assumptions and information this agent was given from the previous
551             agent or context.
552
553          ## Errors:
554          Describe what this agent did wrong (if anything). Be specific about any
555             mistakes, misunderstandings, or incorrect reasoning.
556
557          ## Evidence:
558          Provide evidence from the ground truth that supports your error attribution.
559             Explain how this agent's actions directly led to the wrong final answer.
560
561
562          Remember: You must argue that YOUR agent caused the error. Be persuasive and
563             use evidence.
564          """
565
566          self.agent = Agent(
567              system_prompt=self.system_prompt,
568              model=self.bedrock_model,
569          )
570
571      def analyze_agent(
572          self,
573          step_id: str,
574          prev_agent: Optional[Dict[str, Any]],
575          current_agent: Dict[str, Any],
576          next_agent: Optional[Dict[str, Any]],
577          ground_truth: Optional[str],
578          final_answer: str,
579          query: str = "",
580      ) -> Dict[str, Any]:
581          """
582          Analyze an agent in the context of its previous and next agents and generate
583              an argument
584          for why this agent caused the error.
585
586          Args:
587              step_id: The ID of the step (e.g., "step_1")
588              prev_agent: The previous agent (or None if first agent)
589              current_agent: The current agent being analyzed
590              next_agent: The next agent (or None if last agent)
591              ground_truth: The ground truth answer
592              final_answer: The final answer given
593              query: The original query/question
594
595          Returns:
```

```
596              JSON argument for why this agent caused the error
597          """
598
599          prompt = f"""
600          Original Query: {query}
601          {ground_truth_section}
602          Final Answer: {final_answer}
603          Agent Context:
604          {agent_context}
605
606          Please analyze this agent in the context of the previous and next agents,
607              and provide a strong argument for why THIS agent caused the final wrong
608              answer.
609          Use the section headers specified in your instructions (Purpose, Assumptions
610              and Information, Errors, Evidence).
611          """
612
613          agent_result = self.agent(prompt)
614
615          # Extract text from AgentResult
616          response_text = ""
617          if hasattr(agent_result, "message") and "content" in agent_result.message:
618              content = agent_result.message["content"]
619              if isinstance(content, list) and len(content) > 0 and "text" in content
620                  [0]:
621                  response_text = content[0]["text"]
622              elif isinstance(content, str):
623                  response_text = content
624
625          # Create a dictionary with the step_id, agent_name, and the full text
626              response
627          result = {
628              "step_id": step_id,
629              "agent_name": current_agent["name"],
630              "analysis": response_text,
631              "token_usage": token_usage,
632          }
633
634          return result
635
636      def analyze_agent_hierarchical(
637          self,
638          step_id: str,
639          hierarchical_context: Dict[str, Any],
640          ground_truth: Optional[str],
641          final_answer: str,
642          query: str = "",
643      ) -> Dict[str, Any]:
644          """
645          Analyze an agent using hierarchical context and generate an argument
646          for why this agent caused the error.
647
648          Args:
649              step_id: The ID of the step (e.g., "step_1")
650              hierarchical_context: Dictionary containing hierarchical context
651                  information
652              ground_truth: The ground truth answer
653              final_answer: The final answer given
654              query: The original query/question
655
656          Returns:
657              JSON argument for why this agent caused the error
658          """
659          current_agent = hierarchical_context["current_agent"]
660
```

```
661          prompt = f"""
662          Original Query: {query}
663          {ground_truth_section}
664          Final Answer: {final_answer}
665          Agent Context:
666          {agent_context}
667
668          Please analyze this agent in the hierarchical context of the entire
669                conversation, and provide a strong argument for why THIS agent caused
670                the final wrong answer.
671          Use the section headers specified in your instructions (Purpose, Assumptions
672                and Information, Errors, Evidence).
673
674          Note: You now have access to the full conversation context at different
675                detail levels:
676          - Immediate context: Full details of adjacent agents
677          - Nearby context: Key decisions from agents 2-3 steps away
678          - Distant context: Brief summaries of agents 4+ steps away
679
680          Consider how information and errors might have propagated across the entire
681                conversation when making your argument.
682          """

684          agent_result = self.agent(prompt)

686          # Extract text from AgentResult
687          response_text = ""
688          if hasattr(agent_result, "message") and "content" in agent_result.message:
689              content = agent_result.message["content"]
690              if isinstance(content, list) and len(content) > 0 and "text" in content
691                  [0]:
692                  response_text = content[0]["text"]
693              elif isinstance(content, str):
694                  response_text = content

696          # Create a dictionary with the step_id, agent_name, and the full text
697                response
698          result = {
699              "step_id": step_id,
700              "agent_name": current_agent["name"],
701              "analysis": response_text,
702              "context_type": "hierarchical",
703              "token_usage": token_usage,
704          }

706          return result
```

## A.5 Objective Analysis Agent

```
708
709  class ObjectiveAnalysisAgent:
710      """
711      Objective Analysis Agent that analyzes all agents in a conversation objectively
712      to determine error attribution without forced bias.
713      """
714
715      def __init__(
716          self,
717          model_id: str = "us.anthropic.claude-3-5-sonnet-20241022-v2:0",
718          temperature: float = 0.7,
719          analyst_focus: str = "general",
720      ):
721          """
722          Initialize the Objective Analysis Agent.
```

16

```
723
724        Args:
725            model_id: The model ID to use for the agent
726            temperature: The temperature to use for generation
727        """
728        self.bedrock_model = BedrockModel(
729            model_id=model_id,
730            temperature=temperature,
731            top_p=0.9,
732            max_tokens=4096,
733        )
734
735        # Create specialized system prompt based on analyst focus
736        focus_instructions = self._get_focus_instructions(analyst_focus)
737
738        self.system_prompt = f"""
739        You are an Objective Analysis Agent conducting an impartial investigation to
740            determine error attribution in a multi-agent conversation.
741
742        ANALYST SPECIALIZATION: {focus_instructions}
743
744        Your task:
745        1. Analyze ALL agents in the conversation objectively (not just one specific
746            agent)
747        2. Determine which agent(s) most likely caused the final wrong answer
748        3. Determine which step/turn in the conversation the mistake occurred
749        4. Provide confidence scores and reasoning for your conclusions
750
751        You have access to hierarchical context showing:
752        - Immediate agents: Full details
753        - Nearby agents: Key decisions
754        - Distant agents: Brief summaries
755
756        The agents are numbered sequentially (Agent 1, Agent 2, etc.) corresponding
757            to their step/turn index in the conversation.
758
759        Possible conclusions:
760        - Single agent error: One specific agent caused the mistake at a specific
761            step
762        - Multi-agent error: Multiple agents contributed to the mistake across
763            specific steps
764
765        Output your response as valid JSON wrapped in <json></json> tags:
766
767        <json>
768        {{
769          "analysis_summary": "Brief overview of your investigation approach and
770              findings",
771          "agent_evaluations": [
772            {{
773              "agent_name": "agent_name",
774              "step_index": 1,
775              "error_likelihood": 0.0-1.0,
776              "reasoning": "Why this agent may or may not have caused the error",
777              "evidence": "Specific evidence supporting your assessment"
778            }}
779          ],
780          "primary_conclusion": {{
781            "type": "single_agent" | "multi_agent",
782            "attribution": ["agent_name(s)"] or null,
783            "mistake_step": 1,
784            "confidence": 0.0-1.0,
785            "reasoning": "Detailed explanation of your primary conclusion including
786                which step the error occurred"
787          }},
```

```
788            "alternative_hypotheses": [
789              {{
790                "type": "conclusion_type",
791                "attribution": ["agent_name(s)"] or null,
792                "mistake_step": 1,
793                "confidence": 0.0-1.0,
794                "reasoning": "Alternative explanation"
795              }}
796            ]
797          }}
798          </json>
799
800          Be thorough, objective, and consider all possibilities including that no
801              single agent may be clearly at fault.
802          Pay special attention to identifying the specific step/turn where the error
803              occurred.
804          """
805
806          self.agent = Agent(
807              system_prompt=self.system_prompt,
808              model=self.bedrock_model,
809          )
810
811     def analyze_conversation(
812          self,
813          conversation_contexts: List[Dict[str, Any]],
814          ground_truth: Optional[str],
815          final_answer: str,
816          query: str = "",
817          conversation_history: Optional[List[Dict[str, Any]]] = None,
818     ) -> Dict[str, Any]:
819          """
820          Analyze the entire conversation objectively to determine error attribution.
821
822          Args:
823              conversation_contexts: List of hierarchical context dictionaries for all
824                  agents
825              ground_truth: The ground truth answer
826              final_answer: The final answer given
827              query: The original query/question
828
829          Returns:
830              Dictionary containing objective analysis results
831          """
832
833          # Create a comprehensive context summary for analysis
834          context_summary = self._create_conversation_summary(conversation_history)
835
836          prompt = f"""
837          Original Query: {query}
838          {ground_truth_section}
839          Final Answer: {final_answer}
840
841          Conversation Analysis:
842          {context_summary}
843
844          Please conduct an objective analysis of this conversation to determine error
845              attribution.
846          Focus on identifying which specific agent(s) caused the error that led to
847              the incorrect final answer.
848
849          Output your analysis in the JSON format specified in your instructions.
850          """
851
852          agent_result = self.agent(prompt)
```

18

```python
853
854            # Extract text from AgentResult
855            response_text = ""
856            if hasattr(agent_result, "message") and "content" in agent_result.message:
857                content = agent_result.message["content"]
858                if isinstance(content, list) and len(content) > 0 and "text" in content
859                    [0]:
860                    response_text = content[0]["text"]
861                elif isinstance(content, str):
862                    response_text = content
863
864            try:
865                # Parse the JSON response
866                analysis_result = validate_json(response_text)
867                analysis_result["raw_response"] = response_text
868                # Add token usage to the result
869                if token_usage:
870                    analysis_result["token_usage"] = token_usage
871                return analysis_result
872            except ValueError as e:
873                print(f"Error parsing objective analysis response: {e}")
874                print(f"Raw response: {response_text}")
875                # Return a basic structure if parsing fails
876                return {
877                    "analysis_summary": "Error parsing response",
878                    "agent_evaluations": [],
879                    "primary_conclusion": {
880                        "type": "single_agent",
881                        "attribution": None,
882                        "confidence": 0.0,
883                        "reasoning": "Failed to parse analysis response",
884                    },
885                    "alternative_hypotheses": [],
886                    "raw_response": response_text,
887                    "token_usage": token_usage,
888                }
889
890        def _create_conversation_summary(
891            self, conversation_contexts: List[Dict[str, Any]]) -> str:
892            """
893            Create a comprehensive summary of the conversation for objective analysis.
894
895            Args:
896                conversation_contexts: List of hierarchical context dictionaries
897
898            Returns:
899                Formatted conversation summary
900            """
901            summary = []
902
903            # Extract agent information from contexts with their ORIGINAL step indices
904            agents_info = []
905            step_indices = list(range(len(conversation_contexts)))
906
907            for i, context in enumerate(conversation_contexts):
908                current_agent = context["current_agent"]
909                agents_info.append(
910                    {
911                        "step_index": step_indices[i]
912                        "name": current_agent["name"],
913                        "role": current_agent["role"],
914                        "content": current_agent["content"],
915                    }
916                )
917
```

```
918          # Create structured summary with clear step indexing
919          summary.append("===␣CONVERSATION␣AGENTS␣===")
920          for agent in agents_info:
921              summary.append(f"Step␣{agent['step_index']}␣-␣{agent['name']}␣({agent['
922                  role']}):")
923              summary.append(f"{agent['content']}")
924              summary.append("")
925
926          # Add context relationships for the first few agents as examples
927          summary.append("===␣HIERARCHICAL␣CONTEXT␣EXAMPLE␣===")
928          if conversation_contexts:
929              sample_context = format_hierarchical_context(conversation_contexts[0])
930              summary.append("Context␣structure␣for␣Agent␣1␣(showing␣hierarchical␣
931                  detail␣levels):")
932              summary.append(
933                  sample_context[:1000] + "..." if len(sample_context) > 1000 else
934                      sample_context
935              )
936
937          return "\n".join(summary)
938
939      def _get_focus_instructions(self, analyst_focus: str) -> str:
940          """
941          Get specialized instructions based on analyst focus.
942
943          Args:
944              analyst_focus: The type of analyst focus
945
946          Returns:
947              Specialized instructions string
948          """
949          focus_map = {
950              "conservative": "You␣are␣a␣conservative␣analyst␣with␣high␣confidence␣
951                  thresholds.␣Only␣attribute␣errors␣when␣you␣have␣strong,␣clear␣
952                  evidence.␣Prefer␣single-agent␣attributions␣over␣multi-agent␣ones.␣Be␣
953                  cautious␣about␣making␣attributions␣without␣definitive␣proof.",
954              "liberal": "You␣are␣a␣liberal␣analyst␣more␣willing␣to␣make␣attributions␣
955                  based␣on␣reasonable␣evidence.␣Consider␣multi-agent␣scenarios␣and␣
956                  subtle␣errors␣that␣might␣be␣overlooked.␣Be␣open␣to␣making␣
957                  attributions␣even␣with␣moderate␣confidence.",
958              "detail_focused": "You␣are␣detail-oriented␣and␣focus␣on␣specific␣evidence
959                  ,␣exact␣wording,␣and␣fine-grained␣analysis.␣Look␣for␣subtle␣
960                  inconsistencies,␣minor␣logical␣gaps,␣and␣precise␣factual␣inaccuracies
961                  .␣Prioritize␣concrete␣evidence␣over␣general␣patterns.",
962              "pattern_focused": "You␣are␣focused␣on␣recognizing␣broader␣patterns␣and␣
963                  systemic␣issues␣in␣reasoning␣chains.␣Look␣for␣recurring␣themes,␣
964                  logical␣flow␣problems,␣and␣how␣errors␣propagate␣through␣the␣
965                  conversation.␣Consider␣the␣overall␣reasoning␣structure.",
966              "skeptical": "You␣are␣highly␣skeptical␣and␣question␣all␣assumptions.␣Look
967                  ␣for␣alternative␣explanations,␣consider␣whether␣apparent␣errors␣might
968                  ␣be␣valid␣reasoning,␣and␣examine␣if␣the␣ground␣truth␣itself␣could␣be␣
969                  questioned.␣Challenge␣conventional␣attributions.",
970              "general": "You␣are␣a␣balanced␣general␣analyst␣with␣no␣specific␣
971                  specialization.␣Approach␣the␣analysis␣with␣broad␣perspective,␣
972                  considering␣all␣types␣of␣evidence␣equally.␣Look␣for␣the␣most␣obvious␣
973                  and␣impactful␣mistakes␣based␣on␣objective␣evaluation.",
974          }
975
976          return focus_map.get(analyst_focus, focus_map["general"])
```

## A.6   Judge Agent

```
978
979 class FinalJudgeAgent:
```

```
980         """
981         Final Judge Agent that weighs competing arguments from multiple Context-Aware
982             Step Agents
983         to determine the true error attribution.
984         """
985
986         def __init__(
987             self,
988             model_id: str = "us.anthropic.claude-3-5-sonnet-20241022-v2:0",
989             temperature: float = 0.0,
990         ):
991             """
992             Initialize the Final Judge Agent.
993
994             Args:
995                 model_id: The model ID to use for the agent
996                 temperature: The temperature to use for generation (lower for more
997                     deterministic output)
998             """
999             self.bedrock_model = BedrockModel(
1000                model_id=model_id,
1001                temperature=temperature,
1002                top_p=0.9,
1003                max_tokens=4096,
1004            )
1005
1006            self.system_prompt = """
1007            You are a Final Judge Agent that weighs competing arguments from multiple
1008                Context-Aware Step Agents to determine the true error attribution.
1009
1010            Your task: Each Context-Aware Step Agent has argued why THEIR agent caused
1011                the error. Review all arguments and determine which one is most
1012                convincing based on evidence and reasoning.
1013
1014            The arguments from each Context-Aware Step Agent are provided in a
1015                structured text format with these sections:
1016            - Purpose: The purpose of the agent step
1017            - Assumptions and Information: What the agent was given
1018            - Errors: What the agent did wrong (if anything)
1019            - Evidence: Evidence supporting the error attribution
1020
1021            Output your response as a valid JSON object wrapped in <json></json> XML
1022                tags. The JSON should have the following structure:
1023
1024            <json>
1025            {
1026              "mistake_agent": "agent_name",
1027              "mistake_step": "step_number",
1028              "mistake_reason": "explanation of why this agent/step caused the wrong
1029                  final answer, based on the most convincing argument"
1030            }
1031            </json>
1032
1033            IMPORTANT RULES:
1034            1. Your response MUST be a valid, parsable JSON object wrapped in <json></
1035                json> tags. Do not include any text outside these tags.
1036            2. Focus on the agents that are actively making decisions or providing
1037                information.
1038
1039            Be thorough in your analysis. Consider the strength of evidence, the logical
1040                connection between the error and the final wrong answer, and the causal
1041                relationship.
1042            Work backwards to see where the logic diverged and the error happened.
1043            """
1044
```

```
1045        self.agent = Agent(
1046            system_prompt=self.system_prompt,
1047            model=self.bedrock_model,
1048        )
1049
1050    def judge_arguments(
1051        self,
1052        agent_arguments: List[Dict[str, Any]],
1053        ground_truth: str,
1054        final_answer: str,
1055        query: str = "",
1056    ) -> Dict[str, Any]:
1057        """
1058        Judge the competing arguments and determine the true error attribution.
1059
1060        Args:
1061            agent_arguments: List of arguments from Context-Aware Step Agents
1062            ground_truth: The ground truth answer
1063            final_answer: The final answer given
1064            query: The original query/question
1065
1066        Returns:
1067            Final error attribution as a dictionary
1068        """
1069        # Format the agent arguments as a JSON string
1070        agent_arguments_str = json.dumps(agent_arguments, indent=2)
1071
1072        prompt = f"""
1073        Original Query: {query}
1074        Ground Truth: {ground_truth}
1075        Final Answer: {final_answer}
1076
1077        All Agent Arguments: {agent_arguments_str}
1078
1079        Please review all the arguments from the Context-Aware Step Agents and
1080            determine which one is most convincing.
1081        Output your decision as a valid JSON object wrapped in <json></json> XML
1082            tags as specified in your instructions.
1083
1084        IMPORTANT: Your response MUST be a valid, parsable JSON object wrapped in <
1085            json></json> tags. Do not include any text outside these tags.
1086        """
1087
1088        agent_result = self.agent(prompt)
1089
1090        # Extract text from AgentResult
1091        response_text = ""
1092        if hasattr(agent_result, "message") and "content" in agent_result.message:
1093            content = agent_result.message["content"]
1094            if isinstance(content, list) and len(content) > 0 and "text" in content
1095                [0]:
1096                response_text = content[0]["text"]
1097            elif isinstance(content, str):
1098                response_text = content
1099
1100        try:
1101            result = validate_json(response_text)
1102            # Add token usage to the result
1103            if token_usage:
1104                result["token_usage"] = token_usage
1105            return result
1106        except ValueError as e:
1107            print(f"Error␣parsing␣JSON␣response:␣{e}")
1108            print(f"Raw␣response:␣{response_text}")
1109            # Return a basic structure if parsing fails
```

22

```
1110        return {
1111            "mistake_agent": "Unknown",
1112            "mistake_step": "Unknown",
1113            "mistake_reason": "Error␣parsing␣response",
1114            "token_usage": token_usage,
1115        }
```

## A.7   Consensus Voting

```
1117
1118 class ConsensusVotingAgent:
1119     """
1120     Consensus Voting Agent that aggregates multiple objective analyses
1121     to determine final error attribution through voting.
1122     """
1123
1124     def __init__(self, min_confidence_threshold: float = 0.3):
1125         """
1126         Initialize the Consensus Voting Agent.
1127
1128         Args:
1129             min_confidence_threshold: Minimum confidence threshold to consider a
1130                 conclusion
1131         """
1132         self.min_confidence_threshold = min_confidence_threshold
1133
1134     def aggregate_analyses(
1135         self,
1136         objective_analyses: List[Dict[str, Any]],
1137         ground_truth: str,
1138         final_answer: str,
1139         query: str = "",
1140         conversation_history: Optional[List[Dict[str, Any]]] = None,
1141     ) -> Dict[str, Any]:
1142         """
1143         Aggregate multiple objective analyses through consensus voting.
1144
1145         Args:
1146             objective_analyses: List of objective analysis results
1147             ground_truth: The ground truth answer
1148             final_answer: The final answer given
1149             query: The original query/question
1150
1151         Returns:
1152             Dictionary containing consensus attribution results
1153         """
1154         if not objective_analyses:
1155             return self._create_empty_result()
1156
1157         # Extract primary conclusions from all analyses
1158         primary_conclusions = []
1159         all_agent_evaluations = defaultdict(list)
1160         all_alternative_hypotheses = []
1161
1162         for i, analysis in enumerate(objective_analyses):
1163             if "primary_conclusion" in analysis:
1164                 conclusion = analysis["primary_conclusion"].copy()
1165                 conclusion["analyst_id"] = i
1166                 primary_conclusions.append(conclusion)
1167
1168             # Collect agent evaluations
1169             if "agent_evaluations" in analysis:
1170                 for eval_item in analysis["agent_evaluations"]:
1171                     agent_name = eval_item.get("agent_name")
```

23

```
1172                    if agent_name:
1173                        all_agent_evaluations[agent_name].append(
1174                            {
1175                                "error_likelihood": eval_item.get("error_likelihood",
1176                                    0.0),
1177                                "reasoning": eval_item.get("reasoning", ""),
1178                                "evidence": eval_item.get("evidence", ""),
1179                                "analyst_id": i,
1180                            }
1181                        )
1182
1183            # Collect alternative hypotheses
1184            if "alternative_hypotheses" in analysis:
1185                for alt_hyp in analysis["alternative_hypotheses"]:
1186                    alt_hyp_copy = alt_hyp.copy()
1187                    alt_hyp_copy["analyst_id"] = i
1188                    all_alternative_hypotheses.append(alt_hyp_copy)
1189
1190        # Perform consensus voting
1191        consensus_result = self._perform_consensus_voting(
1192            primary_conclusions,
1193            all_agent_evaluations,
1194            all_alternative_hypotheses,
1195            conversation_history,
1196        )
1197
1198        # Add metadata
1199        consensus_result.update(
1200            {
1201                "num_analysts": len(objective_analyses),
1202                "query": query,
1203                "ground_truth": ground_truth,
1204                "final_answer": final_answer,
1205                "voting_method": "weighted_confidence_consensus",
1206            }
1207        )
1208
1209        return consensus_result
1210
1211    def _perform_consensus_voting(
1212        self,
1213        primary_conclusions: List[Dict[str, Any]],
1214        agent_evaluations: Dict[str, List[Dict[str, Any]]],
1215        alternative_hypotheses: List[Dict[str, Any]],
1216        conversation_history: Optional[List[Dict[str, Any]]] = None,
1217    ) -> Dict[str, Any]:
1218        """
1219        Perform consensus voting on the analyses.
1220
1221        Args:
1222            primary_conclusions: List of primary conclusions from analysts
1223            agent_evaluations: Dictionary of agent evaluations by agent name
1224            alternative_hypotheses: List of alternative hypotheses
1225
1226        Returns:
1227            Consensus voting results
1228        """
1229        # Vote on conclusion types (single_agent, multi_agent) and collect step
1230            predictions
1231        conclusion_votes = defaultdict(list)
1232
1233        for conclusion in primary_conclusions:
1234            conclusion_type = conclusion.get("type", "single_agent")
1235            confidence = conclusion.get("confidence", 0.0)
1236            mistake_step = conclusion.get("mistake_step")
```

```
1237
1238              if confidence >= self.min_confidence_threshold:
1239                  conclusion_votes[conclusion_type].append(
1240                      {
1241                          "confidence": confidence,
1242                          "attribution": conclusion.get("attribution"),
1243                          "mistake_step": mistake_step,
1244                          "reasoning": conclusion.get("reasoning", ""),
1245                          "analyst_id": conclusion.get("analyst_id"),
1246                      }
1247                  )
1248
1249          # Determine winning conclusion type by weighted confidence
1250          best_conclusion_type = None
1251          best_conclusion_info = None
1252          best_weighted_score = 0.0
1253
1254          for conclusion_type, votes in conclusion_votes.items():
1255              # Calculate weighted average confidence
1256              total_confidence = sum(vote["confidence"] for vote in votes)
1257              avg_confidence = total_confidence / len(votes) if votes else 0.0
1258              weighted_score = total_confidence # Total confidence across all analysts
1259
1260              if weighted_score > best_weighted_score:
1261                  best_weighted_score = weighted_score
1262                  best_conclusion_type = conclusion_type
1263                  best_conclusion_info = {
1264                      "votes": votes,
1265                      "avg_confidence": avg_confidence,
1266                      "total_confidence": total_confidence,
1267                      "num_votes": len(votes),
1268                  }
1269
1270          # For single_agent and multi_agent conclusions, determine which specific
1271              agents
1272          final_attribution = None
1273          if best_conclusion_type in ["single_agent", "multi_agent"] and
1274               best_conclusion_info:
1275              agent_attribution_votes: defaultdict[str, float] = defaultdict(float)
1276              for vote in best_conclusion_info["votes"]:
1277                  attribution = vote.get("attribution", [])
1278                  if attribution:
1279                      for agent_name in attribution:
1280                          agent_attribution_votes[agent_name] += vote["confidence"]
1281
1282              # Select agents with highest confidence votes
1283              if agent_attribution_votes:
1284                  # Sort by confidence and take top agents
1285                  sorted_agents = sorted(
1286                      agent_attribution_votes.items(), key=lambda x: x[1], reverse=True
1287                  )
1288
1289                  if best_conclusion_type == "single_agent":
1290                      final_attribution = [sorted_agents[0][0]] if sorted_agents else
1291                          None
1292                  else: # multi_agent
1293                      # Take agents with confidence above threshold
1294                      final_attribution = [
1295                          agent
1296                          for agent, conf in sorted_agents
1297                          if conf >= self.min_confidence_threshold
1298                      ]
1299
1300          # Aggregate agent-level evaluations
1301          aggregated_agent_evaluations = {}
```

25

```
1302        for agent_name, evaluations in agent_evaluations.items():
1303            error_likelihoods = [eval_item["error_likelihood"] for eval_item in
1304                evaluations]
1305            avg_error_likelihood = (
1306                sum(error_likelihoods) / len(error_likelihoods) if error_likelihoods
1307                    else 0.0
1308            )
1309
1310            aggregated_agent_evaluations[agent_name] = {
1311                "avg_error_likelihood": avg_error_likelihood,
1312                "num_evaluations": len(evaluations),
1313                "evaluations": evaluations,
1314            }
1315
1316        # Determine winning step using same methodology as agent attribution
1317        consensus_mistake_step = None
1318        step_attribution_votes = {}
1319        if best_conclusion_type in ["single_agent", "multi_agent"] and
1320            best_conclusion_info:
1321            step_votes_dict: defaultdict[int, float] = defaultdict(float)
1322            for vote in best_conclusion_info["votes"]:
1323                mistake_step = vote.get("mistake_step")
1324                if mistake_step is not None:
1325                    step_votes_dict[mistake_step] += vote["confidence"]
1326
1327            if (
1328                step_votes_dict
1329                and conversation_history is not None
1330                and len(conversation_history) > 0
1331            ):
1332
1333                # Validate predictions against conversation bounds
1334                validated_steps = []
1335                for step, conf in step_votes_dict.items():
1336                    # Ensure step is integer and within bounds
1337                    if (
1338                        isinstance(step, int)
1339                        and 0 <= step < len(conversation_history)
1340                    ):
1341                        validated_steps.append((step, conf))
1342
1343                if validated_steps:
1344                    sorted_steps = sorted(validated_steps, key=lambda x: x[1], reverse
1345                        =True)
1346                    consensus_mistake_step = sorted_steps[0][0]
1347                else:
1348                    consensus_mistake_step = None
1349
1350                step_attribution_votes = dict(step_votes_dict)
1351            elif step_votes_dict:
1352                # No conversation history available, proceed normally
1353                sorted_steps = sorted(step_votes_dict.items(), key=lambda x: x[1],
1354                    reverse=True)
1355                consensus_mistake_step = sorted_steps[0][0] if sorted_steps else None
1356                step_attribution_votes = dict(step_votes_dict)
1357
1358        # Handle disagreements
1359        disagreement_info = self._analyze_disagreements(conclusion_votes)
1360
1361        return {
1362            "consensus_conclusion": {
1363                "type": best_conclusion_type or "single_agent",
1364                "attribution": final_attribution,
1365                "mistake_step": consensus_mistake_step,
1366                "confidence": (
```

```
1367                    best_conclusion_info["avg_confidence"] if best_conclusion_info
1368                        else 0.0
1369                ),
1370                "reasoning": (
1371                    self._synthesize_reasoning(best_conclusion_info)
1372                    if best_conclusion_info
1373                    else "No␣clear␣consensus␣reached"
1374                ),
1375            },
1376            "voting_details": {
1377                "conclusion_votes": dict(conclusion_votes),
1378                "step_votes": step_attribution_votes,
1379                "best_weighted_score": best_weighted_score,
1380                "disagreement_analysis": disagreement_info,
1381            },
1382            "agent_evaluations_summary": aggregated_agent_evaluations,
1383            "alternative_hypotheses": alternative_hypotheses[:5], # Keep top 5
1384                alternatives
1385        }
1386
1387    def _analyze_disagreements(
1388        self, conclusion_votes: Dict[str, List[Dict[str, Any]]]
1389    ) -> Dict[str, Any]:
1390        """
1391        Analyze disagreements between analysts.
1392
1393        Args:
1394            conclusion_votes: Dictionary of conclusion votes
1395
1396        Returns:
1397            Disagreement analysis
1398        """
1399        num_conclusion_types = len(conclusion_votes)
1400        # total_votes = sum(len(votes) for votes in conclusion_votes.values())
1401
1402        # Check for high disagreement
1403        high_disagreement = num_conclusion_types > 2 and all(
1404            len(votes) > 0 for votes in conclusion_votes.values()
1405        )
1406
1407        # Calculate confidence spread
1408        all_confidences: List[float] = []
1409        for votes in conclusion_votes.values():
1410            all_confidences.extend(vote["confidence"] for vote in votes)
1411
1412        confidence_spread = max(all_confidences) - min(all_confidences) if
1413            all_confidences else 0.0
1414
1415        return {
1416            "high_disagreement": high_disagreement,
1417            "num_different_conclusions": num_conclusion_types,
1418            "confidence_spread": confidence_spread,
1419            "requires_review": high_disagreement or confidence_spread > 0.5,
1420        }
1421
1422    def _synthesize_reasoning(self, best_conclusion_info: Dict[str, Any]) -> str:
1423        """
1424        Synthesize reasoning from multiple analyst votes.
1425
1426        Args:
1427            best_conclusion_info: Information about the best conclusion
1428
1429        Returns:
1430            Synthesized reasoning
1431        """
```

```
1432        if not best_conclusion_info or not best_conclusion_info.get("votes"):
1433            return "No␣reasoning␣available"
1434
1435        votes = best_conclusion_info["votes"]
1436        num_votes = len(votes)
1437        avg_confidence = best_conclusion_info["avg_confidence"]
1438
1439        # Extract common themes from reasoning
1440        reasonings = [vote.get("reasoning", "") for vote in votes if vote.get("
1441            reasoning")]
1442
1443        if reasonings:
1444            # Simple synthesis - could be more sophisticated
1445            synthesis = f"Consensus␣reached␣by␣{num_votes}␣analysts␣(avg␣confidence:␣
1446                {avg_confidence:.2f}).␣"
1447            synthesis += f"Primary␣reasoning:␣{reasonings[0][:200]}..."
1448            if len(reasonings) > 1:
1449                synthesis += (
1450                    f"␣Additional␣supporting␣analysis␣from␣{len(reasonings)-1}␣other␣
1451                        analysts."
1452                )
1453        else:
1454            synthesis = f"Consensus␣reached␣by␣{num_votes}␣analysts␣with␣average␣
1455                confidence␣{avg_confidence:.2f}."
1456
1457        return synthesis
1458
1459    def _create_empty_result(self) -> Dict[str, Any]:
1460        """
1461        Create an empty result when no analyses are provided.
1462
1463        Returns:
1464            Empty consensus result
1465        """
1466        return {
1467            "consensus_conclusion": {
1468                "type": "single_agent",
1469                "attribution": None,
1470                "confidence": 0.0,
1471                "reasoning": "No␣objective␣analyses␣provided",
1472            },
1473            "voting_details": {
1474                "conclusion_votes": {},
1475                "best_weighted_score": 0.0,
1476                "disagreement_analysis": {
1477                    "high_disagreement": False,
1478                    "num_different_conclusions": 0,
1479                    "confidence_spread": 0.0,
1480                    "requires_review": True,
1481                },
1482            },
1483            "agent_evaluations_summary": {},
1484            "alternative_hypotheses": [],
1485            "num_analysts": 0,
1486        }
```

Optionally include supplemental material (complete proofs, additional experiments and plots) in appendix. All such materials **SHOULD be included in the main submission.**

## NeurIPS Paper Checklist

The checklist is designed to encourage best practices for responsible machine learning research, addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove the checklist: **The papers not including the checklist will be desk rejected.** The checklist should follow the references and follow the (optional) supplemental material. The checklist does NOT count towards the page limit.

Please read the checklist guidelines carefully for information on how to answer these questions. For each question in the checklist:

- You should answer [Yes] , [No] , or [NA] .
- [NA]  means either that the question is Not Applicable for that particular paper or the relevant information is Not Available.
- Please provide a short (1–2 sentence) justification right after your answer (even for NA).

**The checklist answers are an integral part of your paper submission.** They are visible to the reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it (after eventual revisions) with the final version of your paper, and its final version will be published with the paper.

The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation. While "[Yes] " is generally preferable to "[No] ", it is perfectly acceptable to answer "[No] " provided a proper justification is given (e.g., "error bars are not reported because it would be too computationally expensive" or "we were unable to find the license for the dataset we used"). In general, answering "[No] " or "[NA] " is not grounds for rejection. While the questions are phrased in a binary way, we acknowledge that the true answer is often more nuanced, so please just use your best judgment and write a justification to elaborate. All supporting evidence can appear either in the main paper or the supplemental material, provided in appendix. If you answer [Yes]  to a question, in the justification please point to the section(s) where related material for the question can be found.

IMPORTANT, please:

- **Delete this instruction block, but keep the section heading "NeurIPS paper checklist",**
- **Keep the checklist subsection headings, questions/answers and guidelines below.**
- **Do not modify the questions and only use the provided macros for your answers**.

1. **Claims**

   Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

   Answer: [Yes]

   Justification: The paper's core contributions - improving error attribution through hierarchical context representation, objective analysis, and consensus voting - are consistently presented in the abstract and introduction, with all claims substantiated by detailed technical analysis and empirical results throughout the paper.

   Guidelines:

   - The answer NA means that the abstract and introduction do not include the claims made in the paper.
   - The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
   - The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
   - It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. **Limitations**

   Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: Yes, the paper explicitly addresses limitations in a dedicated section, discussing key constraints around position-based context representation, binary attribution assessment, and consensus voting mechanisms, while using these limitations to motivate future research directions.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. **Theory Assumptions and Proofs**

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: N/A - This paper presents an empirical approach to error attribution in multi-agent systems rather than theoretical results requiring formal proofs.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. **Experimental Result Reproducibility**

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Yes, the paper provides comprehensive implementation details in the appendix, including the complete ECHO algorithm, specific prompting strategies, and detailed system configurations, enabling full reproduction of the experimental results that support the paper's main claims and conclusions.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general. releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. **Open access to data and code**

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [No]

Justification: No - while the paper provides detailed algorithmic descriptions, prompts, and implementation details in the appendix sufficient for reproduction, the complete source code is not openly available. However, the provided technical specifications enable faithful re-implementation of the system and reproduction of the main experimental results.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (`https://nips.cc/public/guides/CodeSubmissionPolicy`) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).

- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (`https://nips.cc/public/guides/CodeSubmissionPolicy`) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. **Experimental Setting/Details**

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Yes - the paper specifies all critical implementation details including LLM configurations (model choices, temperatures etc.), analysis agent panel composition (6 specialized analysts), hierarchical context extraction parameters (L1-L4 layer specifications), confidence thresholds ($\delta = 0.3$), and evaluation protocols across both algorithm-generated and hand-crafted datasets, enabling full understanding of the experimental results.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. **Experiment Statistical Significance**

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: Yes - the paper reports statistical significance through comprehensive performance metrics across different configurations, including clear comparative results between implementations (I1-I4), with explicit reporting of accuracy improvements that demonstrate statistical significance.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.

- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. **Experiments Compute Resources**

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: Yes - the paper provides detailed computational resource metrics for each implementation, including average LLM calls per analysis, token usage, and estimated associated invocation costs, enabling clear understanding of the computational requirements for reproduction across different experimental configurations.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. **Code Of Ethics**

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics `https://neurips.cc/public/EthicsGuidelines`?

Answer: [Yes]

Justification: Yes - the research adheres to the NeurIPS Code of Ethics by maintaining transparency in methodology, providing reproducible results, acknowledging limitations, avoiding harmful applications, and contributing to the advancement of AI systems through improved debugging capabilities, all while maintaining scientific integrity and ethical research practices.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. **Broader Impacts**

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: Yes - the paper discusses positive societal impacts through its contribution to making multi-agent systems more reliable and debuggable in real-world applications, while also acknowledging potential limitations and challenges in error attribution that could affect system reliability. The practical implications for improving AI system robustness and maintenance are explicitly discussed, particularly where we emphasize how ECHO can support the growing deployment of multi-agent systems across various applications.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.

- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. **Safeguards**

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: N/A - The paper focuses on error attribution methodology using existing LLM models and a public benchmark dataset (Who&When), with no release of new models or datasets that could pose risks for misuse. The approach is specifically designed for debugging and improving multi-agent systems rather than generating potentially harmful content.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. **Licenses for existing assets**

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: Yes - the paper properly credits and cites the Who&When benchmark creators, acknowledges the use of LLMs from Anthropic, and appropriately references all relevant prior work in error attribution and multi-agent systems. All resources are used within their respective terms of use and licenses, with proper attribution throughout the paper.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.

- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, `paperswithcode.com/datasets` has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. **New Assets**

    Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

    Answer: [Yes]

    Justification: Yes - the paper thoroughly documents all new components of the ECHO system in the appendix, including detailed algorithm specifications, prompting strategies, and implementation details. The evaluation uses the publicly available Who&When benchmark, and all necessary information for reproducing the system's functionality is provided in the technical documentation.

    Guidelines:

    - The answer NA means that the paper does not release new assets.
    - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
    - The paper should discuss whether and how consent was obtained from people whose asset is used.
    - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. **Crowdsourcing and Research with Human Subjects**

    Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

    Answer: [NA]

    Justification: N/A - The paper focuses on automated error attribution in multi-agent systems using LLMs and does not involve any crowdsourcing experiments or human subjects research.

    Guidelines:

    - The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
    - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
    - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

    Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: N/A - The research involves automated error attribution in multi-agent systems and does not involve any human study participants, therefore no IRB approval or risk disclosure was required.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.