# Lethe: A Local-First Agent-Context Manager for Conversational AI with Adaptive Retrieval and Planning

Anonymous Author
Anonymous Institution
`anonymous@institution.edu`

August 25, 2025

### Abstract

We present LETHE, a local-first agent-context manager designed to understand and retrieve context from conversational AI interactions. Unlike traditional information retrieval systems focused on web search, LETHE specializes in parsing agent conversations—including tool interactions, planning sequences, and multi-turn reasoning traces—to provide contextually relevant information for ongoing agent assistance. Built on Node.js and SQLite for local deployment, LETHE implements adaptive planning (VERIFY/EXPLORE/EXPLOIT), entity-aware diversification, and hybrid retrieval optimized for agent-specific content patterns. We introduce LETHEBENCH-AGENTS, a dataset of agent conversation traces with weak supervision for tool-result recall, action consistency, and planning coherence evaluation. Comprehensive evaluation against six strong local baselines demonstrates significant improvements in agent-relevant metrics: tool-result recall (nDCG@10: +42.3%), planning consistency (+31.7%), and context retrieval quality (+28.9%), while maintaining sub-second response times and complete privacy through local-only operation. These results establish LETHE as a practical foundation for privacy-preserving agent assistance systems.

## 1 Introduction

Modern conversational AI agents engage in complex multi-turn interactions involving tool usage, planning, and reasoning across extended sessions. These agent conversations differ fundamentally from traditional text retrieval scenarios: they contain structured tool interactions, planning sequences, action-observation pairs, and contextual dependencies that span multiple turns. Current retrieval systems, designed primarily for web search or document retrieval [6], fail to capture these agent-specific patterns and relationships.

Consider an AI assistant helping with software development: it executes code, reads error outputs, searches documentation, and iteratively refines solutions across dozens of turns. Critical context includes not just the current query, but the sequence of failed attempts, successful tool invocations, and the evolving understanding of the problem. Traditional RAG systems [9, 4] treat each turn independently and lack awareness of agent-specific content structures like tool calls, planning states, or action dependencies.

Furthermore, agent conversations often contain sensitive information—API keys, personal data, proprietary code—making local processing essential for practical deployment. Cloud-based systems introduce unacceptable privacy risks and latency bottlenecks that disrupt natural conversation flow.

We introduce LETHE, a local-first agent-context manager that addresses these challenges through five key innovations:

1. **Agent-Aware Architecture**: Node.js and SQLite implementation designed for agent conversation patterns and local deployment

2. **Agent-Specific Content Understanding**: Parsing of tool interactions, planning sequences, and action-observation pairs

3. **Adaptive Planning Framework**: VERIFY/EXPLORE/EXPLOIT strategies based on conversation state and agent behavior patterns

4. **Entity-Aware Diversification**: Session-aware information diversity using agent-specific entity recognition

5. **Privacy-First Local Processing**: Complete on-device operation with no cloud dependencies or data leakage

We evaluate LETHE on LETHEBENCH-AGENTS, a comprehensive dataset of agent conversation traces with weak supervision labels for tool-result recall, action consistency, and planning coherence. Our experimental design tests four hypotheses across agent-specific quality metrics, efficiency constraints, coverage requirements, and adaptive planning effectiveness against six strong local baselines.

**Main Contributions:**

- A novel local-first agent-context manager with adaptive planning and agent-specific content understanding

- Agent-aware conversation parsing for tool interactions, planning sequences, and multi-turn reasoning traces

- LETHEBENCH-AGENTS, the first comprehensive evaluation dataset for agent conversation understanding with weak supervision labels

- Rigorous experimental evaluation demonstrating significant improvements in agent-specific metrics: tool-result recall (+42.3%), planning consistency (+31.7%), and context quality (+28.9%)

- Production-ready Node.js/SQLite implementation enabling immediate deployment in privacy-sensitive environments

## 2 Related Work

### 2.1 Agent Memory and Context Management

Contemporary AI agents rely on various memory architectures to maintain context across interactions. Traditional approaches use fixed-size context windows or sliding window strategies, but these fail to capture long-term dependencies and agent-specific patterns [6]. Recent work on memory-augmented agents focuses primarily on cloud-based systems with limited privacy considerations. Our work addresses this gap by providing local-first agent memory with adaptive context retrieval.

### 2.2 Local-First AI and Privacy-Preserving Systems

The local-first paradigm [5] emphasizes user control, privacy, and offline capability—critical requirements for agent deployment in sensitive environments. Recent work on browser-based AI [2] and edge deployment [3] enables sophisticated AI operations locally. However, existing local systems lack the specialized context management needed for multi-turn agent interactions with tool usage and planning.

### 2.3 Tool-Using AI Agents

Modern AI agents increasingly rely on tool interactions—API calls, code execution, database queries—to accomplish complex tasks. These interactions create structured conversation patterns distinct from natural language dialogue [8]. Existing retrieval systems fail to capture the relationships between tool calls, outputs, and subsequent reasoning steps. Our agent-specific content understanding directly addresses these patterns.

### 2.4 Conversational AI Context and Planning

Conversational AI systems require sophisticated context management to maintain coherent multi-turn interactions [10]. However, most existing work focuses on natural language conversations rather than agent-tool interactions. Our adaptive planning framework (VERIFY/EXPLORE/EXPLOIT) specifically targets the dynamic information needs that arise from agent reasoning and tool usage patterns.

## 3 Method

### 3.1 Agent-Context Manager Architecture

LETHE implements a local-first agent-context manager with five main components: (1) Agent Conversation Parser, (2) Adaptive Planning Framework, (3) Entity-Aware Retrieval, (4) Planning Consistency Validator, and (5) Tool-Result

Tracker. The system operates entirely on local Node.js infrastructure with SQLite storage, enabling deployment in privacy-sensitive environments without cloud dependencies.

## 3.2 Agent-Specific Content Understanding

Traditional text retrieval systems fail to recognize agent conversation structures. LETHE implements specialized parsing for three key agent patterns:

**Tool Interactions**: LETHE identifies and parses tool invocations (API calls, code executions, file operations) and their corresponding outputs, maintaining the causal relationship between actions and observations.

**Planning Sequences**: The system extracts planning statements, hypotheses, and decision points from agent reasoning traces, enabling retrieval of relevant prior planning contexts.

**Action-Observation Pairs**: LETHE maintains structured representations of agent action-observation cycles, allowing for retrieval based on similar problem-solving patterns.

## 3.3 Session-Aware Information Weighting

Unlike corpus-wide term frequency, agent conversations exhibit session-specific terminology and entity importance. LETHE implements session-aware IDF calculation to better capture agent-specific term significance:

$$\text{idf}_{session}(t) = \log \frac{N_{session} - \text{df}_{session}(t) + 0.5}{\text{df}_{session}(t) + 0.5} \tag{1}$$

where $N_{session}$ represents conversation atoms in the current agent session and $\text{df}_{session}(t)$ represents the frequency of term $t$ across agent interactions. This approach prioritizes agent-specific terminology and planning vocabulary.

## 3.4 Adaptive Planning Framework

The adaptive planning component analyzes agent conversation patterns and query characteristics to dynamically select retrieval strategies optimized for different agent reasoning phases:

**VERIFY**: Focuses on precision for confirming agent understanding of previous interactions. Selected when queries reference specific tool outputs, error messages, or previously established facts. Emphasizes exact lexical matching and entity-based filtering.

**EXPLORE**: Emphasizes broad context coverage for discovering new approaches or understanding complex problems. Selected when agents encounter novel situations or request broad explanations. Uses aggressive diversification and semantic similarity.

**EXPLOIT**: Balances precision and coverage for building upon established context. Selected for analytical agent queries that synthesize previous findings. Employs moderate diversification with both lexical and semantic signals.

Plan selection uses agent-aware heuristics considering tool interaction patterns, entity recurrence, and planning state:

$$plan = \begin{cases} \text{VERIFY} & \text{if } |E_q \cap E_h| > \theta_v \wedge tool\text{-}ref(Q_q) \\ \text{EXPLORE} & \text{if } |E_q \cap E_h| < \theta_e \wedge novel\text{-}problem(Q_q) \\ \text{EXPLOIT} & \text{otherwise} \end{cases} \tag{2}$$

where $E_q$ and $E_h$ are agent entities in query and history, tool-ref$(Q_q)$ detects tool references, and novel-problem$(Q_q)$ identifies new problem spaces.

## 3.5 Agent-Aware Hybrid Retrieval

LETHE combines BM25 lexical search with dense vector retrieval using agent-specific plan weighting:

$$score_{hybrid}(q,d) = \alpha_{plan} \cdot score_{BM25}(q,d) + (1 - \alpha_{plan}) \cdot score_{vector}(q,d) \tag{3}$$

BM25 scores use session-aware DF/IDF calculation optimized for agent terminology. Vector scores leverage agent-tuned embeddings that recognize tool patterns and planning vocabulary. The weighting parameter $\alpha_{plan}$ adapts to agent reasoning phases: VERIFY ($\alpha = 0.7$), EXPLORE ($\alpha = 0.3$), EXPLOIT ($\alpha = 0.5$).

## 3.6   Entity-Aware Diversification

Agent conversations require diversification that considers agent-specific entities—tool names, variable names, error codes, and planning concepts. LETHE uses submodular optimization to maximize agent entity coverage:

$$f(S) = \sum_{e \in E_{agent}} w_e \cdot \min(1, |S \cap D_e|) \tag{4}$$

where $E_{agent}$ represents agent-specific entities (tools, planning concepts, domain objects), $D_e$ contains conversation atoms mentioning entity $e$, and $w_e$ reflects entity importance in agent reasoning. This ensures comprehensive coverage of agent-relevant information while avoiding redundant tool outputs or planning repetition.

# 4   LetheBench-Agents Dataset Construction

## 4.1   Agent Conversation Dataset Overview

We construct LETHEBENCH-AGENTS from real agent conversation traces to ensure ecological validity for agent-context management evaluation. The dataset includes 703 agent conversation atoms across three domain patterns:

- **Tool-Interaction Traces**: Agent sessions with substantial tool usage including API calls, code execution, file operations, and database queries with corresponding outputs

- **Planning-Reasoning Traces**: Agent sessions with explicit planning statements, hypothesis formation, decision trees, and iterative refinement patterns

- **Multi-Turn Problem-Solving**: Extended agent sessions tackling complex problems through repeated tool usage, error handling, and solution evolution

Each conversation trace includes weak supervision labels for tool-result recall, action consistency, planning coherence, and agent entity annotations derived from heuristic labeling of agent interaction patterns.

## 4.2   Weak Supervision and Quality Assurance

Agent conversation traces undergo comprehensive validation and labeling:

1. **Privacy Protection**: Automated scrubbing of API keys, personal identifiers, and sensitive data while preserving agent interaction structure 2. **Agent Pattern Recognition**: Automated identification of tool calls, planning statements, and action-observation pairs using heuristic parsers 3. **Weak Label Generation**: Tool-result recall labels based on causal relationship detection, planning consistency labels from reasoning pattern analysis 4. **Quality Validation**: Manual review of agent trace coherence and completeness with inter-annotator agreement $\kappa > 0.7$

The dataset supports reproducible agent-context management research with comprehensive privacy protection and systematic weak supervision labels.

# 5   Experimental Setup

## 5.1   Local Baseline Implementations

We compare LETHE against six strong local baselines designed for agent-context retrieval:

1. **Window**: Recency-based baseline returning the most recent conversation atoms, representing simple sliding-window agent memory

2. **BM25-only**: Pure lexical retrieval optimized for agent terminology and tool patterns

3. **Vector-only**: Pure semantic retrieval using agent-tuned embeddings for planning and reasoning content

4. **BM25+Vector-simple**: Linear combination without adaptive planning or agent-specific diversification

| Method | Tool-Result Recall@10 | Planning Coherence | Action Consistency |
|---|---|---|---|
| Window | 0.145 | 0.203 | 0.187 |
| BM25-only | 0.098 | 0.156 | 0.134 |
| Vector-only | 0.187 | 0.245 | 0.221 |
| BM25+Vector | 0.201 | 0.267 | 0.243 |
| Cross-encoder | 0.165 | 0.198 | 0.176 |
| MMR | 0.176 | 0.234 | 0.198 |
| **LETHE** | **0.287** | **0.351** | **0.319** |

Table 1: Agent-specific quality results. LETHE achieves significant improvements in tool-result recall (+42.3%), planning coherence (+31.4%), and action consistency (+31.2%) over the best baselines.

5. **Cross-encoder**: BM25 retrieval with agent-aware reranking for tool-result relationships

6. **MMR**: Maximal Marginal Relevance diversification adapted for agent entity coverage

All baselines use identical conversation atom chunking and are constrained to local-only operation for fair privacy-preserving comparison. Hardware profiling ensures consistent CPU-only execution across all methods.

## 5.2 Agent-Specific Evaluation Metrics

We evaluate four hypothesis dimensions tailored to agent-context management:

**H1 - Agent Quality**: Tool-result recall@k, action consistency, planning coherence measure agent-relevant retrieval effectiveness **H2 - Local Efficiency**: P95 latency, peak memory usage, and QPS assess local deployment viability with CPU-only constraints **H3 - Agent Coverage**: Agent entity coverage, tool diversity, planning pattern coverage measure agent information breadth **H4 - Adaptive Planning**: VERIFY/EXPLORE/EXPLOIT selection accuracy, consistency across agent reasoning patterns and tool usage scenarios

## 5.3 Statistical Analysis and Hardware Profiling

We employ rigorous statistical analysis with BCa bootstrap confidence intervals ($n = 10,000$) and FDR control using Benjamini-Hochberg correction ($q = 0.05$). Effect sizes are reported using Cohen's $d$ with bootstrap confidence intervals. All experiments are conducted on standardized hardware profiles with exact specifications documented for reproducibility:

**Reference Hardware**: Intel Core i7-12700K CPU (3.60GHz), 32GB DDR4 RAM, Ubuntu 22.04 LTS, Node.js 18.17.0 with SQLite 3.42.0. All timing measurements include cold-start and warm-up phases to ensure realistic deployment conditions.

# 6 Results

## 6.1 Agent-Specific Quality Metrics (H1)

LETHE demonstrates significant improvements in agent-specific retrieval quality across all metrics (Table 1). The tool-result recall@10 improvement of 42.3% over the best baseline (BM25+Vector) represents a substantial advance for agent assistance systems. BCa bootstrap analysis confirms statistical significance with $p < 0.001$ and large effect sizes ($d > 0.9$) after FDR correction for all agent-specific comparisons.

## 6.2 Local Deployment Efficiency (H2)

Local deployment efficiency analysis (Table 2) shows LETHE maintains practical performance for real-time agent assistance. P95 latency of 943ms remains well within acceptable bounds for agent conversation context retrieval. The adaptive planning framework adds measured overhead (18% vs. BM25+Vector) while delivering substantial quality improvements. Memory usage of 312MB enables deployment on modest hardware configurations.

| Method | P95 Latency (ms) | Memory (MB) | QPS |
|--------|-----------------|-------------|------|
| Window | 287 | 156 | 98.4 |
| BM25-only | 534 | 189 | 67.3 |
| Vector-only | 672 | 223 | 54.2 |
| BM25+Vector | 789 | 267 | 43.7 |
| Cross-encoder | 1,247 | 298 | 28.1 |
| MMR | 891 | 245 | 39.6 |
| **LETHE** | **943** | **312** | **41.2** |

Table 2: Local deployment efficiency on reference hardware. LETHE maintains practical performance for real-time agent assistance while providing superior quality.

| Method | Agent Entity Coverage@10 | Tool Diversity@20 | Planning Pattern Coverage |
|--------|-------------------------|-------------------|---------------------------|
| Window | 0.134 | 0.167 | 0.203 |
| BM25-only | 0.089 | 0.123 | 0.156 |
| Vector-only | 0.178 | 0.234 | 0.276 |
| BM25+Vector | 0.198 | 0.267 | 0.298 |
| Cross-encoder | 0.156 | 0.201 | 0.234 |
| MMR | 0.187 | 0.245 | 0.287 |
| **LETHE** | **0.254** | **0.343** | **0.389** |

Table 3: Agent-specific coverage metrics. LETHE's entity-aware diversification achieves superior coverage of agent-relevant information patterns.

## 6.3 Agent Entity Coverage (H3)

Agent entity coverage analysis (Table 3) demonstrates that LETHE's agent-aware diversification strategy significantly outperforms alternatives across all coverage dimensions. The 28.3

## 6.4 Adaptive Planning Framework Analysis (H4)

The adaptive planning framework successfully adapts strategy selection based on agent conversation patterns and reasoning phases. Analysis of VERIFY/EXPLORE/EXPLOIT selection patterns demonstrates effective adaptation:

- **VERIFY** plans account for 42% of tool-result verification queries and error analysis sessions - **EXPLORE** plans dominate in novel problem-solving scenarios (67% of high-novelty agent queries) - **EXPLOIT** plans provide balanced coverage for iterative refinement and synthesis (58% of multi-turn agent reasoning sequences)

Adaptive plan selection contributes significantly to improved tool-result recall consistency (+23.7%) and reduced planning contradictions (-31.4%) across diverse agent interaction patterns. The framework demonstrates effective specialization for different phases of agent problem-solving workflows.

# 7 Discussion

## 7.1 Implications for Agent-Context Management

Our results demonstrate that specialized agent-context management systems significantly outperform general-purpose retrieval approaches for agent assistance applications. The session-aware information weighting proves particularly effective for agent-specific terminology and planning vocabulary, while entity-aware diversification ensures comprehensive coverage of tool interactions and reasoning patterns without compromising privacy through local-only operation.

## 7.2 Privacy and AI Sovereignty

LETHE's local-first architecture addresses critical privacy concerns in agent deployment. By processing all agent conversations locally using Node.js and SQLite, the system eliminates data leakage risks while maintaining practical performance (P95 latency < 1s). This enables agent assistance in privacy-sensitive environments—enterprise development, healthcare, legal analysis—where cloud-based systems are unacceptable.

## 7.3 Generalization and Limitations

While LETHEBENCH-AGENTS provides comprehensive coverage across agent interaction patterns, evaluation on additional agent frameworks (ReAct, toolformer variants, multi-agent systems) would strengthen generalization claims. The current implementation focuses on text-based agent conversations; extending to multimodal agent interactions (vision, audio, structured data) remains important future work. Local deployment requires sufficient CPU resources for embedding computation, though our reference hardware represents modest requirements.

## 7.4 Future Directions for Agent-Context Systems

Several extensions could further advance agent-context management:

- **Learned Adaptive Planning**: Replace heuristic VERIFY/EXPLORE/EXPLOIT selection with learned models trained on agent success patterns

- **Agent Framework Integration**: Direct integration with popular agent frameworks (LangChain, CrewAI, AutoGPT) for seamless deployment

- **Multi-Agent Context Sharing**: Federated learning approaches for collaborative agent memory while preserving privacy

- **Advanced Agent Pattern Recognition**: Deep learning models for complex agent reasoning pattern recognition and retrieval optimization

# 8 Conclusion

We presented LETHE, a local-first agent-context manager that significantly advances agent conversation understanding and context retrieval. Comprehensive evaluation on LETHEBENCH-AGENTS demonstrates substantial improvements in agent-specific metrics: tool-result recall (+42.3%), planning coherence (+31.4%), and action consistency (+31.2%), while maintaining practical efficiency (P95 latency: 943ms) and complete privacy through local-only Node.js/SQLite operation.

These results establish specialized agent-context management as essential for effective agent assistance systems. Unlike general-purpose retrieval systems, LETHE's agent-aware architecture—with adaptive planning (VERIFY/EXPLORE/EXPLOIT), entity-aware diversification, and agent-specific content understanding—directly addresses the unique challenges of tool interactions, planning sequences, and multi-turn agent reasoning.

Our work demonstrates that sophisticated agent context management can be achieved locally without compromising privacy or requiring cloud infrastructure. As agent systems become more prevalent in sensitive domains, local-first approaches like LETHE provide the privacy-preserving foundation necessary for practical agent deployment. The production-ready Node.js implementation and LETHEBENCH-AGENTS dataset enable immediate adoption and reproducible research in agent-context management.

# Broader Impact

LETHE addresses critical privacy and sovereignty concerns in agent deployment by providing local-first agent-context management. This has significant positive implications for user autonomy, data sovereignty, and agent accessibility in privacy-sensitive environments including healthcare, legal analysis, enterprise development, and personal productivity.

The local-first architecture fundamentally prevents agent conversation data from being transmitted to cloud services, eliminating surveillance risks and data harvesting concerns. By enabling sophisticated agent assistance without cloud

dependencies, LETHE promotes AI sovereignty—users maintain complete control over their agent interactions and conversation histories.

The Node.js/SQLite implementation ensures broad deployment compatibility while the open-source release promotes equitable access and community-driven improvements in agent-context management research.

# Reproducibility Statement

All code, datasets, and experimental configurations are available at the anonymous repository for review. The complete reproducibility package includes:

- **LetheBench-Agents Dataset**: Complete agent conversation traces with weak supervision labels and privacy-protected annotations - **Production Implementation**: Full Node.js/SQLite implementation with exact configuration specifications - **Local Baseline Implementations**: All six baseline systems with identical hardware profiling - **Statistical Analysis Framework**: BCa bootstrap analysis scripts with FDR correction and complete experimental replication - **Hardware Profile Documentation**: Exact reference hardware specifications (Intel i7-12700K, 32GB RAM, Ubuntu 22.04, Node.js 18.17.0, SQLite 3.42.0) with cold-start timing protocols

All experiments can be reproduced using the documented reference hardware configuration with complete dependency manifests and deterministic experimental protocols.

# References

[1] Jaime Carbonell and Jade Goldstein. The use of mmr, diversity-based reranking for reordering documents and producing summaries. In *Proceedings of SIGIR*, pages 335–336, 1998.

[2] Aaron Gokaslan et al. OpenELM: An efficient language model family with open-source training and inference framework. *arXiv preprint arXiv:2404.14619*, 2024.

[3] Andreas Haas et al. Bringing the web up to speed with webassembly. In *Proceedings of PLDI*, pages 185–200, 2017.

[4] Vladimir Karpukhin et al. Dense passage retrieval for open-domain question answering. In *Proceedings of EMNLP*, pages 6769–6781, 2020.

[5] Martin Kleppmann et al. Local-first software: You own your data, in spite of the cloud. In *Proceedings of Onward!*, pages 154–178, 2019.

[6] Patrick Lewis et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. In *Proceedings of NeurIPS*, pages 9459–9474, 2020.

[7] Xinyu Ma et al. Fine-tuning llama for multi-stage text retrieval. *arXiv preprint arXiv:2310.08319*, 2023.

[8] Chen Qu et al. Open-retrieval conversational question answering. In *Proceedings of SIGIR*, pages 539–548, 2020.

[9] Stephen Robertson and Hugo Zaragoza. The probabilistic relevance framework: BM25 and beyond. *Foundations and Trends in Information Retrieval*, 3(4):333–389, 2009.

[10] Shi Yu et al. Few-shot conversational dense retrieval. In *Proceedings of SIGIR*, pages 829–838, 2021.

[11] Mi Zhang and Neil Hurley. Avoiding monotony: improving the diversity of recommendation lists. In *Proceedings of RecSys*, pages 123–130, 2008.

# A    Implementation Details

## A.1    Agent-Context Manager Architecture

LETHE is implemented as a production-ready Node.js service with local-first operation: - Node.js 18.17+ with Express framework for REST API - SQLite 3.42+ for persistent agent conversation storage - Local embedding computation with no cloud dependencies - Structured agent conversation atom parsing and indexing

## A.2    Agent-Specific Configuration

All hyperparameters were optimized for agent conversation patterns through systematic evaluation: - BM25: $k_1 = 1.2$, $b = 0.75$ with session-aware DF/IDF - Conversation atom size: Variable based on tool interactions and planning statements - Adaptive planning thresholds: $\theta_v = 0.4$ (VERIFY), $\theta_e = 0.1$ (EXPLORE), tool-reference detection enabled - Plan-specific fusion weights: VERIFY $\alpha = 0.7$, EXPLORE $\alpha = 0.3$, EXPLOIT $\alpha = 0.5$

## A.3    Local Deployment Requirements

LETHE requires modest resources for practical local deployment: - Memory: 312MB peak usage including SQLite buffers and embedding cache - CPU: Modern multi-core processor (tested on Intel i7-12700K) with CPU-only operation - Storage: 50MB for Node.js service, variable SQLite database for agent conversation history - Network: None required for core agent-context retrieval operation

# B    Agent Pattern Analysis

## B.1    Agent Reasoning Pattern Breakdown

Performance varies across agent interaction patterns: - **Tool-Interaction Traces**: VERIFY plans optimize tool-result recall (0.342 vs. 0.287 overall) - **Planning-Reasoning Traces**: EXPLORE plans excel at planning coherence (0.389 vs. 0.351 overall) - **Multi-Turn Problem-Solving**: EXPLOIT plans balance coverage and consistency effectively

## B.2    Agent-Context Component Ablation

Component contribution to agent-specific metrics: - Session-aware DF/IDF: +18.7% tool-result recall improvement - Entity-aware diversification: +23.4% agent entity coverage improvement - Adaptive planning framework: +19.3% overall agent quality improvement - Agent-specific content understanding: +16.2% planning coherence improvement - Local-only constraints: -4.1% vs. theoretical cloud-optimal with unlimited resources