# FinTeam: A Multi-Agent Collaborative Intelligence System for Comprehensive Financial Scenarios

Yingqian Wu[1], Qiushi Wang[1], Zefei Long[1], Rong Ye[1], Zhongtian Lu[1], Xianyin Zhang[1], Bingxuan Li[1], Wei Chen[2], Liwen Zhang[3], and Zhongyu Wei[1,4,5*]

[1] School of Data Science, Fudan University, China
[2] School of Software Engineering, Huazhong University of Science and Technology, China
[3] Shanghai University of Finance and Economics, China
[4] Shanghai Innovation Institute, China
[5] Research Institute of Intelligent Complex Systems, Fudan University, China

{wuyq23,qswang23,zflong23,yer23,ztlu22,xianyinzhang22,bxli16}@m.fudan.edu.cn,
lemuria_chen@hust.edu.cn, zhang.liwen@shufe.edu.cn, zywei@fudan.edu.cn

**Abstract.** Financial report generation tasks range from macro- to micro-economics analysis, also requiring extensive data analysis. Existing LLM models are usually fine-tuned on simple QA tasks and cannot comprehensively analyze real financial scenarios. Given the complexity, financial companies often distribute tasks among departments. Inspired by this, we propose FinTeam, a financial multi-agent collaborative system, with a workflow with four LLM agents: *document analyzer*, *analyst*, *accountant*, and *consultant*. We train these agents with specific financial expertise using constructed datasets. We evaluate FinTeam on comprehensive financial tasks constructed from real online investment forums, including macroeconomic, industry, and company analysis. The human evaluation shows that by combining agents, the financial reports generate from FinTeam achieved a 62.00% acceptance rate, outperforming baseline models like GPT-4o and Xuanyuan. Additionally, FinTeam's agents demonstrate a 7.43% average improvement on FinCUGE and a 2.06% accuracy boost on FinEval. Project is available at https://github.com/FudanDISC/DISC-FinLLM/.

**Keywords:** Multi-Agent System · Financial LLM · Agent Instruction Tuning

## 1 Introduction

The financial industry presents unique challenges and opportunities for Large Language Models (LLMs). Recently, domain-specialized LLMs have achieved notable progress in several vertical fields, such as legal reasoning [32] and medical consultation [1]. While LLMs like BloombergGPT [26] and XuanYuan [36] perform well on routine tasks, they struggle with complex, multifaceted financial problems. Single LLM calls often lack the precision needed for detailed financial analysis. In practice, as shown in Figure 1, financial tasks are typically handled by specialists, highlighting the potential of decomposing tasks into sub-tasks managed by dedicated agents.

---

* Corresponding author.

**Fig. 1.** Inspired by the financial companies that assign tasks to specialized teams, FinTeam distributes financial tasks among the *document analyzer*, *analyst*, *accountant*, and *consultant* agents, enabling a more efficient and sophisticated process.

In response to these challenges, we introduce FinTeam, a financial intelligence system composed of multiple collaborating LLM agents, each designed to address specific scenarios in finance. We focus on three key financial scenarios: macroeconomic analysis, industry analysis, and company analysis. These scenarios encompass a range of tasks, from assessing broad economic trends to detailed evaluations of individual companies. Within this framework, four specialized LLM agents collaborate to handle their respective tasks and provide comprehensive solutions.

The four LLM agents focus on processing financial texts such as news and company reports, performing real-time material analysis based on knowledge bases, using computational tools to achieve financial numerical calculation, and professionally answering a variety of financial questions.

To demonstrate FinTeam's effectiveness, we conduct extensive evaluations on both the collaborative system and individual agents. Using a dataset of 150 real investor queries spanning the three scenarios, we compare FinTeam's performance to that of baseline models, including Qwen2.5-7B-Instruct, GPT-4o, ChatGLM3-6B and Xuanyuan-13B. The results, evaluated by GPT-4o, show that FinTeam has an overall score of 4.86 (out of 5), which is significantly better than the other models by 0.03 to 0.82 points. Furthermore, human preference evaluations confirm FinTeam's real-world relevance, with a winning rate of 62.00%. At the same time, the performance of each agent is also significantly better than the other models in each evaluation.

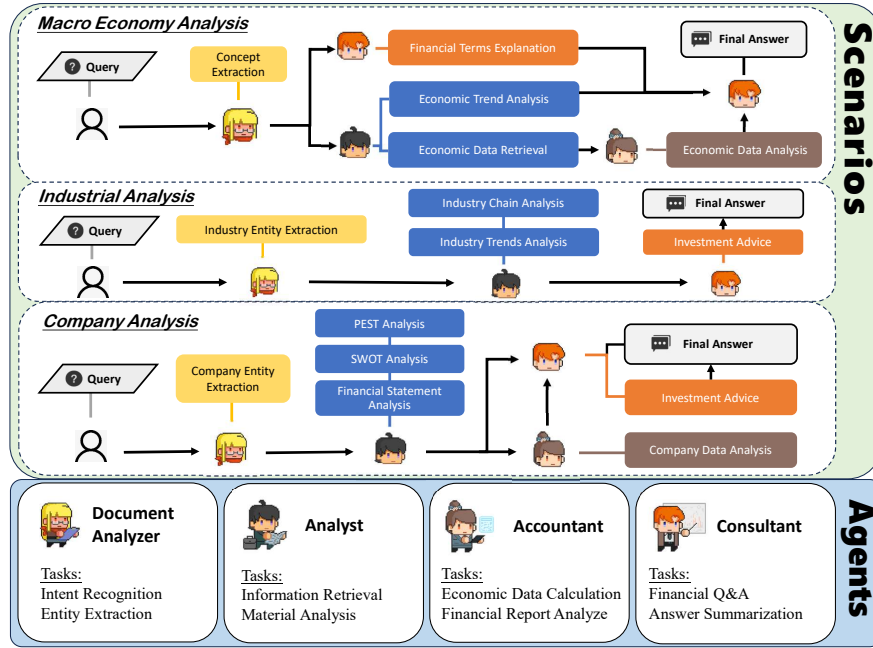In summary, our contributions are as follows:

**Fig. 2.** Overview of FinTeam Multi-Agent Collaborative Intelligence System.

– We design FinTeam, a financial intelligence system comprising multiple collaborating LLM agents: the *analyst*, *document analyzer*, *accountant*, and *consultant*. Each agent addresses specific challenges in various financial scenarios.
– To strengthen their capabilities, we develop the agent training dataset used for fine-tuning. Extensive evaluations demonstrate FinTeam's effectiveness and the value of the dataset.
– Unlike prior evaluations that focus on simple financial tasks, we concentrate on three comprehensive real-world application scenarios: macroeconomic analysis, industry analysis, and company analysis. Both automated evaluations and human assessments demonstrate the effectiveness of FinTeam's workflow.

## 2   Related Work

### 2.1   LLMs for Finance

Large Language Models (LLMs), typically with over a billion parameters, have greatly advanced natural language processing across a wide range of domains. In the financial sector, LLMs offer promising capabilities for understanding complex documents, generating investment insights, and enabling data-driven decision-making. A growing number of specialized financial LLMs have been developed, including BloombergGPT [26], DISC-FinLLM [3], and XuanYuan [36], all of which are trained or adapted on corpora reflecting financial discourse and context.

Beyond general-purpose models, several financial LLMs have been tailored to support specific tasks and modalities. PIXIU [27], fine-tuned from LLaMA [20], targets structured financial tasks such as risk assessment and entity linking. FinVis-GPT [23] incorporates multimodal chart interpretation for visual-grounded analysis. InvestLM [29] focuses on deep financial reasoning using curated QA datasets.

However, these models operate under a single-agent architecture, which limits their capacity to decompose and solve complex, multi-step financial tasks. This motivates exploration into multi-agent systems with modular and collaborative capabilities.

### 2.2   Multi-Agent Collaboration

Multi-agent systems have emerged as promising solutions for tackling complex tasks. These systems employ strategies such as role-playing, collaboration, and task decomposition to improve problem-solving efficiency. For instance, AutoGen [25] provides an open framework that enables agent communication for LLM-based applications, while MetaGPT [7] adopts an assembly-line paradigm where specialized agents execute structured subtasks.

In the financial domain, recent studies have explored agent-based systems to support investment and trading decisions. TradingGPT [11] models agents with different risk profiles and strategies. FinMem [30] incorporates profiling, memory, and decision-making modules to improve cumulative returns. FinAgent [35] integrates image-based financial data into agent interactions to enhance trading decisions. However, most efforts remain trading-focused, lacking broader applications in macroeconomic, industry-level, and company-level financial analysis.

Recent work in other domains has demonstrated that multi-agent collaboration enhances performance in complex tasks. SMART [33] leverages trajectory-based coordination to improve factual consistency in knowledge-intensive scenarios. MASER [31] simulates legal interactions with role-aligned agents. In the medical domain, MDA-gents [10] and AI Hospital [4] design adaptive agent collaborations for clinical reasoning and diagnosis, demonstrating gains in multi-turn interaction and decision-making. These works highlight the effectiveness of structured cooperation and role specialization, principles that motivate our multi-agent design for financial analysis.

## 3   FinTeam

To meet the needs of practical financial scenarios, we propose a multi-agent collaborative financial intelligence system, FinTeam. This system organizes a virtual team of financial agents to handle complex tasks through agent interactions. Four roles are defined: *document analyst*, *analyst*, *accountant*, and *consultant*, each specializing in specific financial skills via supervised training. The construction and statistics of the agent training dataset are shown in Figure 3 and Table 2. Users can deploy these agents individually to handle specific financial tasks or enable collaboration across three main scenarios [17]—macroeconomic, industry, and company analysis—to tackle complex financial challenges. The overall overview of our system is shown in Figure 2.
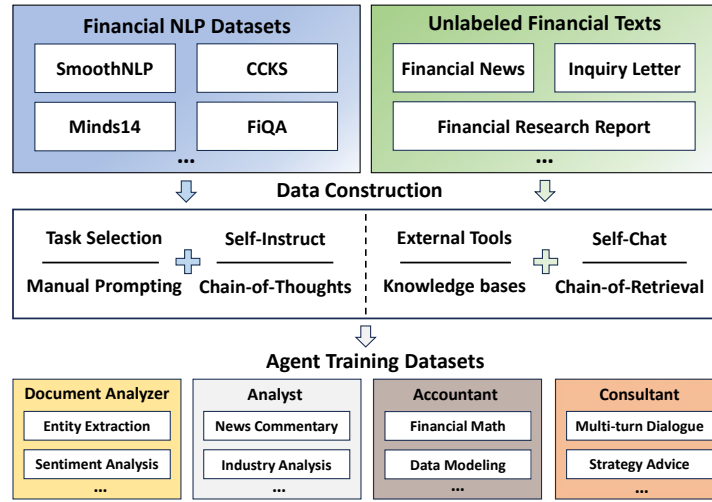
**Fig. 3.** Construction of agent training dataset.

## 3.1 Agent Roles

**Document Analyzer** The financial industry generates vast amounts of data daily, particularly unstructured text data such as news, market commentaries. In our system, the *document analyzer* serves as an agent designed to processing financial texts, capable of performing NLP tasks such as intent recognition, financial entity extraction, and financial sentiment analysis, tailored to specific financial texts and user requirements.

We use domain-specific NLP datasets to train the *document analyzer*. The dataset comprises two parts: labeled open-source datasets and unlabeled texts in financial reports that are automatically annotated by ChatGPT. The details of NLP datasets is shown in Table 9 in Appendix 7.1.

**Analyst** The financial field evolves rapidly, requiring analysis in context of current events. Our system's *analyst* uses Retrieval-Augmented Generation (RAG) to analyze real-time financial materials—news, policies, reports, and data—retrieved from multiple knowledge bases via the m3e-base [24] embedding model.

We propose the Chain-of-Retrieval (CoR) prompting methodology to create a financial analysis instruction dataset, consisting of three main steps: **1) Question Generation:** Formulate financial analysis questions based on financial contexts. **2) Reference Retrieval:** Retrieve relevant documents from knowledge bases. **3) Answer Generation:** Combine questions with retrieved contexts to generate answers. This dataset includes financial questions, reference documents, and their corresponding analyses, all generated by ChatGPT.

**Accountant** Financial texts, particularly reports, often contain complex numerical data, requiring calculations like growth rates or expected earnings. Since LLMs struggle with

accurate calculations [6], we introduce accountant, a tool-augmented agent that utilizes computational tools to address this limitation, following the Toolformer approach [18].

*Accountant* is equipped with four tools, as shown in Table 1, which cover various computing tasks with different commands, inputs, and outputs. For instance, the calculator command is [Calculator(expression)→result]. When a quantitative analysis query is received, accountant generates the appropriate tool call command, halts decoding, executes the calculation, and integrates the result into the response before continuing with the text generation.

For the training set construction of *accountant*, we create a seed task pool with three components: financial calculation questions from exams, arithmetic questions within financial report contexts, and general math questions from BELLE School Math [9]. The answers to these questions include embedded tool commands specifying their usage. Additionally, we use ChatGPT [15] to generate over 50,000 new question-answer pairs using self-instruction [22] and few-shot Chain-of-Thought (CoT) prompting, with answers incorporating calculation tool commands.

| TOOL | DETAIL |
|---|---|
| Expression calculator | Input: expression; Output: result |
| Equation solver | Input: equation system; Output: solution |
| Counter | Input: array of samples; Output: sample size |
| Probability table | Input: number; Output: cumulative standard normal distribution function value at this number |

**Table 1.** Definition of calculation tools

**Consultant**  The finance domain is vast, covering areas like banking, stock trading, and futures. A financial intelligence system needs extensive background knowledge to perform detailed analysis. To address this, we developed an agent called *consultant*, designed to interpret and answer finance-related queries.

To equip the *consultant* with solid financial knowledge and consulting skills, we construct a Chinese-language dataset. We begin by translating the FiQA dataset [13] and use only the Chinese version for training. To enhance domain understanding, we generate QA pairs from 200 finance-specific terms using ChatGPT [15]. We also scrape user queries from the Chinese financial forum JingGuan[6], and apply a self-chat approach [28] to create multi-turn dialogues based on seed topics.

| AGENT | #SAMPLES | INPUT LENGTH | OUTPUT LENGTH |
|---|---|---|---|
| *Consultant* | 63k | 26 | 370 |
| *Document Analyzer* | 95k | 586 | 31 |
| *Accountant* | 64k | 70 | 204 |
| *Analyst* | 20k | 1023 | 521 |
| Total | 241k | 340 | 206 |

**Table 2.** Statistics of agent training dataset. The lengths are the average number of tokens.

---

[6] https://bbs.pinggu.org/

### 3.2   Scenario Settings

To provide comprehensive financial analysis, we design three core scenarios—ranging from macro to micro levels—illustrated in Figure 2. These scenarios enable flexible, tailored responses to user needs.

**Macroeconomic Analysis.** This scenario focuses on macroeconomic theories and trends. FinTeam operates as follows: 1) the *document analyzer* extracts key terms; 2) the *consultant* explains them; 3) the *analyst* gathers and summarizes supplementary data; 4) the *consultant* compiles a final response. This process helps users stay informed about economic developments and make better investment decisions.

**Industry Analysis.** This scenario addresses specific sectors and trends. FinTeam proceeds as follows: 1) the *document analyzer* identifies relevant industries or companies; 2) the *analyst* explores competition, supply chains, and development trends; 3) recent news is included if requested; 4) the *consultant* summarizes insights and offers strategic suggestions. Users gain a clear understanding of industry dynamics and outlook.

**Company Analysis.** This scenario evaluates individual companies using models like PEST and SWOT. FinTeam works as follows: 1) the *document analyzer* extracts key company data; 2) the *analyst* applies PEST, SWOT, and performs sentiment analysis if needed; 3) the *consultant* delivers a synthesized assessment.

Additionally, this scenario includes **financial statement analysis**: 1) the *analyst* retrieves data from balance sheets, income statements, and cash flow reports; 2) the *accountant* calculates key ratios such as profitability, liquidity, and leverage; 3) the *consultant* generates a concise, actionable report.

Overall, this scenario provides structured, in-depth insights into a company's financial health, market position, and strategic outlook, supporting informed investment decisions.

## 4   Experiments

### 4.1   Evaluation Data and Setups

To assess the performance of our FinTeam in real-world scenarios, we gather 150 actual investor inquiries from the popular Chinese online investment forum, NGA Grand Era [7]. For the three major scenarios, we ensure diversity in the collected questions by applying rule-based filtering, ultimately retaining fifty test questions per scenario.

In the macroeconomic scenario, questions address popular topics such as changes in economic indicators, asset price fluctuations, market interest rate variations, and global financial policy news. The industry scenario encompasses inquiries from 27 sub-sectors, including industry news evaluations and investment trends. The company scenario focuses on highly followed publicly listed companies, involving news, earnings reports, and stock price fluctuations.

For evaluation, we employ GPT-4o to score the outputs from our agent system and the other models, ensuring objectivity and accuracy. The evaluation is conducted across four dimensions: **(1) Accuracy**: The model addresses key points directly, avoiding irrelevant details. **(2) Thoroughness**: The model provides a detailed, in-depth answer.

---

[7] https://ngabbs.com/

**(3) Clarity**: The response is clear, concise, and logical. **(4) Professionalism**: The model uses appropriate financial perfessional terms.

We compare model outputs across multiple dimensions, with GPT-4o rating each response from 1 to 5 per category and assigning an overall score. Pairwise significance tests are conducted to confirm statistical improvements. To validate GPT-4o's judgments, human evaluations are also performed by finance undergraduates. They anonymously assess responses and select the best ones. The Acceptance Rate indicates how often a model's output is chosen as the top answer.

## 4.2   Implementation Details

We train the agents using the LoRA mechanism [8] on the Qwen2.5-7B-Chat model with Deepspeed ZeRO-0 [21] on four NVIDIA V100-32G GPUs. The batch size is 1 per GPU, with gradient accumulation steps of 4, a maximum sequence length of 4096, and 2 epochs. The learning rate is $5 \times 10^{-5}$ and follows a cosine annealing schedule. LoRA parameters are set with a target of "all", rank 8, and alpha 16.

## 4.3   Main Results

We compare performance of FinTeam with the baseline models, with results presented in Table 3. From these evaluations, we observe that our model achieve a 0.13 improvement in overall score when answering Chinese financial questions, compared to the baseline. The most significant gains are in thoroughness and professionalism, where the model improve by 0.23 points in both categories. Additionally, our financial agent collaboration system outperform GPT-3.5-turbo and Xuanyuan-13B across all dimensions, demonstrating its effectiveness.The results of significance tests, presented in Table 4, confirm that the improvements in thoroughness, professionalism, and overall score are highly significant, with p-values well below the accepted thresholds, demonstrating the robustness of our system's performance.

| MODEL | ACC. | THO. | CLA. | PRO. | OVERALL |
|---|---|---|---|---|---|
| **FinTeam (ours)** | 4.54 | **4.94** | 4.84 | **4.96** | **4.86** |
| **Qwen2.5-7B-Chat** | 4.51 | 4.69 | **4.99** | 4.80 | 4.78 |
| **GPT-4o** | **4.67** | 4.73 | 4.99 | 4.85 | 4.83 |
| **ChatGLM3-6B** | 3.95 | 3.75 | 4.68 | 3.91 | 4.04 |
| **Xuanyuan-13B** | 4.35 | 4.61 | 4.96 | 4.67 | 4.66 |

**Table 3.** Evaluation results across different dimensions. FinTeam achieves the highest overall score, surpassing the baseline model by 0.08, with notable improvements of 0.25 in thoroughness and 0.16 in professionalism. These results highlight the effectiveness of FinTeam.

| METRIC | ACC. | THO. | CLA. | PRO. | OVERALL |
|---|---|---|---|---|---|
| **t-statistic** | 0.253 | 5.445 | -5.195 | 4.218 | 2.493 |
| **p-value** | 0.8005 | **0.0000** | **0.0000** | **0.0000** | **0.0138** |

**Table 4.** Statistical significance of model evaluation metrics. The results indicate that the improvements in thoroughness, professionalism, and overall score are statistically significant.

The final count of human evaluation picks is shown in Table 5. FinTeam significantly outperforms the other models, achieving an acceptance rate of 62.00%. The results are consistent with the GPT-4o assessment, further validating our system's reliability. It can be concluded that FinTeam is capable of providing professional and thorough answers to users' questions in real financial scenarios, offering users an in-depth understanding across multiple materials.

| Metric | FinTeam (ours) | Qwen2.5-7B-Chat | GPT-4o | ChatGLM3-6B | Xuanyuan-13B |
|---|---|---|---|---|---|
| Acceptance Rate | **62.00%** | 9.33% | 5.33% | 4.00% | 19.33% |

**Table 5.** Human evaluation on the preference of model outputs. Acceptance Rate indicates how often a model's answer is selected as the best among all models. FinTeam is chosen 93 times out of 150 test cases, achieving a selection rate of 62.00%.

## 5    Analysis

What causes our FinTeam to give better generation quality? We specifically analyze how each LLM agent performs in three different aspects.

### 5.1    Evaluation Setup

*Financial NLP Tasks*  We assess the model's NLP ability using the FinCUGE benchmark [12] across six tasks: sentiment analysis (FinFE), event entity (FinQA), causality extraction (FinCQA), summarization (FinNA), relation extraction (FinRE), and entity extraction (FinESE). We create a few-shot evaluation setting with prompts and measure performance using accuracy, F1 score, and ROUGE score.

*Chinese Financial Knowledge Tests*   To evaluate our model's performance on Chinese financial knowledge, we utilize the FinEval [34], which covers 34 financial subcategories, containing a total of 1,151 multiple-choice questions. FinEval is Out-of-Distribution for our dataset, so it can adequately assess the generalization ability of our model and dataset. We measure the performance by calculating the accuracy of multiple-choice questions.

*Data Analysis*  We manually created a dataset of 100 financial calculation problems, adapted from material analysis questions in the Chinese Administrative Aptitude Test, to evaluate our model's capabilities. The dataset is crafted for quality assurance, and the model's performance is assessed based on accuracy in formula construction and result calculation.

### 5.2    Results on Financial NLP Tasks

Table.6 presents the performance of various models across six financial NLP tasks. The Document analyzer demonstrates the highest performance in all task-specific metrics.

With an average score of 47.20, it significantly outperforms the strong baseline Qwen2.5-7B-Instruct, which achieves an average of 39.77. This represents a notable 7.43-point improvement, underscoring the effectiveness and robustness of the Document analyzer in handling diverse financial text understanding and reasoning tasks.

| MODEL | FINFE (ACC) | FINQA (F1) | FINCQA (F1) | FINNA (ROUGE) | FINRE (ACC) | FINESE (F1) | AVG |
|---|---|---|---|---|---|---|---|
| **Document analyzer** | **66.99** | **47.44** | **46.12** | **41.10** | **22.36** | **61.19** | **47.20** |
| **Qwen2.5-7B-Instruct** | 66.19 | 37.33 | 34.50 | 40.80 | 21.02 | 36.76 | 39.77 |
| **ChatGLM3-6B** | 62.57 | 25.69 | 21.41 | 30.60 | 8.39 | 28.44 | 29.18 |
| **Xuanyuan-13B** | 62.48 | 18.07 | 24.94 | 31.10 | 9.13 | 31.08 | 29.47 |

**Table 6.** Performance comparison across six financial tasks on FinCUGE benchmark.*Document analyzer* outperforms the base model Qwen2.5-7B-Instruct by 7.43 points on average.

### 5.3    Results on Financial Knowledge Tests

Table 7 shows the evaluation results of our four LLM agents compared to general and financial LLMs on the FinEval benchmark. This demonstrates their extensive financial knowledge, strong task performance, and adaptability across diverse financial scenarios. Additionally, since FinEval is an out-of-distribution for our dataset, it highlights the universality of our training tasks and dataset.

| MODEL | *Consultant* | QWEN2.5-7B-INSTRUCT | CHATGLM3-6B | GPT-4O | XUANYUAN-13B |
|---|---|---|---|---|---|
| **FinEval(Acc)** | 68.48 | 66.42 | 47.64 | **70.67** | 35.15 |

**Table 7.** Experimental results of average scores on the FinEval benchmark. *Consultant* outperforms the base model Qwen2.5-7B-Instruct by 2.06 points.

### 5.4    Results on Data Analysis

Table. 8 showcases the experiment results on financial computing tasks. The addition of computational plugins to our model generates a notable performance boost compared to the baseline models, surpassing it by 0.09 in formula&results. These results highlight the efficacy of our approach in addressing computational challenges within the financial domain.

| MODEL | FORMULA | FORMULA & RESULT |
|---|---|---|
| *Accountant* | 0.70 | 0.70 |
| Qwen2.5-7B-Instruct | 0.65 | 0.61 |
| GPT-4o | 0.81 | 0.81 |

**Table 8.** Evaluation results of financial calculation. Formula represents the accuracy of the formula in the calculation process, while Formula&Result denotes the accuracy of both formulas and results. *Accountant* is 0.09 higher than the base model in Formula&Result accuracy.

# 6   Conclusion

In this paper, in order to meet the needs of users in actual financial scenarios, we propose a financial intelligence system FinTeam, which connects multiple subtasks through the interaction between LLM agents for enhancing the system's ability to handle complex tasks. We construct the agent training dataset and train four LLM agents using different sub-datasets, which enables the agents to complete complex financial tasks in three scenarios: macroeconomic analysis, industry analysis, and company analysis following collaborative workflows. Within the framework of financial evaluation, we establish multi-dimensional benchmarks and demonstrate their robust capabilities, underscoring the capability of FinTeam to offer substantial support across various financial scenarios.

**Limitations**  Our work has several limitations. The scenario design is limited in scope, leaving out many financial tasks for future exploration. As the system may generate investment-related advice, users should be cautious, as financial outcomes are not guaranteed. The system is also tailored to Chinese financial contexts, and its effectiveness in global markets remains untested.

# Acknowledgements

# References

1. Bao, Z., Chen, W., Xiao, S., Ren, K., Wu, J., Zhong, C., Peng, J., Huang, X., Wei, Z.: Disc-medllm: Bridging general large language models and real-world medical consultation (2023)
2. Biendata: ccksnec2022 (2022), https://www.biendata.xyz/competition/ccks-nec-2022
3. Chen, W., Wang, Q., Long, Z., Zhang, X., Lu, Z., Li, B., Wang, S., Xu, J., Bai, X., Huang, X., et al.: Disc-finllm: A chinese financial large language model based on multiple experts fine-tuning. arXiv preprint arXiv:2310.15205 (2023)
4. Fan, Z., Tang, J., Chen, W., Wang, S., Wei, Z., Xi, J., Huang, F., Zhou, J.: Ai hospital: Benchmarking large language models in a multi-agent medical interaction simulator. arXiv preprint arXiv:2402.09742 (2024)
5. Gerz, D., Su, P.H., Kusztos, R., Mondal, A., Lis, M., Singhal, E., Mrkšić, N., Wen, T.H., Vulić, I.: Multilingual and cross-lingual intent detection from spoken data. In: Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing. pp. 7468–7475 (2021)
6. Golkar, S., Pettee, M., Eickenberg, M., Bietti, A., Cranmer, M., Krawezik, G., Lanusse, F., McCabe, M., Ohana, R., Parker, L., Blancard, B.R.S., Tesileanu, T., Cho, K., Ho, S.: xval: A continuous number encoding for large language models (2023), https://arxiv.org/abs/2310.02989
7. Hong, S., Zhuge, M., Chen, J., Zheng, X., Cheng, Y., Zhang, C., Wang, J., Wang, Z., Yau, S.K.S., Lin, Z., Zhou, L., Ran, C., Xiao, L., Wu, C., Schmidhuber, J.: Metagpt: Meta programming for a multi-agent collaborative framework (2023), https://arxiv.org/abs/2308.00352
8. Hu, E.J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W.: Lora: Low-rank adaptation of large language models (2021), https://arxiv.org/abs/2106.09685

9. Ji, Y., Deng, Y., Gong, Y., Peng, Y., Niu, Q., Zhang, L., Ma, B., Li, X.: Exploring the impact of instruction data scaling on large language models: An empirical study on real-world use cases. arXiv preprint arXiv:2303.14742 (2023)

10. Kim, Y., Park, C., Jeong, H., Chan, Y.S., Xu, X., McDuff, D., Lee, H., Ghassemi, M., Breazeal, C., Park, H.W.: Mdagents: An adaptive collaboration of llms for medical decision-making. Advances in Neural Information Processing Systems **37**, 79410–79452 (2024)

11. Li, Y., Yu, Y., Li, H., Chen, Z., Khashanah, K.: TradingGPT: Multi-Agent System with Layered Memory and Distinct Characters for Enhanced Financial Trading Performance. arXiv e-prints arXiv:2309.03736 (Sep 2023). https://doi.org/10.48550/arXiv.2309.03736

12. Lu, D., Wu, H., Liang, J., Xu, Y., He, Q., Geng, Y., Han, M., Xin, Y., Xiao, Y.: BBT-Fin: Comprehensive Construction of Chinese Financial Domain Pre-trained Language Model, Corpus and Benchmark. arXiv e-prints arXiv:2302.09432 (Feb 2023). https://doi.org/10.48550/arXiv.2302.09432

13. Maia, M., Handschuh, S., Freitas, A., Davis, B., McDermott, R., Zarrouk, M., Balahur, A.: Www'18 open challenge: financial opinion mining and question answering. In: Companion proceedings of the the web conference 2018. pp. 1941–1942 (2018)

14. Malo, P., Sinha, A., Korhonen, P., Wallenius, J., Takala, P.: Good debt or bad debt: Detecting semantic orientations in economic texts. Journal of the Association for Information Science and Technology **65**(4), 782–796 (2014)

15. OpenAI: Chatgpt. https://openai.com/blog/chatgpt (2023)

16. Ren, J., Wang, S., Song, R., Wu, Y., Gao, Y., An, B., Cheng, Z., Xu, G.: Iree: A fine-grained dataset for chinese event extraction in investment research. In: China Conference on Knowledge Graph and Semantic Computing. pp. 205–210. Springer (2022)

17. Robbins, S.P., Coulter, M.: Management. Pearson Prentice Hall, Upper Saddle River, NJ, 11th edn. (2012)

18. Schick, T., Dwivedi-Yu, J., Dessi, R., Raileanu, R., Lomeli, M., Zettlemoyer, L., Cancedda, N., Scialom, T.: Toolformer: Language models can teach themselves to use tools (2023)

19. Tianchi: ccks2022event (2022), https://tianchi.aliyun.com/dataset/dataDetail?dataId=136800

20. Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E., Lample, G.: LLaMA: Open and Efficient Foundation Language Models. arXiv e-prints arXiv:2302.13971 (Feb 2023). https://doi.org/10.48550/arXiv.2302.13971

21. Wang, G., Qin, H., Ade Jacobs, S., Holmes, C., Rajbhandari, S., Ruwase, O., Yan, F., Yang, L., He, Y.: ZeRO++: Extremely Efficient Collective Communication for Giant Model Training. arXiv e-prints arXiv:2306.10209 (Jun 2023). https://doi.org/10.48550/arXiv.2306.10209

22. Wang, Y., Kordi, Y., Mishra, S., Liu, A., Smith, N.A., Khashabi, D., Hajishirzi, H.: Self-instruct: Aligning language model with self generated instructions (2022)

23. Wang, Z., Li, Y., Wu, J., Soon, J., Zhang, X.: Finvis-gpt: A multimodal large language model for financial chart analysis (2023), https://arxiv.org/abs/2308.01430

24. Wang Yuxin, Sun Qingxuan, H.s.: M3e: Moka massive mixed embedding model (2023)

25. Wu, Q., Bansal, G., Zhang, J., Wu, Y., Li, B., Zhu, E., Jiang, L., Zhang, X., Zhang, S., Liu, J., Awadallah, A.H., White, R.W., Burger, D., Wang, C.: Autogen: Enabling next-gen llm applications via multi-agent conversation (2023), https://arxiv.org/abs/2308.08155

26. Wu, S., Irsoy, O., Lu, S., Dabravolski, V., Dredze, M., Gehrmann, S., Kambadur, P., Rosenberg, D., Mann, G.: Bloomberggpt: A large language model for finance. arXiv preprint arXiv:2303.17564 (2023)

27. Xie, Q., Han, W., Zhang, X., Lai, Y., Peng, M., Lopez-Lira, A., Huang, J.: Pixiu: A large language model, instruction data and evaluation benchmark for finance. arXiv preprint arXiv:2306.05443 (2023)

28. Xu, C., Guo, D., Duan, N., McAuley, J.: Baize: An open-source chat model with parameter-efficient tuning on self-chat data. arXiv preprint arXiv:2304.01196 (2023)

29. Yang, Y., Tang, Y., Tam, K.Y.: Investlm: A large language model for investment using financial domain instruction tuning (2023), https://arxiv.org/abs/2309.13064
30. Yu, Y., Li, H., Chen, Z., Jiang, Y., Li, Y., Zhang, D., Liu, R., Suchow, J.W., Khashanah, K.: Finmem: A performance-enhanced llm trading agent with layered memory and character design (2023), https://arxiv.org/abs/2311.13743
31. Yue, S., Huang, T., Jia, Z., Wang, S., Liu, S., Song, Y., Huang, X.J., Wei, Z.: Multi-agent simulator drives language models for legal intensive interaction. In: Findings of the Association for Computational Linguistics: NAACL 2025. pp. 6537–6570 (2025)
32. Yue, S., Liu, S., Zhou, Y., Shen, C., Wang, S., Xiao, Y., Li, B., Song, Y., Shen, X., Chen, W., et al.: Lawllm: Intelligent legal system with legal reasoning and verifiable retrieval. In: International Conference on Database Systems for Advanced Applications. pp. 304–321. Springer (2024)
33. Yue, S., Wang, S., Chen, W., Huang, X., Wei, Z.: Synergistic multi-agent framework with trajectory learning for knowledge-intensive tasks. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 39, pp. 25796–25804 (2025)
34. Zhang, L., Cai, W., Liu, Z., Yang, Z., Dai, W., Liao, Y., Qin, Q., Li, Y., Liu, X., Liu, Z., et al.: Fineval: A chinese financial domain knowledge evaluation benchmark for large language models. arXiv preprint arXiv:2308.09975 (2023)
35. Zhang, W., Zhao, L., Xia, H., Sun, S., Sun, J., Qin, M., Li, X., Zhao, Y., Zhao, Y., Cai, X., Zheng, L., Wang, X., An, B.: A multimodal foundation agent for financial trading: Tool-augmented, diversified, and generalist (2024), https://arxiv.org/abs/2402.18485
36. Zhang, X., Yang, Q., Xu, D.: Xuanyuan 2.0: A large chinese financial chat model with hundreds of billions parameters. arXiv preprint arXiv:2305.12002 (2023)

# 7 Appendix

## 7.1 Information of NLP datasets

To train the *document analyzer* in NLP tasks, we use various datasets, including sentiment analysis, entity extraction, and summarization, to enhance its capabilities. The open-source datasets used are listed in Figure 9. To ensure balance, we randomly sample 10,000 samples from subsets with over 10,000 samples.

| Dataset | Task Type | # Samples |
|---|---|---|
| FPB [14] | Sentiment Analysis | 18690 |
| CCKS-NEC-2022 [2] | Causality Extraction | 7499 |
| SmoothNLP IEE [3] | Event Extraction | 3256 |
| SmoothNLP NHG [3] | Text Generation | 4642 |
| CCKS2022-event [19] | Event Classification | 3578 |
| Minds14 [5] | Intent Prediction | 59143 |
| Financial Report | Entity Extraction | 61705 |
| OpenKG [16] | Entity Extraction | 7672 |
| OpenKG | Entity Extraction | 67921 |
| Wealth-alpaca-lora [4] | Keyword Generation | 41825 |

**Table 9.** Data statistics of our financial NLP datasets.

---

[1] https://github.com/wwwxmu/Dataset-of-financial-news-sentiment-classification

[2] https://github.com/smoothnlp/FinancialDatasets

[3] https://huggingface.co/datasets/gbharti/finance-alpaca