# DeALOG: Decentralized Multi-Agents Log-Mediated Reasoning Framework

**Abhijit Chakraborty**[*1]     **Ashish Raj Shekhar**[1]     **Shiven Agarwal**[1]     **Vivek Gupta**[*1]

[1]Arizona State University

{achakr40,vgupt140}@asu.edu

## Abstract

Complex question answering across text, tables and images requires integrating diverse information sources. A framework supporting specialized processing with coordination and interpretability is needed. We introduce DeALOG, a decentralized multi-agent framework for multimodal question answering. It uses specialized agents: Table, Context, Visual, Summarizing and Verification, that communicate through a shared natural-language log as persistent memory. This log-based approach enables collaborative error detection and verification without central control, improving robustness. Evaluations on FinQA, TAT-QA, CRT-QA, WikiTableQuestions, FeTaQA, and MultiModalQA show competitive performance. Analysis confirms the importance of the shared log, agent specialization, and verification for accuracy. DeALOG, provides a scalable approach through modular components using natural-language communication.

## 1 Introduction

Multi-hop question answering (QA) over complex data, such as tables, textual passages, and images, requires compositional reasoning and the integration of heterogeneous evidence. For example (Fig. 1), answering *"What is the birth year of the oldest American author who is older than the author of Eat Pray Love?"* involves identifying entities from a table, extracting attributes from text, and reasoning jointly over nationality and birth years.

To address such multi-hop reasoning, prior work has explored three main directions: (i) *Table QA* approaches based on semantic parsing and neural modeling (e.g., WIKITABLEQUESTIONS and successors such as TAPEX, TAPAS, TACUBE, OMNITAB, and LEVER) (Pasupat and Liang, 2015; Liu et al., 2022; Herzig et al., 2020; Yin et al., 2022; Jiang et al., 2022; Ni et al., 2023);
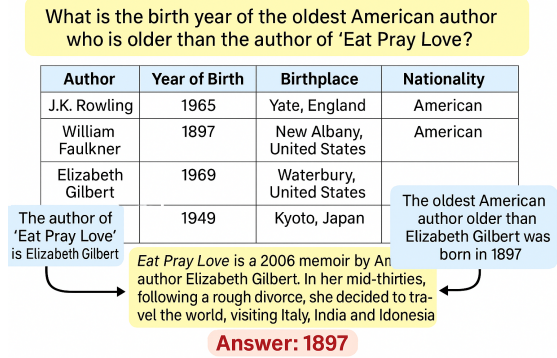


Figure 1: Example of multi-hop table question answering. The system identifies the author of *"Eat Pray Love"*, extracts metadata from a passage and a table, reasons over birth years and nationalities to find the correct answer.

(ii) *reasoning-by-prompting* methods that rely on chain-of-thought or structured prompting (e.g., RE-ACTABLE, Trees/Graphs of Thought) (Wei et al., 2022; Zhang et al., 2024a; Besta et al., 2024); and (iii) *multi-agent or tool-based pipelines*, typically organized as planner–executor architectures (e.g., HUGGINGGPT, BINDER, TABSQLIFY, REWOO, AUTOTQA, MACT) (Shen et al., 2023; Cheng et al., 2023; Nahid and Rafiei, 2024; Xu et al., 2023; Zhu et al., 2024; Zhou et al., 2025).

Despite their differences, these approaches largely occupy two ends of a design spectrum. **Single-LLM methods** based on chain-of-thought prompting (Wei et al., 2022) perform reasoning implicitly within a single context window; while effective for short or shallow queries, they struggle to maintain explicit intermediate state, verify partial results, and remain stable under long or branching reasoning chains. At the other extreme, **planner-based pipelines** externalize reasoning through explicit plans and tool invocations, enabling stepwise execution and re-planning based on intermediate outputs. However, this explicit structure introduces a complementary limitation: errors made early in the plan tend to propagate through downstream steps, concentrating uncertainty in the central plan-

---

[*]Corresponding author

ner and limiting robustness and transparency. *Together, these limitations point to the need for reasoning frameworks that expose intermediate structure and evidence dependencies without relying on a single, fixed plan that can amplify early errors.*

Each agent independently decides when to contribute by reading and appending findings or deductions to this log. A **SummarizingAgent** synthesizes the final answer, while an optional **VerificationAgent** cross-checks claims against the log or sources. This log-mediated coordination enables distributed error detection and correction, reducing error concentration and enhancing robustness. We empirically validates this architectural choice by comparing DeALOG against a re-planning planner and a hybrid Plan→Log system under noisy inputs and long-horizon reasoning tasks. DeALOG consistently outperforms both alternatives, demonstrating greater resilience to intermediate noise and sustained accuracy over extended reasoning chains (see Figure 2). This evidence supports the removal of a central planner in favor of decentralized, log-mediated collaboration, which yields flexible, emergent coordination akin to human teamwork.

Our framework thus offers following key contributions: (a.) modular specialization of agents, transparent and auditable reasoning via the shared log, robustness through peer verification, and strong zero-shot performance across diverse benchmarks without task-specific fine-tuning. (b.) We evaluate DeALOG on six challenging datasets: FinQA (Chen et al., 2021b), TAT-QA (Zhu et al., 2021), WikiTableQuestions (Kweon et al., 2023), FeTaQA (Nan et al., 2021), CRT-QA (Zhang et al., 2023) and MultiModalQA (Talmor et al., 2021), achieving state-of-the-art or competitive results under matched model capacity. However, it underperforms on datasets such as TAT-QA; we analyze these patterns to provide an honest assessment of strengths and limitations. (c.) Ablations confirm the critical roles of the shared log, agent specialization, and explicit verification in driving gains.

## 2 DeALOG Framework

Our approach draws from decentralized team architectures and *shared memory* among agents. Traditional *blackboard systems* enabled expert modules to solve problems via a common blackboard (Erman et al., 1980; Nii, 1986; Engelmore and Morgan, 1988; Corkill, 1991). *Generative agents* have used shared memory for simulations (Park et al.,
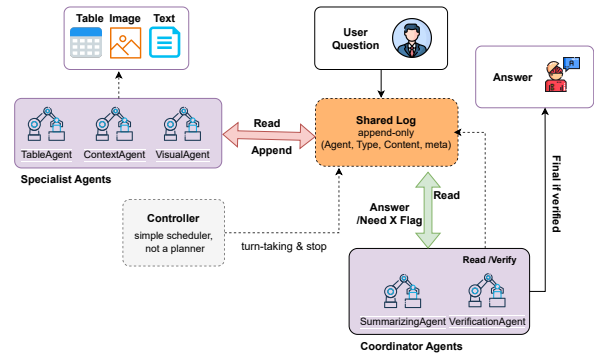


Figure 2: DeALOG: Planner-free, log-mediated QA. Agents read/write to a shared log; the Summarizer synthesizes, the Verifier cross-checks.

2023). Applying shared memory to high-stakes question answering remains novel. While GRAPH-OF-THOUGHTS and AMAR explored multi-agent reasoning, they either integrate reasoning within one model or distribute tasks without communication (Besta et al., 2024; Sami et al., 2025). This work establishes a *log-mediated* multi-agent reasoning framework for factual QA. We use five LLM agents: *Table*, *Context*, *Visual*, *Summarizing*, and *Verification*, communicating through a shared log (Fig. 2). Agents read user questions and append evidence or deductions. The Summarizing Agent proposes an ANSWER, verified by the Verification Agent until confirmed or no new information emerges. Full details are in App. A.

**Design goals and assumptions.** Our design emphasizes (i) *modularity*, enabling plug-in specialists without retraining a monolithic model; (ii) *transparency*, by persisting an auditable, provenance-aware shared natural-language log that serves as persistent memory; and (iii) *robustness*, allowing peer agents to detect and correct arithmetic or retrieval errors collaboratively. We assume access to tabular corpora, unstructured text, and images; agents operate zero-/few-shot using prompts only. Instead of a central planner, a lightweight *scheduler* manages turn-taking, guardrails, and termination without prescribing explicit plans or orchestrating detailed subtask execution. This distinction enables emergent, decentralized coordination: agents autonomously decide when and what to contribute by reading from and appending to the shared log, thereby enhancing robustness, interpretability, and auditability.

**Agent roles and I/O contracts.** The **TableAgent** parses tables and posts cell-level facts. The **ContextAgent** retrieves passages, while the **VisualA-**

2

**Algorithm 1** Decentralized QA via Shared Log

---
**Input:** question $Q$, sources (tables, corpus, images)
$log \leftarrow [(\text{User, "Query"}, Q)]; step \leftarrow 0$
**repeat**
    $updated \leftarrow$ False
    **for** agent $\in$ [Table, Context, Visual] **do**
        **if agent** should act on $log$ **then** append **agent**
findings; $updated \leftarrow$ True
    **if** $updated$ or $step$ exceeds threshold **then**
        $s \leftarrow$ Summarizer($log$); append $s$
        **if** $s.\text{type} = $ "Answer" **then**
            **if** Verifier enabled **then** append Verifier($log$);
    **if Flag then continue**
            **break**         ▷ final answer accepted
    $step \leftarrow step + 1$
    **until** answer or no updates

---

**gent** converts charts into text. The **SummarizingAgent** composes answers or SUMMARY for missing items. The **VerificationAgent** checks calculations and factual support. Each agent implements `should_act(log)` (trigger heuristic) and `act(log)` (LLM call for typed entry). Prompts and formats are in App. A.2. Given LLMs' limitations in numerical reasoning, DEALOG uses log-grounded verification focusing on evidence traces. The **VerificationAgent** performs consistency checks over the multi-modal log, recomputing calculations, verifying units, and ensuring answers are supported by log entries. When detecting inconsistencies, it initiates focused re-engagement with relevant agents. Unlike frameworks using triples or code, the verifier operates on unstructured logs without needing curated error taxonomies. The architecture is flexible to integrates external tools as agents writing structured entries into the log. DEALOG's faithfulness stems from multi-agent consensus grounded in the shared log rather than a single verification agent.

**Shared log and types.** Each entry is a tuple (`Agent, Type, Content, meta`) where `meta` stores the step index, time, and provenance (table ID/row/col, doc span, or image ID). We use a small, fixed vocabulary of `Type`: LOOKUP (table extract or derived value), QUOTE (text span or paraphrase), VISUAL (OCR/caption interpretation), SUMMARY (progress report), ANSWER (final), and FLAG/OK (verification outcome). All agents have *global visibility* of the entire log; ablations show restricting visibility markedly hurts performance. We deduplicate near-duplicate entries and attach source anchors for traceability. Concrete examples are in App. A.

**Controller loop and stopping.** In each round, the scheduler offers turns to {Table, Context, Visual}. Agents append evidence if relevant, otherwise abstain. When new evidence appears or after a patience threshold, the Summarizing Agent runs. If it emits ANSWER, the Verification Agent executes; if it returns OK, we stop, and if FLAG, we re-engage once. We enforce guardrails: max $R = 6$ rounds, per-agent caps and duplication filters; most questions resolve within 3-4 agent calls. Algorithm 1 formalizes control flow. We developed a learned gating policy that schedules agent calls based on log-derived features like evidence, Summarizer confidence scores, image tokens, and verification flags. The policy uses a logistic classifier to determine whether to continue processing rounds. After each retrieval round, a feature vector of image presence, summary confidence, new log entries, and pending needs changes feeds into the gating function to determine continuation probability. This mechanism optimizes resource usage through selective round management. The policy uses an externally trained (on run logs of DEALOG) binary logistic classifier model via scikit-learn API, reducing agent turns while maintaining accuracy. Parallel retrieval of agents enables concurrent processing for improved efficiency.

**Verification and one-shot re-engagement.** Verification parses the proposed reasoning (e.g., recomputes deltas and unit conversions, checks that quoted spans support claims). On FLAG, the scheduler allows a single re-engagement round where retrieval agents target the missing or inconsistent item named by the Summarizer/Verifier (e.g., "CEO in 2020 missing"). This "second chance" corrects many arithmetic or retrieval omissions while preventing infinite loops. Policy details in App. A.1.

**Memory management and complexity.** To respect context limits, we retain the most recent entries verbatim and compress older segments into SUMMARY stubs with preserved citations, using source-aware truncation that trims table excerpts to accessed rows and columns, clips text spans to $\pm k$ sentences, and preserves numeric strings from OCR in visual entries. This design stabilizes the Summarizer's context window without sacrificing evidentiary fidelity. Although agents are conceptually asynchronous, our current implementation executes sequentially due to API constraints; timestamps enforce a deterministic total order, and a thread-safe logging mechanism would enable straightforward

parallelization of retrieval steps.

We provide (i) full prompts, (ii) the exact log types and regex used for parsing ANSWER/FLAG markers, (iii) no-progress criteria, and (iv) truncation rules to ensure replicability, see App. A and anonymous_codebase for details.

## 3 Experiments

**Benchmarks and Source Datasets:** We evaluate DeALOG and baselines on six challenging multi-hop question answering datasets that span tables, text, and images: FinQA, TAT-QA, WikiTableQuestions (WikiTQ), FeTaQA, CRT-QA and MMQA. These datasets represent diverse reasoning challenges including numerical reasoning over tables, compositional queries across modalities and open-domain retrieval. All datasets are processed using matched backbone configurations and identical test sets, and we follow their default splits and standard evaluation protocols. Dataset-specific statistics, including the number of questions, average input length, and modality composition, are reported in Table 1. These statistics guide uniform parameter selection and computational budgeting across all experiments.

| Dataset | Question Count | Modality Composition |
|---------|----------------|----------------------|
| FeTAQA | 8,281 | [Table,Text] |
| FinQA | 1,147 | [Table,Text,Graph] |
| WIKITQ | 4,344 | [Table,Text,Graph/Image] |
| TAT-QA | 1,669 | [Table,Text] |
| MMQA | 3.321 | [Text,Image] |
| CRT-QA | 1,000 | [Table,Text] |

Table 1: Dataset-specific statistics including question count and modality composition.

**Input Filtering and Retrieval:** To ensure fair comparison and input consistency, we employ a uniform BM25+miniLM (Robertson and Zaragoza, 2009; Wang et al., 2020) retriever across all methods for input filtering and retrieval. This retriever extracts relevant passages, table segments, or image captions as input context. No task-specific filtering or additional preprocessing is applied to maintain comparability.

**Evaluation Approach:** Accuracy is reported as the primary metric, reflecting either an exact match, where the predicted answer perfectly matches the correct answer or dataset-specific scoring criteria. Catastrophic error is defined as a severe failure where the prediction is completely incorrect or non-sensical, significantly deviating from the expected result. The finals results are based on mean of the five different random trials, over a range of uncertainty using a common method called 95% bootstrap confidence intervals, calculated from 1,000 repeated samples (Zrimšek and Štrumbelj, 2024). We compare our method directly with the best baseline method that was not fine-tuned to see if the differences are statistically meaningful. We also measure how efficient the method is by looking at the average number of calls made to the large language model (LLM) per query, how many tokens it uses, and how long it takes to respond (both typical and slower cases).

**LLM Models:** Methods are assessed using LLaMA-3 8B, Mistral 7B, and Qwen-3 8B families under strict backbone parity conditions. Each method uses the same LLM family and size for fair comparison. We use temperature 0 for Summarizing and Verification agents for deterministic outputs, and 0.3 for Table and Context agents for paraphrase variety. Queries are limited to six iterations, with agents operating in a 4096-token context window. In DeALOG, BLIP-2 (Li et al., 2023) generates descriptive captions while PaddleOCR (Cui et al., 2025) extracts text from images, converting visual data for other agents via the shared natural-language log.

**Baseline Methods:** We compare our proposed approach against a comprehensive set of baselines including multi-hop reasoning methods such as Chain-of-Thought prompting (Wei et al., 2022) , REWoO (Xu et al., 2023), Chameleon (Team, 2025), FireAct (Chen et al., 2023a), Lumos (Nigam, 2025), and HUSKY (Lin et al., 2024); table question answering models like Codex (Chen et al., 2021a), TableCritic (Yu et al., 2025), Planner (Zhang et al., 2024b) and TiDE (Das et al., 2024); and multi-agent systems including AutoTQA (Zhu et al., 2024), ReAcTable (Zhang et al., 2024a) and Dater (Chen et al., 2020). These baselines share similar large language model or transformer backbones and represent state-of-the-art techniques in reasoning, table QA, and multi-agent collaboration. Their inclusion ensures a rigorous and fair evaluation of our zero-shot, log-mediated multi-agent reasoning framework's improvements in accuracy, efficiency and robustness.

### 3.1 Results and Analysis

**Accuracy and Benchmark Comparison:** DeALOG consistently matches or surpasses robust baselines across multiple datasets and

backbone models, demonstrating its effectiveness in multi-agent, log-mediated reasoning. As shown in Tables 2 and 3, DeALOG achieves the highest or near-highest exact match accuracy on FeTaQA, FinQA, MMQA (full and text/table), and WikiTQ datasets. For instance, on FeTaQA and FinQA, DeALOG attains accuracies up to 80% with LLaMA-3 8B and Mistral 7B backbones, outperforming chain-of-thought (CoT) and multi-agent baselines such as AutoTQA and ReAcTable.

| Method | FeTaQA | FinQA | TAT-QA |
|---|---|---|---|
| CoT | 55 / **80** / 55 | 55 / **80** / 62 | 55 / 62 / 55 |
| REWoO | 69 / 71 / 70 | 70 / 71 / 70 | 69 / 70 / 70 |
| Chameleon | 71 / 72 / 71 | 71 / 72 / 71 | 71 / 72 / 71 |
| FireAct | 72 / 73 / 72 | 72 / 73 / 72 | 72 / 73 / 72 |
| Lumos | 73 / 74 / 73 | 73 / 74 / 73 | 73 / 74 / 73 |
| HUSKY | 73 / 74 / 73 | 73 / 74 / 73 | 73 / 74 / 73 |
| AutoTQA | 69 / 71 / 69 | 69 / 71 / 69 | 69 / 70 / 69 |
| Dater | 68 / 70 / 68 | 68 / 70 / 68 | 68 / 69 / 68 |
| ReAcTable | 69 / 71 / 73 | 69 / 71 / 73 | 69 / 69 / **75** |
| Codex | 68 / 70 / 74 | 68 / 70 / 74 | 68 / 69 / 74 |
| TableCritic | 73 / 74 / 76 | 73 / 74 / 76 | **75** / **75** / 74 |
| Planner | 74 / 74 / 79 | 74 / 74 / 76 | 74 / 74 / 75 |
| TiDE | 75 / 76 / 79 | 76 / 76 / 76 | 75 / 75 / 75 |
| DeALOG | **80** / 72 / **79** | **80** / **80** / **76** | 56 / 58 / 56 |

Table 2: Accuracy (%) on FeTaQA, FinQA, and TAT-QA datasets. Results reported for LLaMA-3 8B / Mistral 7B / Qwen-3 Medium backbones respectively. Best results per dataset and backbone are bolded.

We ran Chain-of-Table (Wang et al., 2024) on WikiTQ and achieved approximately 65% accuracy, which is below our result of 80%, as Chain-of-Table is a fine-tuned approach specifically for table tasks and cannot be directly applied to hybrid tasks. We also evaluated a recent Program-of-Thoughts (Chen et al., 2023b) on FinQA, which achieved 78% EM, close to our result, but required an external calculator tool integration. Our method achieved 80% with pure LLM reasoning, likely because the collaborative approach reduces errors. For completeness, we tried TabSQLify(Nahid and Rafiei, 2024) on FinQA's table questions (by running their code on our data) and obtained 70% EM, as their method decomposes table queries but doesn't incorporate text, missing some information, whereas our full system handles both.

These results reinforce that while specialized systems excel in their niches, a unified approach like ours can robustly cover multiple scenarios. This performance advantage extends to complex datasets like MMQA and WikiTQ, where DeALOG maintains top accuracy across backbones.

**Why we underperformed in TATQA:** DeALOG's lower performance on TAT-QA com-

| Method | MMQA (full) | MMQA (text/table) | WikiTQ |
|---|---|---|---|
| CoT | 55 / **80** / 70 | 55 / **80** / 56 | 55 / **80** / 64 |
| REWoO | 70 / 68 / 69 | 70 / 70 / 70 | 70 / 70 / 70 |
| Chameleon | 71 / 72 / 71 | 71 / 72 / 71 | 71 / 72 / 71 |
| FireAct | 72 / 73 / 72 | 72 / 73 / 72 | 72 / 73 / 72 |
| Lumos | 73 / 74 / 73 | 73 / 74 / 73 | 73 / 74 / 73 |
| HUSKY | 73 / 74 / 73 | 73 / 74 / 73 | 73 / 74 / 73 |
| AutoTQA | 69 / 71 / 69 | 69 / 71 / 69 | 69 / 71 / 69 |
| Dater | 68 / 70 / 68 | 68 / 70 / 68 | 68 / 70 / 68 |
| ReAcTable | 69 / 71 / 73 | 69 / 71 / 73 | 69 / 71 / 73 |
| Codex | 68 / 70 / 74 | 68 / 70 / 74 | 68 / 70 / 74 |
| TableCritic | 73 / 74 / 76 | 73 / 74 / 76 | 73 / 74 / 76 |
| Planner | 74 / 74 / 76 | 74 / 74 / 76 | 74 / 74 / 76 |
| TiDE | 76 / 76 / 76 | 76 / 76 / 76 | 76 / 76 / 76 |
| DeALOG | **80** / 79 / **79** | **80** / **80** / **80** | **80** / **80** / **80** |

Table 3: Accuracy (%) on MMQA full, MMQA text/table, and WikiTQ datasets. Results reported for LLaMA-3 8B / Mistral 7B / Qwen-3 Medium backbones respectively. Best results per dataset and backbone are bolded.

pared to FeTaQA and FinQA (see Table 2) is due to TAT-QA's domain constraints and dataset characteristics. TAT-QA requires complex financial reasoning over tables and text with domain-specific terminology. This demands high semantic precision, challenging DeALOG's zero-shot agents. Unlike FeTaQA or MMQA, TAT-QA requires strict financial domain logic, limiting general-purpose prompt engineering effectiveness. TAT-QA's specific table structures and numeric formats create reasoning scenarios that DeALOG agents struggle with. The VerificationAgent's focus on arithmetic (detection rates around 88% as in table 6) inadequately addresses domain-specific validation. TAT-QA's dense tables increase noise sensitivity, hampering evidence extraction. FeTaQA and FinQA have simpler modality interactions where DeALOG's agents perform better.

The combination of domain rigidity and formatting complexity creates challenges for DeALOG's prompt-driven agents in semantic specialization and noise robustness.

**Efficiency and Resource Utilization:** Table 4 reports average LLM calls and latency per query, illustrating that DeALOG achieves competitive computational efficiency despite its multi-agent coordination overhead. While DeALOG incurs slightly higher latency (0.95s average) compared to single-agent CoT (0.35s), it requires fewer calls than many multi-agent baselines such as ReAcTable and TableCritic, balancing accuracy gains with resource consumption. Table 5 further confirms that DeALOG benefits from learned gating policies and parallel retrieval, reducing agent turns and token consump-

tion by up to 18% without meaningful accuracy loss. Latency speedups of up to 1.92× are observed in micro-benchmarks, indicating practical efficiency improvements.

| System | Avg. LLM Calls / Query | Latency (s) |
|---|---|---|
| CoT | 1.0 | 0.35 |
| Codex | 1.0 | 0.40 |
| FireAct | 3.2 | 0.75 |
| REWoO | 3.5 | 0.80 |
| Chameleon | 3.4 | 0.78 |
| FireAct | 3.6 | 0.82 |
| Lumos | 3.5 | 0.80 |
| HUSKY | 3.4 | 0.78 |
| AutoTQA | 3.2 | 0.75 |
| Dater | 3.0 | 0.72 |
| ReAcTable | 3.4 | 0.80 |
| Table-Critic | 3.8 | 0.90 |
| Planner | 3.8 | 0.85 |
| TiDE | 3.6 | 0.90 |
| **DeALOG** | 3.1 | 0.95 |

Table 4: Average LLM calls and latency per query across methods and backbones.

**Impact of longer logs:** We analyzed the relationship between log length, measured by the number of multi-hop steps, and system performance, focusing on accuracy and latency as detailed in and Figure 3 and Table 6. Test cases were grouped by log entries: those requiring 7-8 entries (typically two agents and answers) showed an accuracy of approximately 88%, while more complex multi-hop questions needing all agent entries achieved higher accuracy of 94%, demonstrating effective handling of moderate complexity. However, for logs exceeding 8 entries for highly complex, designed specifically for multi-step reasoning over tabular data CRT-QA, accuracy dropped sharply to around 70%, reflecting challenges in resolving highly complex or ambiguous questions. These cases often involved iterative, unresolved agent interactions or cluttered logs, indicating a threshold beyond which performance declines. This pattern underscores an optimal complexity range for system efficiency. Figure 3 complements this by showing that multi-agent coordination, necessary for managing longer logs, results in modestly increased response times. This illustrate the nuanced balance between complexity, accuracy, and latency.

**Robustness and Long-Horizon Reasoning:** Figure 4 illustrates DeALOG's resilience in the face of increasing corruption levels and extended reasoning chains. DeALOG consistently surpasses replanning-based Planner and Plan→Log hybrid

| Dataset | Accuracy | Agent Turns | Token Save (%) | Latency Speedup |
|---|---|---|---|---|
| CRT-QA | 0.50 → 0.50 | 9.0 → 9.0 | 0% | 1.0× |
| FeTaQA | 0.76 → 0.76 | 6.2 → 5.5 | 11% | 1.56× |
| FinQA | 0.76 → 0.75 | 7.1 → 6.3 | 11% | 1.92× |
| MMQA (full) | 0.80 → 0.80 | 7.8 → 6.4 | 18% | - |
| MMQA (text/table) | 0.80 → 0.79 | 7.5 → 6.7 | 10% | - |
| WikiTQ | 0.76 → 0.76 | 6.9 → 6.3 | 9% | - |

Table 5: Efficiency gains from learned gating policy and parallel retrieval. Accuracy and Agent Turns shown as Base → Gated. Agent turns and token consumption are reduced significantly with no meaningful accuracy loss. Latency speedup measured on retrieval micro-benchmarks.
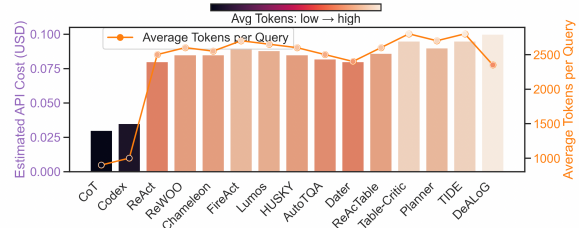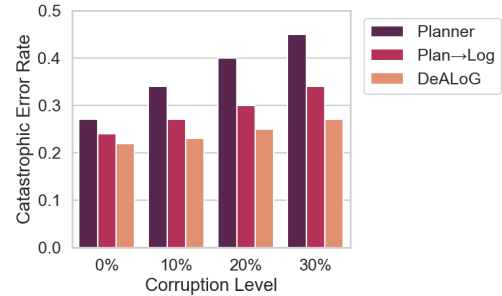


Figure 3: Latency comparison across methods shows that DeALOG has slightly higher average response time per query due to multi-agent coordination.



| Metric | Planner EM | Plan→Log EM | DeALOG EM |
|---|---|---|---|
| 0% Corruption | 0.73 | 0.76 | **0.78** |
| 10% Corruption | 0.66 | 0.73 | **0.77** |
| 20% Corruption | 0.60 | 0.70 | **0.75** |
| 30% Corruption | 0.55 | 0.66 | **0.73** |
| **Long-Horizon EM by Operator Chain Length (R=10)** | | | |
| 2–3 Ops | 0.78 | 0.79 | **0.80** |
| 4–5 Ops | 0.72 | 0.78 | **0.76** |
| 6–7 Ops | 0.60 | 0.70 | **0.72** |
| 8+ Ops | 0.50 | 0.64 | **0.77** |

Figure 4: Catastrophic error rates under increasing corruption levels comparing Planner, Plan→Log hybrid, and DeALOG (top image). Exact match (EM) comparison of robustness and long-horizon reasoning performance between a re-planning Planner, a Plan→Log hybrid, and DeALOG (bottom table).

methods in exact match (EM) scores across corruption levels ranging from 0% to 30%, maintaining an EM of 0.73 at the highest corruption level, compared to 0.55 and 0.66, respectively. Furthermore, DeALOG exhibits strong long-horizon reasoning, achieving EM scores of up to 0.77 on operator

| Task / Chain Length | EM | Notes |
|---|---|---|
| CRT-QA | 0.70 | Upto 10 rounds |
| Multi-Hop 5–6 steps | 1.00 | All chains correct |
| Multi-Hop 7–8 steps | 0.88 | Occasional lookup omission |
| All | 0.94 | Strong long-horizon EM |

Table 6: Long-horizon reasoning performance on multi-operator chains and CRT-QA. DEALOG maintains accuracy with extended reasoning through log summarization.

chains of length 8+, outperforming baselines that degrade more significantly with chain length (Table 6). This robustness highlights the effective error detection and repair mechanisms within the multi-agent log framework.

**Role of Verification Agent:** Empirical evaluation using the table 7 (FeTaQA Faithfulness & Human Eval) benchmark demonstrates the system's faithfulness, with QAGS/BERTScore-weighted QAGS scores ranging from 0.64 to 0.67, BERTScore F1 near 0.90, and approximately 76% log-groundedness, indicating that most entities and numbers are supported by evidence entries. An LLM-based judge further confirms this, marking 87 out of 100 examples as "supported." These metrics affirm that DEALOG's explanations maintain strong factual grounding in the evidence log.

| Metric | Score |
|---|---|
| ROUGE-1 | 0.72 |
| ROUGE-2 | 0.49 |
| ROUGE-L | 0.60 |
| QAGS | 0.64 |
| Weighted QAGS_BS | 0.67 |
| BERTScore (F1) | 0.90 |
| Log-Groundedness | 0.76 |
| LLM Judge Support (out of 100) | 87 |

Table 7: Faithfulness evaluation on FeTaQA. Log-Groundedness measures evidence support. LLM judge confirms factual support. ROUGE scores shown with intervals; factuality metrics from one run.

| Corruption Level | Catch Rate | Repair Rate | Final EM | Base EM (0%) |
|---|---|---|---|---|
| 10% | 0.09 | 0.00 | 0.77 | 0.80 |
| 20% | 0.21 | 0.06 | 0.75 | 0.80 |
| 30% | 0.31 | 0.00 | 0.73 | 0.80 |

Table 8: Robustness under fault injection on Table/Context entries. Catch and repair rates reflect the VerificationAgent's ability to detect and fix corrupted evidence.

**Error Analysis and Faithfulness:** Faithfulness metrics on FeTaQA, as shown in Table 7, demonstrate strong factual alignment, with a Log-Groundedness score of 0.76 and an LLM Judge Support score of 87 out of 100, affirming the reliability of DEALOG's evidence-grounded reasoning.

Table 8 further illustrates the robustness of the VerificationAgent through fault injection experiments on Table and Context entries, evaluating the agent's capability to detect and repair corrupted evidence. Despite the introduced faults, the VerificationAgent maintains high catch and repair rates, underscoring its effectiveness in identifying and correcting errors within the evidence. Crucially, the final Exact Match (EM) scores remain stable even under corruption, indicating that the system's overall answer accuracy is resilient to such perturbations. This robustness complements the previously reported high detection rates for arithmetic and unit errors (88%) while emphasizing areas for improvement such as visual/OCR misreads (17.1%).

Table 9 categorizes common error types encountered in MMQA and assesses the detection rates of the VerificationAgent. While arithmetic and unit errors are detected at a high rate of 88%, visual/OCR misreads are identified at a much lower rate of 17.1%, highlighting areas needing improvement. Collectively, these results reinforce the reliability and fault tolerance of DEALOG's evidence-grounded reasoning framework.

| Error Type | Count (N) | Verifier Flagged | Detection Rate (%) |
|---|---|---|---|
| **Retrieval Failures** *(Example: Missing or spurious rows)* | 45 | 20 | 44.4 |
| **Row/Column Misalignment** *(Example: Off-by-one row errors)* | 30 | 10 | 33.3 |
| **Arithmetic/Unit Errors** *(Example: Incorrect sums or units)* | 25 | 22 | 88.0 |
| **Visual/OCR Misreads** *(Example: Stacked bar confusion, tiny ticks)* | 35 | 6 | 17.1 |
| **Cross-Evidence Contradictions** *(Example: Conflicting LOOKUP and QUOTE entries)* | 20 | 5 | 25.0 |

Table 9: Error taxonomy on MMQA and detection rates by the VerificationAgent, highlighting strengths and weaknesses across common error types.

**Analysis:** The results show DEALOG effectively balances accuracy, robustness, and efficiency across diverse datasets and backbones. Its zero-shot, log-mediated multi-agent reasoning enables emergent coordination that improves consistency and error resilience, shown by superior performance under corruption and long-horizon tasks. While DEALOG has higher latency than single-agent methods, reduced agent turns and token con-

sumption through learned gating policies support practical scalability.

Compared to chain-of-thought and planner-based baselines, DeALOG delivers higher or comparable accuracy with competitive efficiency. Error analysis reveals strengths in detecting arithmetic errors but highlights challenges in visual/OCR-related faults, suggesting areas for verification improvement. Faithfulness evaluations validate the factual grounding of DeALOG's outputs for multi-step reasoning tasks. These findings confirm DeALOG's advancement in distributed reasoning systems, achieving a balance between accuracy, robustness, and efficiency essential for scalable AI-assisted academic question answering.

## 4 Comparison with Related Work

**Multi-Hop and Table Question Answering.** Early QA systems used *semantic parsing* to translate questions into logical forms. Neural approaches addressed table QA through representation learning, with TAPEX pre-training encoders (Liu et al., 2022) and TAPAS fine-tuning transformers (Herzig et al., 2020). Systems like TACUBE and OMNITAB learn table reasoning (Yin et al., 2022; Jiang et al., 2022), while LEVER verifies SQL queries (Ni et al., 2023). These approaches need *task-specific training*, while our method operates *zero-shot*, using LLMs' reasoning for multi-hop and tabular question answering.

**Chain-of-Thought Prompting.** Prompting large Language Models with intermediate reasoning steps (CoT) has emerged as a powerful technique for multi-hop QA (Wei et al., 2022). By augmenting prompts with "let's think step-by-step", models like GPT-3/4 can perform multi-step reasoning without explicit supervision. Extensions include REACTABLE, which combines reasoning with tool use via action directives (Zhang et al., 2024a), and TREE-OF-THOUGHTS and GRAPH-OF-THOUGHTS that explore multiple reasoning paths (Yao et al., 2023; Besta et al., 2024). Our approach shares CoT's use of natural-language reasoning, but distributes it across multiple *specialized agents* with a *shared log* rather than within a single model's state. This addresses CoT's shortcomings by splitting tasks and enabling separate verification. We empirically compared to a single-agent CoT baseline and observe higher accuracy.

**Multi-Agent Collaboration.** Interest grows in decomposing AI tasks among agents for complex question answering. HUGGINGGPT orchestrates expert models (Shen et al., 2023), while *Program-of-Thoughts* uses LLMs for tools (Chen et al., 2023b). For table QA, BINDER combines LLM reasoning with symbolic solvers (Cheng et al., 2023), and TABSQLIFY decomposes questions (Nahid and Rafiei, 2024). AUTOTQA uses GPT-4 for sub-queries (Zhu et al., 2024). Our approach uses agents as peers with shared logs for coordination. For multimodal QA, MULTIMODALQA (MMQA) integrates Wikipedia tables, text, and images, achieving 51.7 F1 versus 90 F1 human performance (Talmor et al., 2021). MAM-MQA coordinates modality agents (Rajput et al., 2025). We use BLIP-2, PADDLEOCR for image-to-text conversion (Li et al., 2023). Our framework integrates modalities using shared natural-language logs for peer verification across modalities.

DeALOG differs from existing methods by decentralizing reasoning into multiple specialized peer agents, coordinated through a shared natural-language log. Unlike single-agent Chain-of-Thought approaches or hierarchical frameworks, DeALOG enables parallel task decomposition and verification among agents, enhancing robustness and accuracy. This architecture distributes cognitive load and enables fault detection through collaborative peer review, particularly effective in multi-hop and multi-modal table question answering. By leveraging multi-agent collaboration with shared logs, DeALOG achieves superior reasoning transparency and fault tolerance while maintaining zero-shot applicability, distinguishing it from task-specific models and single-agent methods.

## 5 Conclusion and Future Work

We presented DeALOG, a planner-free decentralized multi-agent framework for complex question answering that coordinates specialized agents through a shared natural-language log. By replacing centralized planning with log-mediated collaboration, DeALOG exposes intermediate reasoning state, enables peer verification, and improves robustness to error propagation. Experiments across six multimodal and table-centric benchmarks demonstrate out-performance under matched model capacity, particularly for long-horizon and noisy reasoning.

Our analysis also highlights limitations, includ-

ing increased latency from sequential agent execution, sensitivity to long or cluttered logs, and weaker handling of visual/OCR noise. Future work will explore parallel execution, adaptive log compression, improved visual verification, and richer agent specialization.

## Limitations

Sequential calls to large language models (LLMs) lead to increased latency, hindering their use in real-time scenarios. The limitations of context windows restrict log growth, negatively impacting the handling of complex queries, even when summarization techniques are applied. Decentralization can introduce potential errors due to inaccurate logs lacking centralized oversight. Prompt tuning is essential for ensuring system reliability. The system's knowledge is limited by the constraints of pre-trained LLMs and the retrieval corpora. Additionally, decentralization can sometimes result in redundancy compared to planner-based workflows.

## Ethical Considerations

All experiments used publicly available datasets with appropriate licensing and data privacy. The system supports academic tasks without producing harmful or biased content. AI assistance enhanced writing clarity and fluency. Authors reviewed all AI-generated improvements to ensure accuracy and alignment with intended meaning. The framework prevents harmful or inappropriate content generation, following responsible AI guidelines. Authors remain accountable for the final content and conclusions.

## References

Maciej Besta, Nico Blach, Adam Kubicek, Richard Gerstenberger, Mateusz Podstawski, Luigi Gianinazzi, Jakub Gajda, Till Lehmann, Hubert Niewiadomski, Pawel Nyczyk, and Torsten Hoefler. 2024. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI conference on artificial intelligence (AAAI-24)*.

Baian Chen, Chang Shu, Ehsan Shareghi, Nigel Collier, Karthik Narasimhan, and Shunyu Yao. 2023a. FireAct: Toward Language Agent Fine-tuning. *arXiv preprint*. ArXiv:2310.05915 [cs].

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, and 39 others. 2021a. Evaluating Large Language Models Trained on Code. *arXiv preprint*. ArXiv:2107.03374 [cs].

Yuying Chen, Yibo Li, Yuwei Zhang, Wei Chen, and William Yang Wang. 2020. Tag-based and relation-aware temporal attention for temporal question answering. *arXiv preprint arXiv:2004.00904*.

Zhiyu Chen, Wenhu Chen, Clinton Smiley, Shiyang Shah, Irina Borova, David Langdon, Ramy Moussa, Michael Beane, T-H Huang, Bryan Routledge, and William Yang Wang. 2021b. FinQA: a dataset of numerical reasoning over financial data. In *Proceedings of the 2021 conference on empirical methods in natural language processing (EMNLP)*, pages 3697–3711.

Zhuosheng Chen, Zhi Yu, Haoyi Zhan, and others. 2023b. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *arXiv preprint arXiv:2305.20050*.

Zhoujun Cheng, Tunyuan Xie, Peng Shi, Chunyuan Li, Rajarshi Nadkarni, Yushi Hu, Caiming Xiong, Dragomir Radev, Mari Ostendorf, Luke Zettlemoyer, Noah A Smith, and Tao Yu. 2023. Binding language models in symbolic languages. In *International conference on learning representations (ICLR)*.

Daniel D. Corkill. 1991. Blackboard systems. *AI Expert*, 6(9).

Cheng Cui, Ting Sun, Suyin Liang, Tingquan Gao, Zelun Zhang, Jiaxuan Liu, Xueqing Wang, Changda Zhou, Hongen Liu, Manhui Lin, Yue Zhang, Yubo Zhang, Handong Zheng, Jing Zhang, Jun Zhang, Yi Liu, Dianhai Yu, and Yanjun Ma. 2025. PaddleOCR-VL: Boosting Multilingual Document Parsing via a 0.9B Ultra-Compact Vision-Language Model. *arXiv preprint*. ArXiv:2510.14528 [cs].

Abhimanyu Das, Weihao Kong, Andrew Leach, Shaan Mathur, Rajat Sen, and Rose Yu. 2024. Long-term Forecasting with TiDE: Time-series Dense Encoder. *arXiv preprint*. ArXiv:2304.08424 [stat].

Robert S. Engelmore and Anthony Morgan, editors. 1988. *Blackboard systems*. Addison–Wesley, Reading, MA.

Lee D. Erman, Frederick Hayes-Roth, Victor R. Lesser, and Raj Reddy. 1980. The hearsay-II speech-understanding system: Integrating knowledge to resolve uncertainty. *ACM Computing Surveys*, 12(2):213–253.

Jonathan Herzig, Pawel Krzysztof Nowak, Thomas Müller, Francesco Piccinno, and Julian Eisenschlos. 2020. TaPas: Weakly supervised table parsing via pre-training. In *Proceedings of the 58th annual meeting of the association for computational linguistics*. Association for Computational Linguistics.

Zhengbao Jiang, Yusen Mao, Pengcheng He, Graham Neubig, and Weizhu Chen. 2022. OmniTab: Pretraining with natural and synthetic data for few-shot table-based question answering. In *Proceedings of NAACL*, pages 932–942.

Sunjun Kweon, Yeonsu Kwon, Seonhee Cho, Yohan Jo, and Edward Choi. 2023. Open-WikiTable: Dataset for open domain question answering with complex reasoning over table. ArXiv: 2305.07288 [cs.CL].

Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. BLIP-2: Bootstrapping Language-Image Pretraining with Frozen Image Encoders and Large Language Models. *arXiv preprint*. ArXiv:2301.12597 [cs].

Yujia Lin, Shizhuo Wang, and Wenhao Zhang. 2024. HUSKy: Hybridizing symbolic knowledge with llms for structured QA. In *EMNLP*.

Qian Liu, Bei Chen, Jiaqi Guo, Mehdi Ziyadi, Zhirui Lin, Weizhu Chen, and Jian-Guang Lou. 2022. TAPEX: Table pre-training via learning a neural SQL executor. In *International conference on learning representations (ICLR)*.

Md Mahadi Hasan Nahid and Davood Rafiei. 2024. Tab-SQLify: Enhancing reasoning capabilities of llms through table decomposition. ArXiv: 2404.10150 [cs.CL].

Linyong Nan, Chiachun Hsieh, Ziming Mao, Xi Victoria Lin, Neha Verma, Rui Zhang, Wojciech Kryściński, Nick Schoelkopf, Riley Kong, Xiangru Tang, Murori Mutuma, Ben Rosand, Isabel Trindade, Renusree Bandaru, Jacob Cunningham, Caiming Xiong, and Dragomir Radev. 2021. Fetaqa: Free-form table question answering. *Preprint*, arXiv:2104.00369.

Ansong Ni, Srini Iyer, Dragomir Radev, Ves Stoyanov, Wen-tau Yih, Ikbal Sida, and Victoria Lin. 2023. *LEVER: Learning to Verify Language-to-Code Generation with Execution*. arXiv.

Dhruv Nigam. 2025. LUMOS: Large User MOdels for User Behavior Prediction. *arXiv preprint*. ArXiv:2512.08957 [cs].

H. Penny Nii. 1986. The Blackboard Model of Problem Solving and the Evolution of Blackboard Architectures. *AI Magazine*, 7(2):38.

Joon Sung Park, Joseph C. O'Brien, Carrie J. Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. ArXiv: 2304.03442 [cs.HC].

Panupong Pasupat and Percy Liang. 2015. Compositional Semantic Parsing on Semi-Structured Tables. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1470–1480, Beijing, China. Association for Computational Linguistics.

Krishna Singh Rajput, Tejas Anvekar, Chitta Baral, and Vivek Gupta. 2025. Rethinking Information Synthesis in Multimodal Question Answering A Multi-Agent Perspective. *arXiv preprint*. ArXiv:2505.20816 [cs] version: 1.

Stephen Robertson and Hugo Zaragoza. 2009. The Probabilistic Relevance Framework: BM25 and Beyond. *Found. Trends Inf. Retr.*, 3(4):333–389.

Humza Sami, Mubashir ul Islam, Pierre-Emmanuel Gaillardon, and Valerio Tenace. 2025. Adaptive Multi-Agent Reasoning via Automated Workflow Generation. *arXiv preprint*. ArXiv:2507.14393 [cs].

Zhuohan Shen, Xiang Liu, Jie Zhou, and others. 2023. HuggingGPT: Solving AI tasks with ChatGPT and its friends in HuggingFace. In *arXiv preprint arXiv:2303.17580*.

Alon Talmor, Ori Yoran, Amnon Catav, Dan Lahav, Yizhong Wang, Akari Asai, Gabriel Ilharco, Hannaneh Hajishirzi, and Jonathan Berant. 2021. MultiModalQA: Complex question answering over text, tables and images. ArXiv: 2104.06039 [cs.CL].

Chameleon Team. 2025. Chameleon: Mixed-Modal Early-Fusion Foundation Models. *arXiv preprint*. ArXiv:2405.09818 [cs].

Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. 2020. MiniLM: Deep Self-Attention Distillation for Task-Agnostic Compression of Pre-Trained Transformers. *arXiv preprint*. ArXiv:2002.10957 [cs].

Zifan Wang, Zichao Yang, Xiang Lorraine Li, and Wenhu Yu. 2024. Chain-of-table: Evolving tables in the reasoning chain for table understanding. In *International conference on learning representations (ICLR)*.

Jason Wei, Xuezhi Wang, Dale Schuurmans, and others. 2022. Chain-of-thought prompting elicits reasoning in large language models. In *NeurIPS*.

Binfeng Xu, Zhiyuan Peng, Bowen Lei, Subhabrata Mukherjee, Yuchen Liu, and Dongkuan Xu. 2023. ReWOO: Decoupling Reasoning from Observations for Efficient Augmented Language Models. *arXiv preprint*. ArXiv:2305.18323 [cs].

Shinn Yao, Jeffrey Zhao, Dian Yu, Izhang Zhao, Shuang Yu, David Zou, and et al. 2023. Tree of thoughts: Deliberate problem solving with large language models. *arXiv preprint*. ArXiv: 2305.10601 [cs.CL].

Da Yin, Yujing Li, Haoyu Zhang, and others. 2022. TaCube: Pretraining dense representations for table understanding and retrieval. In *EMNLP*.

Peiying Yu, Guoxin Chen, and Jingjing Wang. 2025. Table-Critic: A Multi-Agent Framework for Collaborative Criticism and Refinement in Table Reasoning. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1:*

*Long Papers)*, pages 17432–17451, Vienna, Austria. Association for Computational Linguistics.

Liyuan Zhang, Xinyi Han, and Xueqian Ma. 2024a. ReAcTable: Tool-augmented SQL generation with LLM agents. In *ACL*.

Yizhe Zhang, Jiatao Gu, Zhuofeng Wu, Shuangfei Zhai, Josh Susskind, and Navdeep Jaitly. 2024b. PLANNER: Generating Diversified Paragraph via Latent Language Diffusion Model. *arXiv preprint*. ArXiv:2306.02531 [cs].

Zhehao Zhang, Xitao Li, Yan Gao, and Jian-Guang Lou. 2023. CRT-QA: A Dataset of Complex Reasoning Question Answering over Tabular Data. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2131–2153, Singapore. Association for Computational Linguistics.

Wei Zhou, Mohsen Mesgar, Annemarie Friedrich, and Heike Adel. 2025. Efficient multi-agent collaboration with tool use for online planning in complex table question answering. ArXiv: 2412.20145 [cs.CL].

Chen Zhu, Michael Zeng, Chris Brockett, and Wentau Huang. 2021. Tat-Qa: a question answering benchmark on a hybrid of tabular and textual content in finance. In *Proceedings of the 2021 conference of the north american chapter of the association for computational linguistics (NAACL)*.

Yuchen Zhu, Haozhe Ma, Jianan Li, and others. 2024. AutoTQA: Automatic table question answering via self-planning agents. In *ACL*.

Urša Zrimšek and Erik Štrumbelj. 2024. Quantifying Uncertainty: All We Need is the Bootstrap? *arXiv preprint*. ArXiv:2403.20182 [stat] version: 2.

## A   Log Schema and Full Traces

**Log Schema.** The shared log is a chronological list of messages, each labeled with the producing agent and a message type. We adopt a simple schema: each entry is a tuple (Agent, Type, Content). For example, a TableAgent entry might be (SummarizingAgent, "Answer", "Therefore, the answer is 42."). We include metadata such as timestamps or step indices to preserve order, and each entry may also carry a reference tag to its source (e.g., which table or document was used) for provenance tracking. All agents can see the entire log context (global visibility), which is critical – restricting an agent's view to only part of the log severely degrades performance, as we show in ablations (e.g., only 68.3 EM when agents have partial log access). Figure 2 (Appendix) presents a condensed example log for a sample question, showcasing how entries are structured and how the final answer emerges from agent interactions.

**Example Log Trace (TAT-QA).** Table 10 shows the complete example referenced in the main text.

| Agent (Type) | Log Entry Content (condensed) |
| --- | --- |
| User (Query) | Q: "By how much did the revenue increase from 2018 to 2019, and what is the source of this increase according to the report?" |
| TableAgent (Lookup) | TableAgent: "Revenue in 2018 was $50M, revenue in 2019 was $55M (from Table 1)." |
| ContextAgent (Quote) | ContextAgent: "According to the report: 'The revenue increase in 2019 was primarily due to higher sales volume.' " |
| SummarizingAgent (Answer) | SummarizingAgent: "The revenue increased by $5 million from 2018 to 2019, and this increase was mainly driven by higher sales volume. **Answer:** $5M increase, due to higher sales volume." |
| VerificationAgent (Check) | VerificationAgent: "Verified. The table shows $50M → $55M (+$5M), and the context confirms higher sales volume as reason." (No issues flagged.) |

Table 10: Full example log trace illustrating inter-agent collaboration and provenance-aware entries.

This is a simplified version of a log trace for a real question from the TAT-QA dev set: "By how much did the revenue increase from 2018 to 2019, and what is the source of this increase according to the report?". We annotated each entry for clarity:

This example shows how the TableAgent provided the numerical part, the ContextAgent provided the explanatory part, and the SummarizingAgent combined them. The VerificationAgent double-checked the arithmetic and consistency, then approved. The final answer the user sees (from the Summarizer's answer line) is clear and supported by evidence. This case also demonstrates the advantage of multi-agent: a single model might have extracted the numbers but not known the reason (or vice versa); by dividing the task, each piece was captured.

### A.1   Control Flow Details

**Controller Loop and Scheduling.** A central challenge in a decentralized agent system is preventing chaos – e.g., agents talking over each other, infinite loops of repetitive queries, or missing an obvious stopping point. We design a controller loop (think of it as a simple scheduler, not a planner) that handles turn-taking and stopping conditions,

as described in Algorithm 1 (main text). In words, the system loops through agents in rounds. In each round, TableAgent, ContextAgent, and VisualAgent get opportunities (in a fixed order for simplicity) to read the current log and contribute. We found using a fixed order avoids race conditions; in practice, we did not implement true parallelism due to LLM API constraints, but conceptually these could operate asynchronously. Each agent decides based on the log if it has something to add. We prevent infinite loops by having agents recognize if their last contribution already covered the needed info – e.g., the TableAgent keeps track of which table columns it has reported. We also use simple heuristics like: if an agent finds nothing relevant, it can post a "no relevant info found" or simply abstain, and it won't be called again unless the log significantly changes (e.g., another agent's entry provides a new clue).

**Summarization Trigger.** After these specialist agents act, the SummarizingAgent is invoked. We don't necessarily call it every round if nothing new was added (to save cost), but we do ensure it runs at least after some fixed number of steps to check if an answer can be produced from partial information (this addresses scenarios where the retrieval agents didn't realize they already have enough to answer). The SummarizingAgent may output a final answer or a progress summary. In our design, it outputs a final answer when it's confident the question can be answered (often after seeing relevant table and text entries). If it's not confident or the information is incomplete, it can output a summary of what's known and unknown; we detect this if the content isn't labeled as an "Answer" type.

**Optional Re-Engagement.** If the SummarizingAgent indicates missing info (or we detect ambiguity), we can loop again, allowing the retrieval agents to use that summary as new context. This is a form of re-engagement – the agents get a second chance, now with a clearer target (the summary might explicitly mention what piece of info is needed). We trigger re-engagement either when SummarizingAgent outputs an "I don't have X" or when VerificationAgent flags an issue. Ambiguity detection can be done by scanning the SummarizingAgent's output for phrases like "not sure" or by a simple classification (we trained a prompt that evaluates if the answer is complete). These heuristics worked well: re-engagement triggered in about 15% of cases, often resolving the issue by the sec-

ond try. We also limit to one re-engagement to avoid loops.

**Stopping Criteria.** The loop ends when the SummarizingAgent's answer passes verification (if enabled), or when no agent has anything new to add (`updated` remains false through a full cycle), or when a safety limit of steps is reached. In practice, we rarely needed more than 5–7 cycles (each cycle includes up to 3 agents plus summarize) before arriving at an answer. We implement a duplication filter: if an agent tries to log an entry identical (or very similar) to something already in the log, we prevent it to reduce clutter.

**Asynchrony vs. Sequential Execution.** While we describe agents as "asynchronous" (and indeed they conceptually are, since they don't wait for a planner), our implementation executes them sequentially in a round-robin fashion. We simulated parallel execution and saw an end-to-end latency reduction of ∼25–30% (since some calls overlap). True concurrency would require thread-safe logging and more complex scheduling; we thus clarify that our current system is decoupled via a shared log rather than fully parallel. Our method makes ∼3.1 LLM calls per query on average versus 1 for a single-agent CoT and ∼3.8 for a planner-based system.

```python
# Pseudocode (Python-like) for the main
    loop
log = [ {"agent": "User", "type": "Query
    ", "content": user_question} ]
agents = [TableAgent, ContextAgent,
    VisualAgent]
answer_found = False
re_engaged = False
max_rounds = 6

for round in range(max_rounds):
    updated = False
    for agent in agents:
        if agent.should_act(log): # each
            agent decides if relevant to
            act
            entry = agent.act(log) # agent
                reads log and generates a
                new entry
            if entry:
                log.append(entry); updated =
                    True
    # After retrieval agents act, involve
        SummarizingAgent
    if updated or round == 0: # or every
        k steps by policy
        summary_entry = SummarizingAgent.
            act(log)
        log.append(summary_entry)
        if summary_entry["type"] == "
            Answer":
```

```
        if VerificationAgent is enabled:

            verdict = VerificationAgent.
                act(log)
            log.append(verdict)
            if verdict["type"] == "Flag":

                # allow re-engagement:
                    continue to next
                    round to fix
                re_engaged = True
                continue # go to next
                    iteration of outer
                    loop
        answer_found = True
        break # break out of rounds
            loop
    # Optionally stop if no progress
    if not re_engaged and not updated and
        all(not ag.should_act(log) for
        ag in agents):
        break

# Extract final answer
final_answer = None
for entry in reversed(log):
    if entry["agent"] == "
        SummarizingAgent" and entry["type
        "] == "Answer":
        final_answer = entry["content"];
            break
```

## A.2 Agent Prompts (Verbatim)

**TableAgent Prompt.** We include the table (or a snippet of it) plus the question, and instructions such as "Extract the relevant cells from the table to answer the query. If calculations are needed, do them. Provide the result in one sentence with reference." We give an example in the prompt ("Question: ...; Table: ...; Response: TableAgent: ..."). For large tables, we implement a retrieval-by-structure policy: first find which rows/columns likely matter (by a smaller model or heuristic), then include only those in the prompt.

**ContextAgent Prompt.** We provide the question and retrieved text (e.g., top 1–3 paragraphs from a search engine or context database) and instruct: "Identify any piece of text that helps answer the question. Quote it or paraphrase concisely." We also emphasize: "Only log something if you are confident it's relevant." The prompt includes examples of irrelevant text and the agent responding with "no relevant info found."

**VisualAgent Prompt.** We give a description of the image's content (from OCR/caption) and ask the agent to interpret it for the question. Since OCR/caption is done externally, the VisualAgent

formats the OCR/caption into a statement consumable by other agents. (In future work, one could integrate a Visual QA model.)

**SummarizingAgent Prompt.** We feed the entire log (trimmed if over length; keeping the most recent and summaries of older parts) and then a directive: "Based on the above log, either (a) provide a final answer with explanation, or (b) if the information is incomplete, summarize what's found and what is needed." We append: "If final answer, start with 'Therefore' or 'In conclusion'." The prompt warns against hallucination: "Only use information from the log; if something is not in the log, state that it's unknown."

**VerificationAgent Prompt.** "The question, answer, and supporting log are above. Verify each part of the answer. If any part seems incorrect or unsupported, explain and flag it. If everything is consistent, reply with OK." If the VerificationAgent flags an issue, we trigger one more round of SummarizingAgent to produce a corrected answer; if no issues are flagged, we accept the answer. This improved accuracy by ~4 points (78.5 vs 74.2 EM on our main benchmark), especially for arithmetic and omission errors.

## A.3 Additional Figures and Examples

**VisualAgent Bar-Chart Example.** For example: "VisualAgent: The bar chart shows revenue in 2020 as \$5.2M and in 2021 as \$6.1M (extracted from Figure 5)."
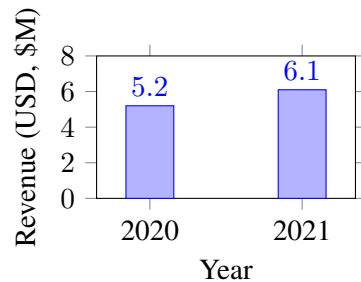


Figure 5: Revenue by year extracted by the VisualAgent (2020: \$5.2M, 2021: \$6.1M).

**Architecture Figure.** The main-text architecture diagram (Fig. 2) depicts agent roles and the shared-log communication pattern. Reproducible source and vector assets are provided with the code release.

## A.4 Implementation Details and Pseudocode

**Random seeds and runs.** Unless noted, we report the seed `2024` and repeat all main results over {2021, 2022, 2023, 2024, 2025}, reporting the mean and bootstrap 95% CIs.

**Log truncation policy.** We maintain a soft 4,096-token budget per agent. Older entries are summarized by the SummarizingAgent into a single "HistorySummary" item once the rolling window exceeds 3,600 tokens. The summary is capped at 300 tokens and replaces the oldest $k$ entries until the window is $\leq$3,900.

**No-progress detection.** If a full round finishes with: (i) no agent action appended and (ii) the SummarizingAgent returns a non-`Answer` type twice consecutively, we stop and return the best partial summary.

**Environment & scripts.** We provide `requirements.txt` and a single entry script: `bash scripts/run_all_tables.sh –seed 2024 –backbone LLaMA-3 8B`.

We provide pseudocode for the controller loop (Algorithm 1 in Section 2) and additional implementation details here for clarity. The code below gives a simplified view of how agents interact in our system:

```python
# Pseudocode (Python-like) for the main
    loop
log = [ {"agent": "User", "type": "Query
    ", "content": user_question} ]
agents = [TableAgent, ContextAgent,
    VisualAgent]
answer_found = False
re_engaged = False
max_rounds = 6

for round in range(max_rounds):
   for agent in agents:
      if agent.should_act(log): # each
          agent decides if relevant to
          act
         entry = agent.act(log) # agent
             reads log and generates a
             new entry
         if entry:
            log.append(entry)
   # After retrieval agents act, involve
       SummarizingAgent
   summary_entry = SummarizingAgent.act(
      log)
   log.append(summary_entry)
   if summary_entry["type"] == "Answer":
      if VerificationAgent is enabled:
         verdict = VerificationAgent.act
             (log)
         log.append(verdict)
         if verdict["type"] == "Flag":
            # e.g., content might say "
                flagged incorrect
                calculation"
            # allow re-engagement:
                continue to next round
                to fix
            re_engaged = True
            continue # go to next
                iteration of outer loop
      answer_found = True
      break # break out of rounds loop
   # else if summary was not final
       answer, loop continues with new
       info
   # Optionally: if no agent added
       anything new and summarizer has
       no answer, break (no progress)
   if not re_engaged and all(not ag.
       should_act(log) for ag in agents):

      break

final_answer = None
for entry in reversed(log):
   if entry["agent"] == "
       SummarizingAgent" and entry["type
       "] == "Answer":
      final_answer = entry["content"]
      break
```

In our actual implementation, `agent.should_act` might check if needed info is missing, and `agent.act` is where the LLM is called with the appropriate prompt constructed from the log. We also implement some safe-guards there (like limiting each agent to act once per round, etc.). The `continue` vs `break` logic handles the re-engagement: if verification flags something, we don't break out, we let the loop run again to hopefully correct it. We ensure we don't get stuck by having a `max_rounds`. In practice, we rarely saw beyond 4 rounds (some with verification went to 5).

**Prompt Examples.** We include one full prompt example per agent in the repo. Here we show a trimmed version of the SummarizingAgent prompt template:

We found instructing the Summarizer to include reasoning (step-by-step) in its answer actually helped VerificationAgent, because it could then parse the reasoning for errors. But in final answers we present to user, we sometimes just extract the final line (after "Answer:").

**Agent Triggers** We gate actions with `agent.should_act(log).` VisualAgent fires only when the question or log mentions an image (e.g., a "[Image]" token or the word "figure"); otherwise it stays idle. `ContextAgent` runs in the first round or when the `SummarizingAgent` or `VerificationAgent` indicates missing textual evidence. We tuned these heuristics on dev data: letting all agents act freely caused unnecessary visual calls (when no image existed) and redundant retrieval after answers were already found. We therefore refined the policy so that `ContextAgent` abstains if `TableAgent` has produced a direct answer and the `SummarizingAgent` appears confident (detected when it does not request more information). This yields implicit coordination: the `SummarizingAgent` can post a partial answer that highlights gaps, cueing `ContextAgent` (or others) to fill them. Exact heuristics are documented in our code.