

SHEETMIND: AN END-TO-END LLM-POWERED MULTI-AGENT FRAMEWORK FOR SPREADSHEET AUTOMATION*

Ruiyan Zhu^{1†}, Xi Cheng^{1‡}, Ke Liu², Brian Zhu¹, Daniel Jin¹, Neeraj Parihar¹, Zhoutian Xu¹, Oliver Gao¹

¹Cornell University, Ithaca, NY, USA

²University of California, Berkeley, CA, USA

ABSTRACT

We present *SheetMind*, a modular multi-agent framework powered by large language models (LLMs) for spreadsheet automation via natural language instructions. The system comprises three specialized agents: a *Manager Agent* that decomposes complex user instructions into subtasks; an *Action Agent* that translates these into structured commands using a Backus–Naur Form (BNF) grammar; and a *Reflection Agent* that validates alignment between generated actions and the user’s original intent. Integrated into Google Sheets via a Workspace extension, SheetMind supports real-time interaction without requiring scripting or formula knowledge. Experiments on benchmark datasets demonstrate an 80% success rate on single-step tasks and approximately 70% on multi-step instructions, outperforming ablated and baseline variants. Our results highlight the effectiveness of multi-agent decomposition and grammar-based execution for bridging natural language and spreadsheet functionalities. The project code is available at <https://github.com/colonel-aureliano/Excel-Agent>.

Keywords Spreadsheet Automation · Multi-Agent Systems · Large Language Models · End-to-End Execution · Natural Language Interfaces

1 Introduction

Spreadsheets are among the most widely used software tools for data analysis and management. Yet, effective use of spreadsheet functionalities often requires users to master formulas, macros, and structured syntax, creating a barrier for non-technical users. Natural Language Interfaces (NLIs) [1] offer a path toward greater usability. Recent advances in large language models (LLMs), such as GPT-family models [2, 3], have revitalized interest in this goal by enabling powerful reasoning and code-generation capabilities for table-based tasks [4, 5].

However, designing robust NLIs for spreadsheets remains non-trivial. Spreadsheet manipulation often involves complex operations with latent constraints, ambiguous goals, and multiple interdependent steps. Prior work such as OmniTab [6], and SpreadsheetCoder [7] have focused on one-shot code generation or formula prediction, which limits their ability to handle tasks requiring multi-step reasoning. More recent systems like SheetCopilot [5] and SheetAgent [8] begin to integrate LLMs with spreadsheet environments, enabling end-to-end automation. Yet, these systems rely on monolithic prompting pipelines and lack the modularity or robustness to adapt to multi-phase workflows or incorporate feedback during execution.

SheetMind addresses these limitations by introducing a *multi-agent architecture* for LLM-powered spreadsheet reasoning and manipulation. Inspired by general-purpose agentic frameworks like MetaGPT [9] and PC-Agent [10], SheetMind decomposes complex natural language instructions into executable subtasks through specialized agents. Our system includes (1) a *Manager Agent* that interprets and decomposes user intent; (2) an *Action Agent* that generates structured spreadsheet actions using a Backus-Naur Form (BNF) grammar [11]; and (3) a *Reflection Agent* that validates and revises actions to ensure alignment with the original intent.

*This is a preprint under review.

[†]Both authors contributed equally to this research.

[‡]Corresponding author: xc557@cornell.edu

We deploy SheetMind as a Google Workspace extension, enabling real-world spreadsheet manipulation through conversational interfaces. Our framework demonstrates improved performance in both accuracy and transparency compared to single-agent baselines. SheetMind advances the use of LLMs as effective collaborators in structured environments through a modular and extensible framework for intelligent spreadsheet assistance.

2 Related Work

2.1 Large Language Models for Spreadsheets

Recent advancements in spreadsheet automation leverage natural language instructions via Large Language Models (LLMs). SheetCopilot introduced atomic spreadsheet actions and a state-machine planner to handle multi-step tasks, significantly surpassing previous code-generation methods [5]. SheetAgent further refined spreadsheet reasoning with a modular agent architecture comprising Planner, Informer, and Retriever components, effectively addressing long-horizon spreadsheet tasks through iterative reasoning [8]. TableTalk enhanced spreadsheet development by scaffolding user interactions through an observation-planning-action loop, significantly reducing cognitive load and improving spreadsheet construction efficiency [12]. These systems typically convert user requests into a sequence of formula edits or scripting calls to accomplish the task [13].

2.2 General-Purpose Multi-Agent LLM Systems

Beyond spreadsheets, multi-agent LLM frameworks such as MetaGPT and PC-Agent demonstrated that explicit role assignments and hierarchical decomposition significantly enhance task automation. MetaGPT encoded structured operating procedures for software development roles [9], while PC-Agent employed hierarchical agents specialized for GUI tasks [10]. Similarly, AFLOW automated agentic workflow generation by exploring optimized LLM execution graphs through systematic refinement and feedback loops [14]. Likewise, frameworks such as AgentVerse [15] demonstrate that a team of collaborating LLM agents with complementary expertise can outperform a lone agent, exhibiting emergent cooperative behavior in reasoning, coding, and tool-use tasks. These approaches highlight the benefits of coordination strategies – from role specialization to search-based planning – in improving the reliability and scope of autonomous LLM-driven systems.

In contrast, our proposed system, SheetMind, uniquely integrates structured action generation using BNF grammar with a hierarchical multi-agent architecture (Manager, Action, and Reflection agents). SheetMind builds on PC-Agent [10] by coupling modular agents with structured action grammars in the spreadsheet domain.

3 Methodology

3.1 Problem Statement

Despite their widespread use, spreadsheet applications pose significant usability challenges for non-technical users. Complex formula syntax, conditional logic, and data manipulation functions often create a gap between user intent and execution.

SheetMind addresses this gap by enabling natural language interaction with spreadsheets. Users can issue instructions such as “Delete any element from the fifth column that starts with a number” without needing to write formulas or scripts.

Focusing on Google Sheets, SheetMind operates on typical spreadsheet data types: text, numbers, dates, and formulas, across various tasks, including data cleaning, analysis, modification, validation, structural changes, and cross-sheet operations. The system bridges intent and execution through structured multi-agent coordination and grammar-based action translation.

3.2 Multi-Agent Design

SheetMind employs a hierarchical multi-agent architecture inspired by general-purpose collaboration frameworks such as PC-Agent [10], while adapting the principles to structured spreadsheet environments. The system decomposes user instructions into granular actions via a sequence of specialized agents (i.e., **Manager Agent**, **Action Agent**, and **Reflection Agent**), each handling a distinct yet interdependent stage in the end-to-end workflow (Figure 1).

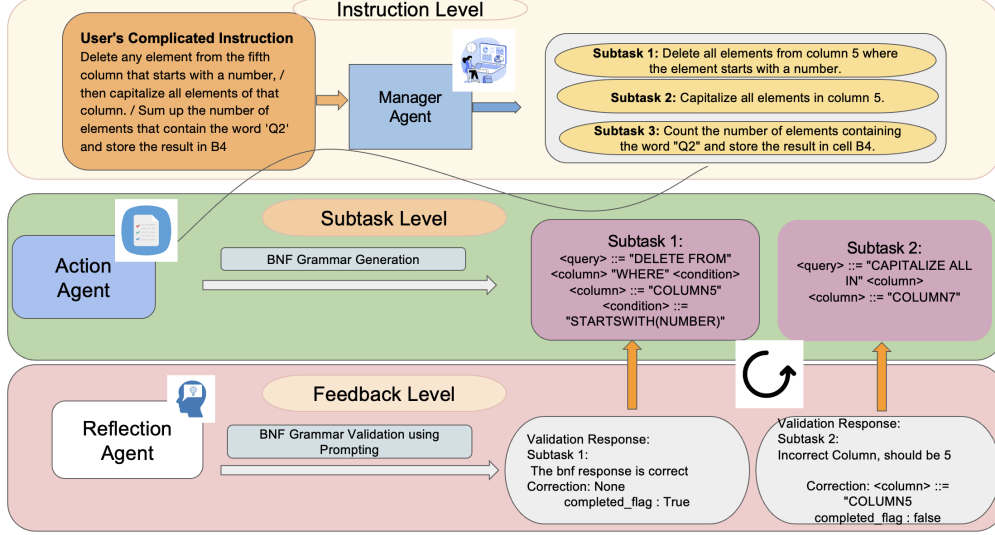


Figure 1: Multi-agent workflow for spreadsheet automation

Manager Agent The Manager Agent interprets user-issued natural language instructions $I = \{i_1, i_2, \dots, i_n\}$ and decomposes them into an ordered sequence of semantically coherent subtasks $T = \{t_1, t_2, \dots, t_k\}$. These subtasks reflect atomic spreadsheet operations, such as deleting rows, formatting cells, or computing column sums. The Manager also resolves dependencies among subtasks, enabling context sharing and parameter propagation. This parallels the task planning function in PC-Agent but is streamlined for deterministic, table-centric workflows.

Action Agent For each subtask t_i , the Action Agent generates a spreadsheet command a_i using a BNF-constrained grammar:

$$t_i \xrightarrow{\text{LLM} + \text{BNF Grammar}} a_i = (op, args, cond)$$

where op is a spreadsheet operation (e.g., SELECT, SET, DELETE), $args$ represent ranges or values, and $cond$ encodes optional filters (e.g., regular expressions or value predicates). This grammar-bound approach reduces hallucination and enhances alignment between the generated command and spreadsheet semantics. Unlike the GUI-level decision layer in PC-Agent, our Action Agent emits API-ready, deterministic commands optimized for structured environments.

Reflection Agent The Reflection Agent functions as the system’s critical validation and feedback module, overseeing both the Action Agent and, when necessary, the Manager Agent. It implements a two-stage evaluation process to ensure that generated actions remain semantically faithful to user intent:

- **Pre-execution semantic validation:** Before execution, each command a_i is evaluated for alignment with its subtask t_i using a validation function $v_i = \mathbb{V}(t_i, a_i) \in \{\text{valid}, \text{invalid}\}$. If invalid, the Reflection Agent triggers iterative regeneration via the Action Agent until a semantically consistent action is found.
- **Post-execution effect monitoring:** After action execution, the Reflection Agent compares spreadsheet states (O_{i-1}, O_i) to assess whether the outcome reflects the intended effect of t_i . If no change or an incorrect effect is detected, the Reflection Agent feeds back to the Action Agent to refine the command. For persistent failures, it escalates to the Manager Agent, requesting subtask reformulation.

This two-stage validation loop adapts PC-Agent’s reflection framework [10] to a symbolic domain, leveraging structured API feedback over visual perception.

Inter-Agent Feedback Loop. The interplay between agents forms a robust feedback-driven pipeline:

$$I \xrightarrow{\text{Manager}} T \xrightarrow{\text{Action}} A \xrightarrow{\text{Reflection}} \hat{A} \xrightarrow{\text{Execution}} \text{Spreadsheet}$$

The Reflection Agent acts as a central auditor, refining low-level actions and, when needed, guiding high-level task revisions, which enhances robustness and preserves user intent.

3.3 System Implementation

SheetMind’s architecture (Figure 2) comprises three coordinated agents in the **Back-End**: Manager, Action, and Reflection, leveraging LLMs and BNF grammar for structured task execution.

Front-End: The front-end offers a chat-based interface for natural language interaction within Google Sheets. It captures user instructions and dynamically extracts contextual information from the active spreadsheet, including cell values, types, and structural metadata. This contextualized request is packaged and sent to the back-end agent pipeline. Once the back-end returns a validated sequence of actions, the front-end executes them in real time and presents a clear, human-readable summary of the changes, ensuring transparency and user feedback.

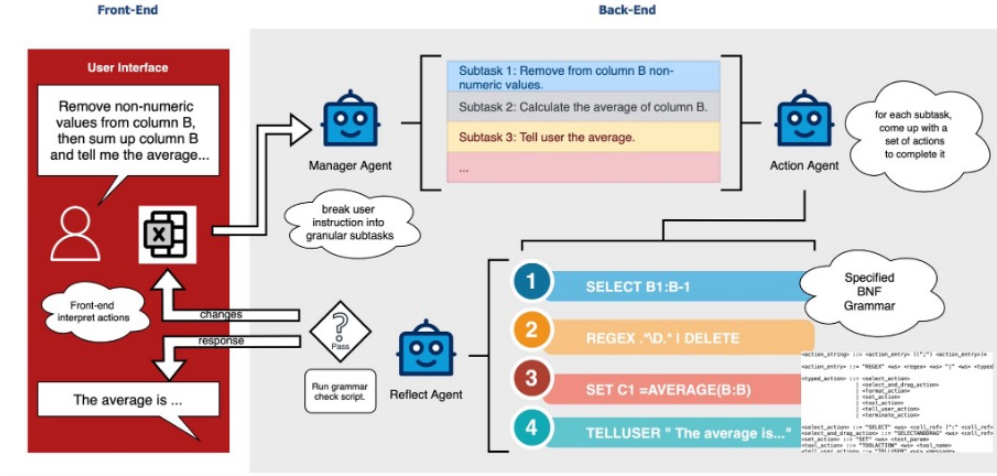


Figure 2: SheetMind system architecture

To summarize, SheetMind’s implementation incorporates the following key technologies:

- **Large Language Models** (via Google Gemini APIs) for natural language understanding and instruction parsing.
- **BNF-style grammars** to formalize and constrain spreadsheet operations into executable commands.
- **Google Apps Script integration** for seamless interaction with Google Sheets, enabling real-time content manipulation and feedback.

The workflow follows a clear sequence: user instructions are received, parsed by the Manager Agent into subtasks, transformed into actionable commands by the Action Agent, validated by the Reflection Agent, and finally executed on the spreadsheet. This modular approach allows for component-wise optimization and facilitates debugging and enhancement of the system.

4 Experimental Results and Discussion

To evaluate the effectiveness of *SheetMind*, we conduct experiments on self-curated benchmark datasets consisting of both simple and complex spreadsheet manipulation tasks. As illustrated in Figure 3, SheetMind significantly outperforms baseline configurations, validating the utility of its multi-agent architecture.

We compare the full system (i.e., **Manager**, **Action**, and **Reflection** agents) with ablated versions removing one or more components. The evaluation focuses on task success rate as the primary metric. Key findings include:

- **Simple Tasks (single-step instructions):** The full SheetMind system achieved a success rate of approximately 80%, substantially outperforming the Action-Only variant at 20%. This underscores the value of structured task decomposition and validation, even for basic operations.
- **Complex Tasks (multi-step instructions):** As instruction complexity increased, the advantage of SheetMind became more pronounced. While the full system maintained a 70% success rate, removing the Manager Agent led to a sharp decline to 15%, highlighting the importance of hierarchical planning for multi-step execution.

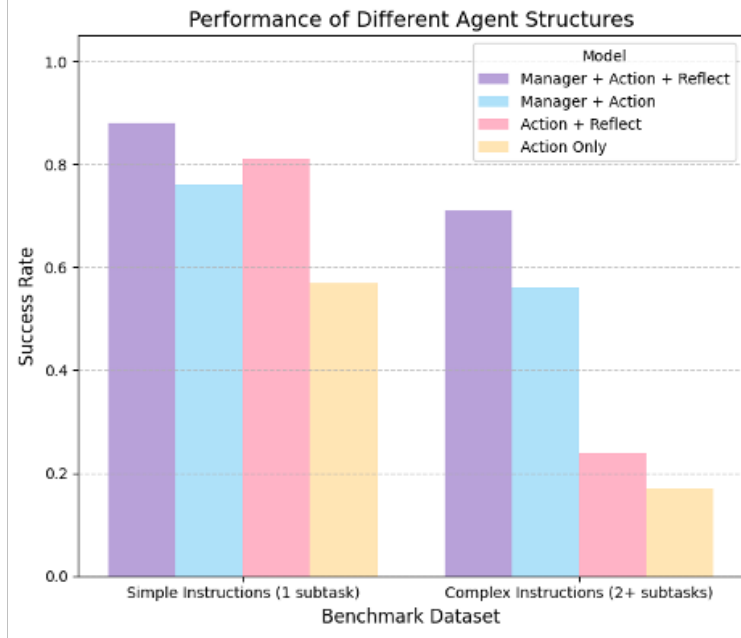


Figure 3: Task Success Rates Across Multi-Agent Architectures

- **Ablation Analysis:** Excluding the Reflection Agent resulted in a notable drop in performance due to the absence of runtime validation. This demonstrates the critical role of feedback and self-correction in enhancing execution robustness.

Beyond strong quantitative performance, SheetMind achieves key design objectives: accurate natural language interpretation, conditional task automation, real-time validation, and seamless Google Sheets integration. Its modular architecture enables easy extension to a wide range of spreadsheet operations.

The use of BNF-based grammar allows precise translation of natural language into structured commands (e.g., SELECT, SET, DELETE) with support for conditional logic and regex patterns. This abstraction cleanly separates intent parsing from execution, enhancing system robustness and maintainability. Overall, the results validate the effectiveness of combining multi-agent coordination with structured grammar for scalable, reliable spreadsheet automation.

5 Conclusions and Limitations

This work introduces *SheetMind*, an end-to-end multi-agent framework for translating natural language instructions into executable spreadsheet operations. Through hierarchical decomposition, grammar-based command synthesis, and reflective validation, SheetMind lowers the technical barrier to spreadsheet automation. Key design insights include:

- **Hierarchical task decomposition** via a Manager Agent enables robust handling of multi-step instructions.
- **BNF-style formal grammars** provide a reliable intermediate representation between natural language and executable actions.
- **Reflective validation** improves execution reliability through dynamic correction and intent preservation.

By integrating LLM-driven planning, structured execution, and agentic feedback, SheetMind delivers strong performance across a range of spreadsheet tasks while maintaining interpretability and extensibility.

While SheetMind demonstrates strong performance across diverse spreadsheet tasks, it currently relies on prompt-based coordination between agents and assumes deterministic spreadsheet environments. Future work may explore more autonomous agent collaboration, formal task benchmarking, and robustness under noisy or ambiguous user input.

References

- [1] Derek Flood, Kevin Mc Daid, and Fergal Mc Caffery. Nlp-sir: A natural language approach for spreadsheet information retrieval. *arXiv preprint arXiv:0908.1193*, 2009.
- [2] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Nee-lakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [3] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [4] Jonathan Herzig, Pawel Krzysztof Nowak, Thomas Mueller, Francesco Piccinno, and Julian Eisenschlos. Tapas: Weakly supervised table parsing via pre-training. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4320–4333, 2020.
- [5] Hongxin Li, Jinran Su, Yuntao Chen, Qing Li, and Zhaoxiang Zhang. Sheetcopilot: Bringing software productivity to the next level through large language models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [6] Zhengbao Jiang, Yi Mao, Pengcheng He, Graham Neubig, and Weizhu Chen. Omnitab: Pretraining with natural and synthetic data for few-shot table-based question answering. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 932–942, 2022.
- [7] Xinyun Chen, Petros Maniatis, Rishabh Singh, Charles Sutton, Hanjun Dai, Max Lin, and Denny Zhou. Spreadsheetcoder: Formula prediction from semi-structured context. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 1661–1672. PMLR, 18–24 Jul 2021.
- [8] Yibin Chen, Yifu Yuan, Zeyu Zhang, YAN ZHENG, Jinyi Liu, Fei Ni, Jianye HAO, Hangyu Mao, and Fuzheng Zhang. Sheetagent: Towards a generalist agent for spreadsheet reasoning and manipulation via large language models. In *THE WEB CONFERENCE 2025*, 2025.
- [9] Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, Chenyu Ran, Lingfeng Xiao, Chenglin Wu, and Jürgen Schmidhuber. MetaGPT: Meta programming for a multi-agent collaborative framework. In *The Twelfth International Conference on Learning Representations*, 2024.
- [10] Haowei Liu, Xi Zhang, Haiyang Xu, Yuyang Wanyan, Junyang Wang, Ming Yan, Ji Zhang, Chunfeng Yuan, Changsheng Xu, Weiming Hu, et al. Pc-agent: A hierarchical multi-agent collaboration framework for complex task automation on pc. *arXiv preprint arXiv:2502.14282*, 2025.
- [11] Daniel D McCracken and Edwin D Reilly. Backus-naur form (bnf). In *Encyclopedia of computer science*, pages 129–131. 2003.
- [12] Jenny T Liang, Aayush Kumar, Yasharth Bajpai, Sumit Gulwani, Vu Le, Chris Parnin, Arjun Radhakrishna, Ashish Tiwari, Emerson Murphy-Hill, and Guastavo Soares. Tabletalk: Scaffolding spreadsheet development with a language agent. *arXiv preprint arXiv:2502.09787*, 2025.
- [13] Zeyao Ma, Bohan Zhang, Jing Zhang, Jifan Yu, Xiaokang Zhang, Xiaohan Zhang, Sijia Luo, Xi Wang, and Jie Tang. Spreadsheetbench: Towards challenging real world spreadsheet manipulation. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024.
- [14] Jiayi Zhang, Jinyu Xiang, Zhaoyang Yu, Fengwei Teng, Xiong-Hui Chen, Jiaqi Chen, Mingchen Zhuge, Xin Cheng, Sirui Hong, Jinlin Wang, Bingnan Zheng, Bang Liu, Yuyu Luo, and Chenglin Wu. AFlow: Automating agentic workflow generation. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [15] Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chi-Min Chan, Heyang Yu, Yaxi Lu, Yi-Hsin Hung, Chen Qian, Yujia Qin, Xin Cong, Ruobing Xie, Zhiyuan Liu, Maosong Sun, and Jie Zhou. Agentverse: Facilitating multi-agent collaboration and exploring emergent behaviors. In *The Twelfth International Conference on Learning Representations*, 2024.