

EnterpriseBench: Benchmarking LLM Agents on Enterprise-Level Strategic Reasoning and Decision-Making

Anonymous ACL submission

Abstract

As Large Language Model (LLM) agents demonstrate increasingly strong reasoning capabilities, they are being progressively adopted to support complex decision-making in enterprise environments. However, existing benchmarks primarily evaluate reasoning through static, single-shot tasks with well-defined objectives and immediate correctness, whereas real-world decision-making is inherently interactive, open-ended, and shaped by delayed consequences and competing goals. To address this challenge, we introduce **EnterpriseBench**, a comprehensive benchmark for assessing LLM agents in realistic enterprise contexts. EnterpriseBench covers a hierarchy of reasoning demands, ranging from information extraction, numerical and domain knowledge reasoning to high-fidelity interactive decision-making, including management consulting cases and serious games. We also introduce an agent-oriented taxonomy that organizes tasks by capability domains and intrinsic difficulty. Empirical evaluations across nine state-of-the-art agent architectures reveal a natural division of labor, where simpler tasks favor lightweight models and complex decision-making benefits from stronger reasoning agents. Building on this insight, we present **AOA**, an Agent-of-Agents framework that integrates the complementary capabilities of such specialized agents. Our codes and benchmark implementation are publicly available¹.

1 Introduction

Recent advances in Large Language Models (LLMs) and LLM-based agents have substantially enhanced their reasoning capabilities, enabling them to address increasingly complex decision-making problems (Liu et al., 2024; Yang et al., 2025; Zhang et al., 2025a, 2024; Yu et al., 2024;

Li et al., 2025; Xing, 2025). However, most existing benchmarks (Zhang et al., 2025c; Mohammadi et al., 2025; Zhu et al., 2024; Chen et al., 2021) still focus on static and well-defined tasks such as mathematical reasoning (Glazer et al., 2024; Fan et al., 2024), programming (Jimenez et al., 2023), and question answering (Chen et al., 2021), where conditions are fixed and correctness can be verified against a single ground-truth answer. However, real-world decision-making is far more complex, a gap that becomes especially evident when LLM agents are deployed in enterprise settings. For example, when a firm considers expanding production capacity, it must reason under uncertain future demand, anticipate competitors’ responses, allocate limited capital across alternative investments, and balance short-term costs against long-term profitability. Such decisions are ubiquitous in real-world enterprise settings, shape future states of the environment, involve delayed and potentially irreversible consequences, and rarely lead to a single correct solution.

Fortunately, recent LLM agents have begun to demonstrate surprising competence in such complex decision-making scenarios (Raptis et al., 2025). LLM agents have been explored for investment decision-making tasks (Xiao et al., 2024; Chen et al., 2025b). In these settings, agents must reason over incomplete and noisy information, make sequential decisions, and revise their strategies as market conditions evolve. However, although LLM agents have demonstrated such capabilities and the potential to simulate firms engaging in complex reasoning and decision-making, existing benchmarks remain inadequate for rigorously evaluating their performance in real-world scenarios.

Evaluating real-world decision-making is inherently challenging, as such tasks often lack well-defined datasets, objective ground-truth answers, or clear evaluation boundaries, in contrast to math-

¹<https://anonymous.4open.science/r/EnterpriseBench>

emational (Glazer et al., 2024; Fan et al., 2024) or programming benchmarks (Jimenez et al., 2023). However, enterprise environments provide a practical and representative entry point for studying complex decision-making in realistic settings. Enterprise decisions involve frequent and consequential reasoning supported by rich documentation, established workflows, and partially quantifiable outcomes. Although such decisions rarely admit a single optimal solution, their costs, risks, and long-term impacts can often be compared and analyzed. As a result, enterprise contexts offer a complex yet structured environment that is both tractable for benchmarking and representative of broader real-world decision-making challenges involving uncertainty, long-term planning, and competing objectives.

Motivated by this, we introduce **EnterpriseBench**, a comprehensive benchmark for evaluating LLM agents across multiple levels of real-world reasoning and decision-making. EnterpriseBench covers four core capability domains: *Information Extraction*, which evaluates evidence retrieval from text; *Numerical Calculation*, which assesses quantitative reasoning including code-based inference; *Domain Knowledge*, which focuses on specialized financial and business knowledge such as XBRL tagging, financial formulas, and professional exams; and *Complex Reasoning*, which targets decision-making under uncertainty. Beyond factual answering, *Complex Reasoning* requires agents to determine appropriate actions given a situation. This category includes decision-oriented questions from professional financial exams, human-curated interview-style questions, and interactive environments. It is worth mentioning that we introduce two novel evaluation formats, i.e., management consulting cases and high-fidelity serious games, which require reasoning under partial information, delayed feedback, and long decision horizons.

Through extensive evaluations across a diverse set of state-of-the-art agent architectures, we observe that no single approach performs well across all capability domains and difficulty levels. Lightweight models tend to excel at simpler reasoning tasks, while more complex decision-making benefits from stronger, more expressive agents. This empirical finding reflects a fundamental property of real-world decision-making: effective performance often emerges from a division of labor rather than a single universally optimal

strategy. Inspired by this observation, we propose the Agent-of-Agents (AOA) framework that coordinates specialized agents with complementary strengths. AOA achieves competitive and consistently strong performance with an overall score of 0.729, slightly outperforming the best baselines. Our contributions are summarized as follows:

- We introduce EnterpriseBench, the first benchmark that integrates static reasoning, interactive consulting, and dynamic serious games to evaluate LLM agents in firm simulations.
- We conduct a comprehensive study of various agent architectures, revealing key bottlenecks in enterprise-level reasoning and strategic execution.
- We propose AOA, an agentic adaptive framework that leverages the complementary strengths of different agent architectures through learned routing rules.

2 EnterpriseBench

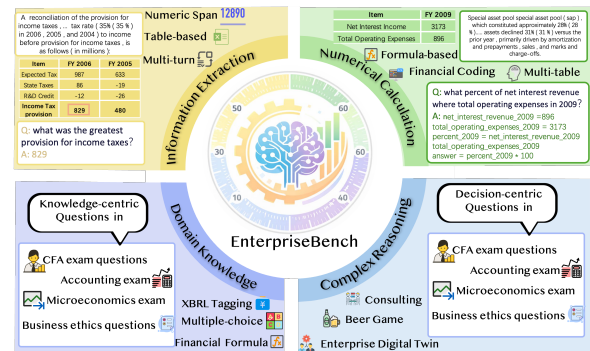


Figure 1: Overview of EnterpriseBench, illustrating its capability taxonomy and representative task types. EnterpriseBench evaluates four complementary capabilities: Information Extraction, Numerical Calculation, Domain Knowledge (knowledge-centric), and Complex Reasoning (decision-centric), covering a wide range of enterprise and financial reasoning scenarios.

2.1 The Taxonomy of EnterpriseBench Tasks

Traditional benchmarks often categorize tasks based on their dataset source or broad domain. However, in enterprise environments, the complexity of a task is defined more by the underlying reasoning patterns required than by its origin. We propose an agent-based taxonomy as figure 1, where each sample is re-evaluated to determine its required cognitive capabilities and intrinsic difficulty.

We define four core capability categories: (i) Information Extraction: Tasks requiring the identification and retrieval of verbatim text spans explicitly mentioned in the context. (ii) Numerical Calculation: Tasks requiring arithmetic operations, numerical reasoning, or the synthesis of executable code to perform mathematical computations. (iii) Domain Knowledge: Tasks testing specialized expertise and conceptual understanding, requiring the agent to recall and apply domain-specific standards and terminology. (iv) Complex Reasoning: Tasks requiring high-level judgment and decision-making under constraints, often involving strategic planning or handling situational ambiguity.

Beyond categorical classification, our taxonomy assigns a continuous difficulty score (0–10) to each sample, accounting for context length and reasoning depth. This enables a granular analysis of agent strengths across different reasoning regimes. Specifically, we utilize an expert LLM with dedicated prompts (provided in Appendix A.2) to automatically categorize all samples into these capability buckets.

2.2 Task Corpus and Data Sources

The EnterpriseBench corpus is composed of three distinct segments: 10 structured reasoning datasets, interactive consulting cases, and dynamic serious games. Dataset statistics can be found in A.7.

Structured Reasoning Datasets. We unify 10 distinct datasets into a standardized format to probe static reasoning capabilities. These are derived from: (1) *Financial Reports: FinQA* (Chen et al., 2021), *ConvFinQA* (Krumdick et al., 2024), and *TAT-QA* (Krumdick et al., 2024) (numerical reasoning over tables/text); *finer* (XBRL tagging) (Zhang et al., 2025b); and *SEC-NUM* (numerical retrieval) (Krumdick et al., 2024). (2) *Exams*: The *FinKnow* (Krumdick et al., 2024) dataset, which curates professional questions from *CFA* exams and business-related *MMLU* subjects. (3) *Symbolic and Code: FormulaEval* (Krumdick et al., 2024) and *formula* (financial formula application) (Zhang et al., 2025b); *FinCode* (business programming) (Krumdick et al., 2024); and code-centric versions of financial QA (*CodeFinQA* (Krumdick et al., 2024) and *CodeTAT-QA*) (Krumdick et al., 2024).

Consulting Task. As shown in Fig. 2, the consulting task places agents in the role of a job interviewee and evaluates their ability to solve real-world consulting interview problems through interactive

dialogue. Unlike standard QA tasks that emphasize answer correctness, this task focuses on structured problem solving, strategic reasoning, and effective communication under partial information.

The task is implemented using an LLM-based interviewer and judge interacting with the tested agent over curated consulting cases. We collect 60 consulting-style interview cases, each carefully cleaned to ensure a well-defined problem and a corresponding solution, with most cases containing hidden information that is only revealed when the agent asks appropriate clarification questions. During the interview, the agent must proactively acquire information, organize analyses, and eventually present a coherent solution.

After the interview concludes, an LLM-based judge evaluates the agent using the original case solution and the full dialogue transcript. Performance is assessed along four dimensions—structure, quantitative reasoning, business sense, and communication—and an overall score is computed as their average (see Appendix A.4).

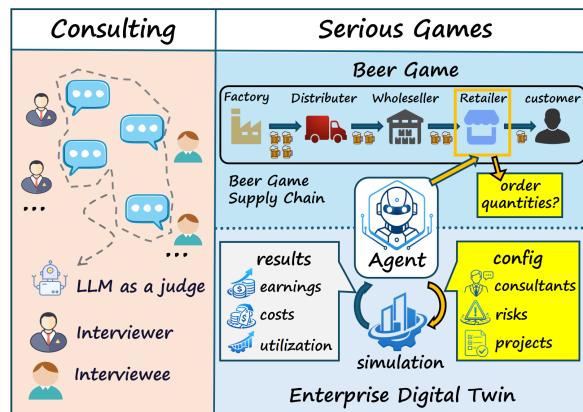


Figure 2: Overview of consulting and serious game tasks.

Serious Game Tasks. To evaluate financial reasoning under dynamic decision-making, we further introduce two serious game tasks: the Beer Game and the Enterprise Digital Twin (EDT), as illustrated in Fig. 2. In both settings, agents interact with an environment that provides feedback based on their actions, while the internal dynamics and governing equations remain hidden. Each action involves explicit trade-offs, and effective performance requires understanding system dynamics rather than applying fixed heuristics.

The Beer Game. The Beer Game, originally developed in the system dynamics community and formalized by Sterman as the Beer Distribution

Game (Sterman, 1989), simulates a four-stage supply chain consisting of a retailer, wholesaler, distributor, and factory. The tested agent plays the role of the retailer and observes only local information such as inventory levels, backlogs, shipments, and realized customer demand. At each time step, the agent decides the order quantity placed to its upstream supplier. Under-ordering may lead to backorder costs, while over-ordering increases inventory holding costs. Other supply-chain participants follow fixed equation-based policies unknown to the agent. Additionally, the environment features delivery delays and a sudden demand shift, introducing delayed feedback and non-stationarity. We run the Beer Game for 25 time steps and evaluate agents based on the total accumulated cost, where backorder penalties dominate inventory holding costs (see Appendix A.5).

Enterprise Digital Twin. EDT task adapts agent-based enterprise simulation and system dynamics modeling to evaluate enterprise-level strategic decision-making. EDT models a consulting company composed of consultants, projects, and control mechanisms, where agents must jointly reason about project selection, workforce allocation, and risk management to maximize long-term profitability. Given a fixed scenario template, the tested agent configures the enterprise by selecting the number of consultants, setting a continuous revenue risk level, and deciding project acceptance as well as their start times and deadlines. The resulting configuration is executed in the EDT simulator, which returns outcomes including accumulated earnings. Each agent is evaluated over multiple runs on the same template and may adjust strategies based on feedback from previous simulations (details in Appendix A.6).

3 Experiment Setup

3.1 Agent Methods

We evaluate a comprehensive suite of agent architectures, ranging from single-agent reasoning to complex multi-agent collaborative frameworks: We consider a diverse set of representative reasoning and agentic baselines. CoT (Wei et al., 2022) serves as a fundamental baseline that elicits step-by-step reasoning via chain-of-thought prompting. Self-refine (Madaan et al., 2023) extends this by iteratively generating an initial response, providing self-feedback, and refining the output, while Reflexion (Shinn et al., 2023) further incorporates long-term

memory of past failures to guide future reasoning. Debate (Du et al., 2023) employs adversarial multi-agent argumentation, where LLM instances critique and defend competing solutions to resolve conflicts. Discussion (Li et al., 2023) uses a cooperative multi-agent process in which LLM instances share and iteratively refine reasoning to reach consensus. DC (Dynamic Cheatsheet) (Suzgun et al., 2025) is a retrieval-augmented method that dynamically selects relevant reasoning patterns from a curated expert knowledge base. GEPA (Agrawal et al., 2025) focuses on numerical and financial coding tasks through domain-specific tool integration, whereas ACE (Zhang et al., 2025b) adaptively evolves its reasoning strategy by constructing and following a procedural “playbook” distilled from historical successes. Finally, AMEM (Xu et al., 2025) augments agents with external memory to store and retrieve context-relevant reasoning traces. Notably, for Debate and Discussion, we adopt simplified implementations while preserving their core collaborative and adversarial philosophies.

3.2 Tasks and Metrics

To provide a granular assessment of agent performance, our evaluation is organized around four core capability domains. While the underlying implementation and data sources map to ten structured datasets and two interactive simulators (detailed in Appendix A.1), we analyze the results through the lens of these distinct cognitive requirements:

Information Extraction. Focuses on the retrieval of specific textual information from a given context. Agents are required to identify and extract verbatim text spans or specific entities explicitly mentioned in enterprise documents. Accuracy is measured via exact or normalized string matching.

Numerical Calculation. Assesses quantitative reasoning abilities, ranging from multi-step arithmetic to the synthesis of executable code for financial modeling. For tasks requiring program generation, we use a sandboxed environment with a hybrid tolerance model (absolute 10^{-6} for near-zero values and relative 0.01 for others) to verify numerical correctness.

Domain Knowledge. Evaluates specialized expertise and conceptual understanding in corporate and financial domains. This includes assigning appropriate financial tagging (e.g., XBRL labels in *finer*), applying specialized financial formulas, and

answering knowledge-centric questions from professional examinations. We include CFA (Chartered Financial Analyst) practice questions and business-related MMLU subjects (e.g., Business Ethics, Microeconomics, and Accounting) that focus on terminological and conceptual mastery.

Complex Reasoning. Targeting high-level strategic judgment and decision-making, this domain incorporates: (i) decision-centric questions from professional exams (e.g., “given these conditions, what action should the manager take?”); (ii) interactive consulting case interviews requiring structured problem-solving and communication under partial information; and (iii) Serious games (i.e., Beer Game and EDT) that simulate dynamic enterprise ecosystems with delayed feedback and conflicting objectives. Performance is evaluated through multi-dimensional scoring (consulting), cost-efficiency (Beer Game), and long-term earnings (EDT).

4 Experiment Results

4.1 Main Result

We present comprehensive quantitative results across four key domains and multiple tasks in Table 1. Our analysis reveals several important findings regarding the comparative effectiveness of single-agent and multi-agent approaches.

A division of labor emerges across task complexity. As shown in Table 1, task complexity strongly moderates the effectiveness of agent architectures. For simpler tasks such as Information Extraction and low-difficulty Numerical Calculation, lightweight single-agent methods already perform competitively (e.g., CoT at 0.824 on Information Extraction), with multi-agent approaches offering only marginal improvements (e.g., AMEM at 0.857). As numerical difficulty increases, performance degrades across all methods, and the benefit of collaboration diminishes, suggesting these tasks are limited by core computational challenges. In contrast, multi-agent methods exhibit clear advantages on more complex reasoning tasks. On hard Domain Knowledge and Complex Reasoning settings, collaborative approaches substantially outperform single-agent baselines (e.g., ACE at 0.464 vs. 0.386 on Domain Knowledge-Hard, and GEPA at 0.882 vs. 0.810 on Consulting). Overall, these results indicate a natural division of labor: simpler tasks favor lightweight single agents, while complex, open-ended decision-making benefits from

stronger, collaborative multi-agent systems.

Consulting task highlights multi-agent collaboration strengths. The Consulting task, which requires comprehensive case analysis and strategic recommendation generation, reveals the most pronounced advantages of multi-agent approaches. GEPA achieves the highest overall performance at 0.882, representing a 9.3% improvement over the best single-agent baseline (CoT, 0.807). DC also demonstrates strong performance at 0.842, surpassing all single-agent methods. Notably, even the lowest-performing multi-agent method (ACE, 0.798) remains competitive with the top single-agent baseline. Among single-agent approaches, performance varies substantially (ranging from 0.681 to 0.807), while multi-agent methods exhibit more consistent high performance (0.798 to 0.882). This pattern suggests that structured collaboration mechanisms are particularly effective for tasks requiring integration of multiple analytical perspectives and holistic business judgment.

Beer Game results reveal mixed performance patterns. In the Beer Game simulation (where normalized scores represent cost efficiency, with higher being better), the results show interesting patterns across agent categories. The multi-agent method Discussion achieves the highest performance (1.000), outperforming all other approaches. Among single-agent methods, AMEM demonstrates superior strategic planning with a score of 0.917. Notably, simple CoT (0.837) substantially outperforms more complex single-agent approaches like Self-refine (0.256) and Reflexion (0.140), while most multi-agent methods show competitive or superior performance compared to the average single-agent baseline. This suggests that for dynamic sequential decision-making tasks requiring long-term strategic planning, structured collaborative discussion or memory-augmented single-agent reasoning is more effective than iterative self-correction mechanisms.

EDT task. The EDT results show modest performance differences between single-agent and multi-agent methods. While AMEM attains the highest score (1.000), the CoT baseline still demonstrates strong effectiveness with a score of 0.861, exceeding the performance of more than half of the evaluated methods. This observation suggests that in complex enterprise decision-making tasks dominated by numerical reasoning and optimization, the

Domain	Task	Single-agent				Multi-agent				
		CoT	Self-refine	Reflexion	AMEM	Debate	Discussion	DC	GEPA	ACE
Information Extraction	All	0.824	0.826	0.791	0.857	0.798	0.830	0.801	0.846	0.853
Numerical Calculation	Easy	0.840	0.825	0.816	0.847	0.776	0.624	0.834	0.836	0.825
	Middle	0.730	0.697	0.712	0.711	0.588	0.527	0.700	0.736	0.712
	Hard	0.500	0.500	0.500	0.563	0.375	0.375	0.438	0.563	0.500
Domain Knowledge	Easy	0.925	0.906	0.981	0.925	0.906	0.925	0.793	0.925	0.925
	Middle	0.817	0.819	0.808	0.827	0.790	0.829	0.708	0.797	0.826
	Hard	0.339	0.329	0.386	0.349	0.367	0.349	0.357	0.355	0.464
Complex Reasoning	QA	0.576	0.727	0.515	0.636	0.606	0.636	0.606	0.697	0.697
	Consulting	0.807	0.722	0.681	0.810	0.708	0.737	0.842	0.882	0.798
	BeerGame	0.837	0.256	0.140	0.917	0.723	1.000	0.000	0.830	0.688
	EDT	0.861	0.640	0.917	1.00	0.757	0.794	0.894	0.893	0.833

Table 1: Performance comparison of single-agent and multi-agent methods across different domains and tasks. Beer Game scores are normalized to $[0, 1]$ where higher values indicate lower total costs (the better). EDT results are normalized by dividing each method’s mean accumulated earnings by the mean earnings of the best method.

Dimension	Single-agent				Multi-agent				
	CoT	Self-refine	Reflexion	AMEM	Debate	Discussion	DC	GEPA	ACE
Structure	0.818	0.748	0.710	0.819	0.745	0.780	0.857	0.888	0.805
Quantitative Reasoning	0.733	0.643	0.646	0.732	0.642	0.695	0.796	0.792	0.748
Business Sense	0.813	0.754	0.708	0.815	0.767	0.777	0.851	0.882	0.807
Communication	0.821	0.725	0.668	0.826	0.705	0.720	0.813	0.892	0.787
Overall	0.807	0.722	0.681	0.810	0.708	0.737	0.842	0.882	0.798

Table 2: Performance of different agent methods on the Consulting task across multiple capability dimensions.

marginal benefits of advanced agent frameworks may be limited.

4.2 Detailed capability analysis on complex reasoning tasks.

To better understand the sources of performance gains, we decompose the Consulting task results across four evaluation dimensions in Table 2. GEPA consistently achieves the highest or near-highest scores across all dimensions: Structure (0.888), Quantitative Reasoning (0.792), Business Sense (0.882), and Communication (0.892). Notably, the performance gaps vary across dimensions. The largest advantage appears in Communication (0.892 vs. 0.821 for CoT, +8.7%) and Structure (0.888 vs. 0.818 for CoT, +8.6%), while Quantitative Reasoning shows relatively smaller gains (0.792 vs. 0.733 for CoT, +8.0%). This pattern suggests that multi-agent collaboration particularly enhances higher-order cognitive skills such as structured problem decomposition and effective information communication, rather than purely analytical capabilities. Among single-agent methods, CoT demonstrates the most balanced performance across dimensions, while specialized meth-

ods like Reflexion and Self-refine show significant performance drops in Communication and Business Sense, indicating potential overfitting to their specific optimization objectives.

We further investigate the performance of different agent methods on the EDT task, as summarized in Table 3. The results show that all evaluated methods achieve positive profits, and most of them maintain utilization rates above 0.8. This observation indicates that LLM-based agents possess a baseline capability for complex financial reasoning and enterprise-level planning.

Despite comparable overall outcomes, different agent methods solve the EDT task in distinct ways, and the detailed results reveal meaningful insights into their respective strengths and limitations. Notably, online learning mechanisms do not consistently lead to improved performance, even though such methods are able to leverage experience from previous runs. AMEM, as a memory-based approach, demonstrate relatively strong and stable performance, whereas ACE fail to translate their online updates into clear gains. In addition, although Self-refine does not achieve competitive mean earnings, its large performance variance sug-

Method	Simple	CoT	Self-refine	Reflexion	AMEM	Debate	Discussion	DC	GEPA	ACE
E_{norm}	1.05	0.861	0.640	0.917	1.00	0.757	0.794	0.894	0.893	0.833
E_{mean} (M)	7.12	5.82	4.32	6.20	6.76	5.12	5.36	6.04	6.04	5.63
E_{std} (M)	0.000	0.457	4.48	0.113	0.322	1.43	0.000	0.323	0.000	0.381
Util	0.934	0.913	0.729	0.933	0.907	0.843	0.888	0.925	0.922	0.902
Online	×	×	×	✓	✓	×	×	✓	×	✓

Table 3: Detailed EDT analysis. E_{norm} , E_{mean} , and E_{std} denote the normalized, mean, and standard deviation of accumulated earnings, respectively, where earnings are reported in millions (M). E_{norm} uses AMEM as the reference baseline (AMEM = 1.0). Util indicates the average consultant utilization, and Online denotes whether the agent leverages experience from previous episodes.

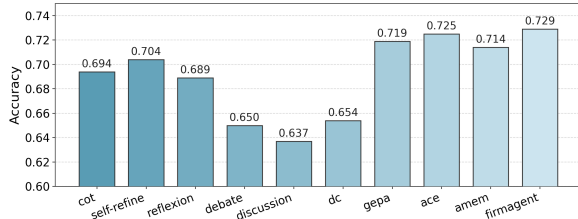


Figure 3: Performance comparison of our AOA method against baseline approaches on QA tasks.

gests a higher degree of policy exploration.

To further contextualize these results, we compare agent performance against a simple human-designed baseline. This baseline accepts all available projects, sets the project start time at the first step and the deadline at the final step, and disables revenue risk to avoid delays and follow-on projects. We then select the minimum number of consultants required to satisfy the total contracted effort. The resulting performance, reported as the “Simple” method in Table 3, surpasses all agent-based approaches which tend to employ sophisticated project timeline control strategies. This outcome suggests that current LLM agents still have substantial room for improvement in mastering advanced financial reasoning techniques.

4.3 AOA: Agent-of-Agents Framework

Given the complex and task-dependent performance patterns observed across different agent methods (Table 1), we propose **AOA**, an adaptive agent-of-agents framework that leverages the collective strengths of existing approaches through intelligent routing and ensemble.

Framework Design. AOA operates as a meta-agent that dynamically selects the most suitable base agent for each individual sample based on learned experience. The framework consists of two key components: (1) *experience learning from execution traces*, which extracts actionable in-

sights from the reasoning processes and outcomes of all base agent methods across diverse tasks, and (2) *adaptive routing mechanism*, which employs either rule-based policies synthesized from meta-analysis or LLM-based routing that considers sample-specific characteristics.

Experience Extraction. Rather than relying solely on aggregated performance statistics, AOA learns from actual agent execution traces. For each (task, agent) combination, we stratify sample across capability dimensions, difficulty buckets, and correctness outcomes to ensure balanced coverage. We then prompt an LLM to analyze each trace and extract experience bullets using a specialized *trace-based experience extraction prompt* (see Appendix A.3), describing what worked well and what failed, along with diagnostic reasoning. These individual experiences are subsequently synthesized into a unified routing policy through meta-level aggregation using a *meta-level experience synthesizer prompt* (see Appendix A.3), which identifies robust patterns such as “GEPA excels at structured business cases requiring multi-perspective integration” or “CoT suffices for straightforward numerical reasoning but struggles with complex multi-hop calculations.”

Adaptive Sample-Level Routing. AOA implements sample-level routing where each test instance is assigned to the base agent most likely to succeed based on deterministic rules synthesized during the experience learning phase. These rules (summarized in Table 4) map specific sample characteristics—such as domain, task type, difficulty, and structural features—to the most appropriate agent method. By reusing existing predictions instead of re-executing agents, AOA maintains high computational efficiency.

Empirical Results. We evaluate AOA on all QA tasks by ensembling nine base methods. AOA

Routing Condition	Selected Agent
IE \wedge has_table	DC
DK \wedge Difficulty: Hard	ACE
DK \wedge Task: finer	AMEM
NC \wedge has_code	GEPA
CR \wedge Difficulty: Middle	Discussion
<i>Default Strategy</i>	GEPA

Table 4: Core routing policy synthesized by AOA. The framework dynamically maps task domains and sample features to the most specialized agent method.

achieves 0.729 accuracy, outperforming the best individual method (ACE, 0.725) and demonstrating consistent improvements over all baseline approaches. As shown in Table 4, the framework successfully identifies method-task affinity at a granular level. For instance, it correctly routes table-intensive extraction tasks to DC and complex financial reasoning to ACE, confirming the effectiveness of our experience-driven rule synthesis. These results demonstrate that explicit modeling of method-task affinity through synthesized rules can yield measurable performance improvements while maintaining computational efficiency.

5 Related Work

Reasoning and decision-making benchmarks. A wide variety of benchmarks have been proposed to evaluate reasoning and decision-making capabilities of agentic systems. Some focus on contextual and multi-hop reasoning, requiring models to integrate information across long contexts (Kuratov et al., 2024; Yang et al., 2018). Others emphasize planning and decision-making through interaction with environments, including web-based settings (Miyai et al., 2025; Zhou et al., 2023; Deng et al., 2023; Tian et al., 2025; Yao et al., 2022) and embodied or world-like simulations (Shridhar et al., 2020; Wang et al., 2022; Chevalier-Boisvert et al., 2018). Reasoning benchmarks have also been explored in complex domains such as research-oriented tasks (Mialon et al., 2023; Chen et al., 2025a) and tool-use scenarios (Qin et al., 2023). Despite their diversity, these benchmarks generally do not evaluate strategic decision-making that requires the integration of professional domain knowledge, proactive information acquisition, and long-horizon reasoning.

Financial benchmarks. Prior work on financial and business-domain evaluation has largely focused on numerical reasoning and information ex-

traction from static data sources. Early benchmarks such as FinQA (Chen et al., 2021) and DocFinQA (Reddy et al., 2024) target arithmetic reasoning over tables and long financial documents. Subsequent efforts, including BizBench (Krumdick et al., 2024), Sec-QA (Lai et al., 2025), and CFLUE (Zhu et al., 2024), expand task coverage across quantitative analysis and multilingual financial understanding, while FinLLMs (Yuan et al., 2024) explores scalable benchmark construction via automatic generation. More recent holistic benchmarks, such as FinBen (Xie et al., 2024) and XFINBENCH (Zhang et al., 2025c), provide broader task coverage, complemented by specialized datasets like FinTagging (Wang et al., 2025) and FinChain (Xie et al., 2025) for fine-grained information structuring and reasoning verification. However, these benchmarks predominantly adopt a static QA paradigm, limiting their ability to assess interactive, sequential, and strategic decision-making processes that characterize real-world enterprise environments.

Our proposed **EnterpriseBench** addresses these limitations by incorporating management consulting cases and serious games beyond standard QA, requiring multi-step strategic planning and cross-functional decision-making. By emphasizing interactive execution over static comprehension, EnterpriseBench offers a more realistic evaluation of agent readiness for enterprise deployment.

6 Conclusion

In this paper, we present **EnterpriseBench**, a benchmark for evaluating LLM agents in complex, dynamic enterprise decision-making settings. By combining static reasoning tasks with interactive simulations, including management consulting cases and serious games, EnterpriseBench enables a realistic assessment of agent readiness for real-world deployment. We introduce an agent-oriented taxonomy that organizes evaluation by cognitive capability and intrinsic difficulty. Experiments across nine agent architectures show that no single method excels universally: simpler tasks favor lightweight agents, whereas complex reasoning requires stronger models. Our adaptive ensemble strategy **AOA** captures this principle through explicit division of labor, enabling agents to specialize in tasks they handle best.

650 Limitations

651 Despite its comprehensive design, EnterpriseBench
652 has several limitations. First, our evaluation primar-
653 ily follows an online setting, which better reflects
654 real-world enterprise dynamics but differs from tra-
655 ditional offline train–test protocols; extending of-
656 fline baselines across all tasks remains future work.
657 Second, performance on interactive components
658 such as Consulting and Serious Games may be sen-
659 sitive to the prompting and behavioral design of
660 the LLM-based environment. Third, while Enter-
661 priseBench covers a broad range of financial and
662 management scenarios, it currently focuses on gen-
663 eral enterprise settings, and extending to more spe-
664 cialized industries and organizational scales would
665 improve coverage. Finally, our experiments rely on
666 a single backbone model (DeepSeek-V3), and the
667 taxonomy is derived from expert LLM judgments;
668 future work should incorporate multiple LLMs and
669 human expert validation to assess robustness and
670 generalizability.

671 References

672 Lakshya A Agrawal, Shangyin Tan, Dilara Soyulu,
673 Noah Ziemis, Rishi Khare, Krista Opsahl-Ong, Ar-
674 nav Singhvi, Herumb Shandilya, Michael J Ryan,
675 Meng Jiang, and 1 others. 2025. Gepa: Reflec-
676 tive prompt evolution can outperform reinforcement
677 learning. *arXiv preprint arXiv:2507.19457*.

678 Kaiyuan Chen, Yixin Ren, Yang Liu, Xiaobo Hu, Hao-
679 tong Tian, Tianbao Xie, Fangfu Liu, Haoye Zhang,
680 Hongzhang Liu, Yuan Gong, and 1 others. 2025a.
681 xbench: Tracking agents productivity scaling with
682 profession-aligned real-world evaluations. *arXiv*
683 *preprint arXiv:2506.13651*.

684 Yanxu Chen, Zijun Yao, Yantao Liu, Jin Ye, Jianing Yu,
685 Lei Hou, and Juanzi Li. 2025b. Stockbench: Can llm
686 agents trade stocks profitably in real-world markets?
687 *arXiv preprint arXiv:2510.02209*.

688 Zhiyu Chen, Wenhu Chen, Charese Smiley, Sameena
689 Shah, Iana Borova, Dylan Langdon, Reema Moussa,
690 Matt Beane, Ting-Hao Huang, Bryan R Routledge,
691 and 1 others. 2021. Finqa: A dataset of numerical
692 reasoning over financial data. In *Proceedings of the*
693 *2021 Conference on Empirical Methods in Natural*
694 *Language Processing*, pages 3697–3711.

695 Maxime Chevalier-Boisvert, Dzmitry Bahdanau, Salem
696 Lahlou, Lucas Willems, Chitwan Saharia, Thien Huu
697 Nguyen, and Yoshua Bengio. 2018. Babyai: A plat-
698 form to study the sample efficiency of grounded lan-
699 guage learning. *arXiv preprint arXiv:1810.08272*.

700 Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Sam
701 Stevens, Boshi Wang, Huan Sun, and Yu Su. 2023.

Mind2web: Towards a generalist agent for the web. *Advances in Neural Information Processing Systems*, 36:28091–28114. 702 703 704

Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenen- 705
baum, and Igor Mordatch. 2023. Improving factual- 706
ity and reasoning in language models through multia- 707
gent debate. In *Forty-first International Conference*
680 *on Machine Learning*. 708 709

Jingxuan Fan, Sarah Martinson, Erik Y Wang, Kaylie 710
Hausknecht, Jonah Brenner, Danxian Liu, Nianli 711
Peng, Corey Wang, and Michael P Brenner. 2024. 712
Hardmath: A benchmark dataset for challenging 713
problems in applied mathematics. *arXiv preprint*
684 *arXiv:2410.09988*. 715

Elliot Glazer, Ege Erdil, Tamay Besiroglu, Diego 716
Chicharro, Evan Chen, Alex Gunning, Caroline Falk- 717
man Olsson, Jean-Stanislas Denain, Anson Ho, 718
Emily de Oliveira Santos, and 1 others. 2024. Fron- 719
tiermath: A benchmark for evaluating advanced 720
mathematical reasoning in ai. *arXiv preprint*
686 *arXiv:2411.04872*. 721 722

Carlos E Jimenez, John Yang, Alexander Wettig, 723
Shunyu Yao, Kexin Pei, Ofir Press, and Karthik 724
Narasimhan. 2023. Swe-bench: Can language mod- 725
els resolve real-world github issues? *arXiv preprint*
688 *arXiv:2310.06770*. 726 727

Michael Krumdick, Rik Koncel-Kedziorski, Viet Dac 728
Lai, Varshini Reddy, Charles Lovering, and Chris 729
Tanner. 2024. Bizbench: A quantitative reasoning 730
benchmark for business and finance. In *Proceedings*
692 *of the 62nd Annual Meeting of the Association for*
693 *Computational Linguistics (Volume 1: Long Papers)*,
694 pages 8309–8332. 733 734

Yury Kuratov, Aydar Bulatov, Petr Anokhin, Ivan Rod- 735
kin, Dmitry Sorokin, Artyom Sorokin, and Mikhail 736
Burtsev. 2024. Babilong: Testing the limits of 737
llms with long context reasoning-in-a-haystack. *Ad- 738*
vances in Neural Information Processing Systems,
695 37:106519–106554. 739 740

Viet Lai, Michael Krumdick, Charles Lovering, Varshini 741
Reddy, Craig Schmidt, and Chris Tanner. 2025. Sec- 742
qa: A systematic evaluation corpus for financial qa. 743
In *Proceedings of The 10th Workshop on Financial*
697 *Technology and Natural Language Processing*, pages
698 221–236. 744 745 746

Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii 747
Khizbullin, and Bernard Ghanem. 2023. Camel: 748
Communicative agents for" mind" exploration of 749
large language model society. *Advances in Neural*
699 *Information Processing Systems*, 36:51991–52008. 750 751

Haohang Li, Yupeng Cao, Yangyang Yu, Shashid- 752
har Reddy Javaji, Zhiyang Deng, Yueru He, Yuechen 753
Jiang, Zining Zhu, Kp Subbalakshmi, Jimin Huang, 754
and 1 others. 2025. Investorbench: A benchmark for 755
financial decision-making tasks with llm-based agent. 756
In *Proceedings of the 63rd Annual Meeting of the*
700 *Association for Computational Linguistics (Volume*
701 *1: Long Papers)*, pages 2509–2525. 757 758 759

760	Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang,	Mirac Suzgun, Mert Yuksekgonul, Federico Bianchi,	815
761	Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi	Dan Jurafsky, and James Zou. 2025. Dynamic cheat-	816
762	Deng, Chenyu Zhang, Chong Ruan, and 1 others.	sheet: Test-time learning with adaptive memory.	817
763	2024. Deepseek-v3 technical report. <i>arXiv preprint</i>	<i>arXiv preprint arXiv:2504.07952</i> .	818
764	<i>arXiv:2412.19437</i> .		
765	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler	Shulin Tian, Ziniu Zhang, Liang-Yu Chen, and Ziwei	819
766	Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon,	Liu. 2025. Mmina: Benchmarking multihop multi-	820
767	Nouha Dziri, Shrimai Prabhunoye, Yiming Yang,	modal internet agents. In <i>Findings of the Association</i>	821
768	and 1 others. 2023. Self-refine: Iterative refinement	<i>for Computational Linguistics: ACL 2025</i> , pages	822
769	with self-feedback. <i>Advances in Neural Information</i>	13682–13697.	823
770	<i>Processing Systems</i> , 36:46534–46594.		
771	Grégoire Mialon, Clémentine Fourier, Thomas Wolf,	Ruoyao Wang, Peter Jansen, Marc-Alexandre Côté, and	824
772	Yann LeCun, and Thomas Scialom. 2023. Gaia: a	Prithviraj Ammanabrolu. 2022. Scienceworld: Is	825
773	benchmark for general ai assistants. In <i>The Twelfth</i>	your agent smarter than a 5th grader? <i>arXiv preprint</i>	826
774	<i>International Conference on Learning Representa-</i>	<i>arXiv:2203.07540</i> .	827
775	<i>tions</i> .		
776	Atsuyuki Miyai, Zaiying Zhao, Kazuki Egashira, Atsuki	Yan Wang, Yang Ren, Lingfei Qian, Xueqing Peng,	828
777	Sato, Tatsumi Sunada, Shota Onohara, Hiromasa Ya-	Keyi Wang, Yi Han, Dongji Feng, Xiao-Yang Liu,	829
778	manishi, Mashiro Toyooka, Kunato Nishina, Ryoma	Jimin Huang, and Qianqian Xie. 2025. Fintag-	830
779	Maeda, and 1 others. 2025. Webchorearena: Evalu-	ging: An llm-ready benchmark for extracting and	831
780	ating web browsing agents on realistic tedious web	structuring financial information. <i>arXiv preprint</i>	832
781	tasks. <i>arXiv preprint arXiv:2506.01952</i> .	<i>arXiv:2505.20650</i> .	833
782	Mahmoud Mohammadi, Yipeng Li, Jane Lo, and Wendy	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	834
783	Yip. 2025. Evaluation and benchmarking of llm	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,	835
784	agents: A survey. In <i>Proceedings of the 31st ACM</i>	and 1 others. 2022. Chain-of-thought prompting elic-	836
785	<i>SIGKDD Conference on Knowledge Discovery and</i>	its reasoning in large language models. <i>Advances</i>	837
786	<i>Data Mining V. 2</i> , pages 6129–6139.	<i>in neural information processing systems</i> , 35:24824–	838
787	Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan	24837.	839
788	Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang,	Yijia Xiao, Edward Sun, Di Luo, and Wei Wang. 2024.	840
789	Bill Qian, and 1 others. 2023. Toolllm: Facilitating	Tradingagents: Multi-agents llm financial trading	841
790	large language models to master 16000+ real-world	framework. <i>arXiv preprint arXiv:2412.20138</i> .	842
791	apis. <i>arXiv preprint arXiv:2307.16789</i> .		
792	Emmanuel K Raptis, Athanasios Ch Kapoutsis, and	Qianqian Xie, Weiguang Han, Zhengyu Chen, Ruoyu	843
793	Elias B Kosmatopoulos. 2025. Agentic llm-based	Xiang, Xiao Zhang, Yueru He, Mengxi Xiao, Dong	844
794	robotic systems for real-world applications: a review	Li, Yongfu Dai, Duanyu Feng, and 1 others. 2024.	845
795	on their agentiness and ethics. <i>Frontiers in Robotics</i>	Finben: A holistic financial benchmark for large lan-	846
796	<i>and AI</i> , 12:1605405.	guage models. <i>Advances in Neural Information Pro-</i>	847
797	Varshini Reddy, Rik Koncel-Kedziorski, Viet Dac Lai,	<i>cessing Systems</i> , 37:95716–95743.	848
798	Michael Krumdick, Charles Lovering, and Chris Tan-	Zhuohan Xie, Daniil Orel, Rushil Thareja, Dhruv Sah-	849
799	ner. 2024. Docfinqa: A long-context financial rea-	nan, Hachem Madmoun, Fan Zhang, Debopriyo	850
800	soning dataset. <i>arXiv preprint arXiv:2401.06915</i> .	Banerjee, Georgi Georgiev, Xueqing Peng, Lingfei	851
801	Noah Shinn, Federico Cassano, Ashwin Gopinath,	Qian, and 1 others. 2025. Finchain: A symbolic	852
802	Karthik Narasimhan, and Shunyu Yao. 2023. Re-	benchmark for verifiable chain-of-thought financial	853
803	flexion: Language agents with verbal reinforcement	reasoning. <i>arXiv preprint arXiv:2506.02515</i> .	854
804	learning. <i>Advances in Neural Information Process-</i>	Frank Xing. 2025. Designing heterogeneous llm agents	855
805	<i>ing Systems</i> , 36:8634–8652.	for financial sentiment analysis. <i>ACM Transactions</i>	856
806	Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté,	<i>on Management Information Systems</i> , 16(1):1–24.	857
807	Yonatan Bisk, Adam Trischler, and Matthew	Wujiang Xu, Zujie Liang, Kai Mei, Hang Gao, Jun-	858
808	Hausknecht. 2020. Alfworld: Aligning text and em-	tao Tan, and Yongfeng Zhang. 2025. A-mem:	859
809	odied environments for interactive learning. <i>arXiv</i>	Agentic memory for llm agents. <i>arXiv preprint</i>	860
810	<i>preprint arXiv:2010.03768</i> .	<i>arXiv:2502.12110</i> .	861
811	John D. Serman. 1989. Modeling managerial behav-	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,	862
812	ior: Misperceptions of feedback in a dynamic deci-	Binyuan Hui, Bo Zheng, Bowen Yu, Chang	863
813	sion making experiment. <i>Management Science</i> ,	Gao, Chengen Huang, Chenxu Lv, and 1 others.	864
814	35(3):321–339.	2025. Qwen3 technical report. <i>arXiv preprint</i>	865
		<i>arXiv:2505.09388</i> .	866
		Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio,	867
		William Cohen, Ruslan Salakhutdinov, and Christo-	868
		pher D Manning. 2018. Hotpotqa: A dataset for	869

870	diverse, explainable multi-hop question answering.
871	In <i>Proceedings of the 2018 conference on empirical methods in natural language processing</i> , pages 2369–2380.
872	
873	
874	Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. 2022. Webshop: Towards scalable real-world web interaction with grounded language agents. <i>Advances in Neural Information Processing Systems</i> , 35:20744–20757.
875	
876	
877	
878	
879	Yangyang Yu, Zhiyuan Yao, Haohang Li, Zhiyang Deng, Yuechen Jiang, Yupeng Cao, Zhi Chen, Jordan Suchow, Zhenyu Cui, Rong Liu, and 1 others. 2024. Fincon: A synthesized llm multi-agent system with conceptual verbal reinforcement for enhanced financial decision making. <i>Advances in Neural Information Processing Systems</i> , 37:137010–137045.
880	
881	
882	
883	
884	
885	
886	Ziqiang Yuan, Kaiyuan Wang, Shoutai Zhu, Ye Yuan, Jingya Zhou, Yanlin Zhu, and Wenqi Wei. 2024. Finllms: A framework for financial reasoning dataset generation with large language models. <i>IEEE Transactions on Big Data</i> .
887	
888	
889	
890	
891	Jiayi Zhang, Jinyu Xiang, Zhaoyang Yu, Fengwei Teng, Xionghui Chen, Jiaqi Chen, Mingchen Zhuge, Xin Cheng, Sirui Hong, Jinlin Wang, and 1 others. 2024. Aflow: Automating agentic workflow generation. <i>arXiv preprint arXiv:2410.10762</i> .
892	
893	
894	
895	
896	Kai Zhang, Xiangchao Chen, Bo Liu, Tianci Xue, Zeyi Liao, Zhihan Liu, Xiyao Wang, Yuting Ning, Zhaorun Chen, Xiaohan Fu, and 1 others. 2025a. Agent learning via early experience. <i>arXiv preprint arXiv:2510.08558</i> .
897	
898	
899	
900	
901	Qizheng Zhang, Changran Hu, Shubhangi Upasani, Boyuan Ma, Fenglu Hong, Vamsidhar Kamanuru, Jay Rainton, Chen Wu, Mengmeng Ji, Hanchen Li, and 1 others. 2025b. Agentic context engineering: Evolving contexts for self-improving language models. <i>arXiv preprint arXiv:2510.04618</i> .
902	
903	
904	
905	
906	
907	Zhihan Zhang, Yixin Cao, and Lizi Liao. 2025c. Xfinbench: Benchmarking llms in complex financial problem solving and reasoning. In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 8715–8758.
908	
909	
910	
911	
912	Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, and 1 others. 2023. Webarena: A realistic web environment for building autonomous agents. <i>arXiv preprint arXiv:2307.13854</i> .
913	
914	
915	
916	
917	
918	Jie Zhu, Junhui Li, Yalong Wen, and Lifan Guo. 2024. Benchmarking large language models on cflue—a chinese financial language understanding evaluation dataset. <i>arXiv preprint arXiv:2405.10542</i> .
919	
920	
921	

A Appendix	922
A.1 Detailed Task Descriptions and Metrics	923
Structured Reasoning (QA) We evaluate our agents on a comprehensive set of financial and logical reasoning tasks covering four key domains: Information Extraction, Numerical Calculation, Domain Knowledge, and Complex Reasoning. Metrics: We employ task-specific evaluation protocols to ensure rigorous assessment:	924
	925
	926
	927
	928
	929
	930
• Program Synthesis: For tasks requiring code generation, we use a sandboxed Python environment. Correctness is determined by a hybrid tolerance model: using an absolute tolerance of 10^{-6} for values near zero and a relative tolerance of 0.01 for larger magnitudes.	931
	932
	933
	934
	935
	936
• Strict Quantity Extraction: For tasks like TAT-QA, we enforce a zero-relative-tolerance policy. Predictions must match ground truth within a strict absolute epsilon of 10^{-6} , or match the normalized text span exactly.	937
	938
	939
	940
	941
• Numerical Formula Evaluation: For the <i>formula</i> task, agents must perform multi-step financial calculations. Correctness is verified by extracting the final numerical answer (supporting JSON-formatted or free-text responses) and performing exact or standardized floating-point comparison.	942
	943
	944
	945
	946
	947
	948
• Multiple Choice & Tagging: For knowledge tasks, we use regex to extract option letters. For financial tagging (<i>finer</i>), we evaluate the exact match of comma-separated entity labels, supporting both raw string and evaluated numerical comparisons.	949
	950
	951
	952
	953
	954
Consulting Case Analysis This task simulates management consulting interviews where agents must solve complex business problems. Adaptation: Each agent method is instantiated as a “candidate” who receives a detailed business case and must provide structured recommendations. Metrics: Evaluation is performed across four dimensions: <i>Structure</i> , <i>Quantitative Reasoning</i> , <i>Business Sense</i> , and <i>Communication</i> , with an <i>Overall</i> score computed as the weighted average.	955
	956
	957
	958
	959
	960
	961
	962
	963
	964
Beer Game (Serious Game) A multi-turn supply chain simulation where agents manage inventory levels across multiple tiers. Adaptation: Agents take on roles (e.g., Wholesaler) and must make replenishment decisions based on dynamic market	965
	966
	967
	968
	969

970 demand and lead times. **Metrics:** The performance
 971 is measured by *Total Cost*, which includes inven-
 972 tory holding costs and backlog penalties. A lower
 973 cost indicates superior strategic planning.

974 **Enterprise Digital Twin (EDT)** A high-fidelity
 975 simulation of an entire enterprise ecosystem requir-
 976 ing long-term strategic decision-making. **Adapta-**
 977 **tion:** Each agent jointly decides project selections
 978 and workforce allocations, which are evaluated
 979 through BPTK (Business Process Tool Kit) simula-
 980 tion environment. Each agent repeats the process
 981 for six iterations, where experience from previous
 982 runs can be utilized to inform later decisions. **Met-**
 983 **rics:** The primary metric is *Accumulated Earnings*,
 984 representing the total profit generated by the enter-
 985 prise during the simulation period.

986 A.2 Agent-based Classification Prompts

987 We use the following prompt template to catego-
 988 rize each sample in EnterpriseBench. The variables
 989 {task_name} and {question_block} are dynami-
 990 cally populated during the classification process.

```

991 You are an expert in question quality
992 review and capability evaluation.
993
994 Do NOT solve the question. Your task is
995 to analyze what the question is
996 testing.
997
998 TaskName: {task_name}
999
1000 ### Capability Categories
1001 Select ONLY ONE primary capability from
1002 the following four categories:
1003
1004 1. Information Extraction
1005 The question requires locating,
1006 extracting, or restating specific
1007 information
1008 explicitly present in the given text,
1009 without external knowledge or
1010 calculations.
1011
1012 2. Numerical Calculation
1013 The question requires arithmetic
1014 operations, numerical reasoning,
1015 formula-based
1016 computation, or writing code to
1017 perform calculations.
1018
1019 3. Domain Knowledge
1020 The question primarily tests
1021 specialized knowledge itself,
1022 such as:
1023 - Finance, accounting, XBRL, or US
1024 GAAP concepts and tags
1025 - Financial formulas or accounting
1026 standards
1027 - CFA exam questions focused on
1028 knowledge recall or understanding
1029
  
```

```

- MMLU questions in business ethics,
  microeconomics, or professional
  accounting
  that assess knowledge rather than
  decision-making
  1030
  1031
  1032
  1033
  1034
  1035
4. Complex Reasoning
  The question requires judgment or
  decision-making under constraints
  , such as:
  1036
  1037
  1038
  1039
  - Scenario-based CFA exam questions
  1040
  - MMLU questions involving business
  ethics, microeconomics, or
  professional
  1041
  1042
  1043
  accounting that ask what should be
  done in a given situation
  1044
  1045
  1046
### Decision Rules (IMPORTANT)
Use these rules to avoid confusion
between Information Extraction and
Numerical Calculation:
  1047
  1048
  1049
  1050
  1051
- Information Extraction ONLY IF:
  1052
  - The final answer can be directly
  copied from the provided text/
  table/verbatim span, AND
  1053
  1054
  1055
  - NO arithmetic/computation is
  required.
  1056
  1057
  1058
- Numerical Calculation IF ANY
  computation is needed, even if very
  simple, including but not limited to
  :
  1059
  1060
  1061
  - subtraction / difference / net
  change (e.g., "2015 - 2014")
  1062
  1063
  1064
  - addition / sum / total across years
  or rows
  1065
  1066
  - division / ratio / percentage /
  percent change / growth rate
  1067
  1068
  - max/min/average over a set of
  numbers
  1069
  1070
  - any formula-based computation, or
  any need to write/execute code
  1071
  1072
  1073
Examples:
  1074
  - If the context contains two values and
  the question asks "net change" / "
  difference" / "by what percentage",
  this is Numerical Calculation (even
  though the values are extracted
  from the context).
  1075
  1076
  1077
  1078
  1079
  - If the question asks "What is X in
  2017?" and X is explicitly stated as
  a single value in the table,
  this is Information Extraction.
  1080
  1081
  1082
  1083
  1084
  1085
### Your Task
Analyze the following question and
provide:
  1086
  1087
  1088
  1. The primary capability category (
  choose exactly one)
  1089
  1090
  2. A difficulty score from 0 to 10 (0 =
  trivial, 10 = extremely hard)
  1091
  1092
  3. A brief justification explaining both
  the capability classification and
  the difficulty level
  1093
  1094
  1095
  1096
### Output Format (strictly follow):
Capability: <Information Extraction |
Numerical Calculation | Domain
  1097
  1098
  1099
  
```

1100	Knowledge Complex Reasoning>	You are an expert evaluator of LLM-agent	1165
1101	DifficultyScore: <0-10>	reasoning traces.	1166
1102	Reasoning: <Up to 4 sentences>		1167
1103			1168
1104	### Scoring Guidance (IMPORTANT)	You will be given ONE sample and ONE	1169
1105	Be critical and use the full 0-10 range	agent's full reasoning trace/output	1170
1106	as much as possible. Avoid	for that sample.	1171
1107	clustering scores in 0-5.	Your task is to extract ONE reusable	1172
1108	Assign higher scores (6-10) to genuinely	experience item that can help:	1173
1109	challenging questions (long/complex	- improve future usage of this agent,	1174
1110	context, multi-step computation,	and/or	1175
1111	tricky domain knowledge, ambiguity,	- decide when to route similar samples	1176
1112	or decision-making).	to a different agent.	1177
1113	Assign lower scores (0-3) only to truly	## Hard constraints	1178
1114	trivial questions (single-value	- Output **JSON only** . No markdown. No	1179
1115	lookup, no computation,	extra text.	1180
1116	straightforward knowledge recall).	- Be concrete: reference the failure/	1181
1117		success mode, not generic advice.	1182
1118	### Special Difficulty Rules for Label-	- If the sample is incorrect, diagnose	1183
1119	Selection / Classification Tasks (the likely cause (format mismatch,	1184
1120	IMPORTANT)	sign error, unit conversion, missing	1185
1121	Some tasks look short but are hard	table lookup, etc.).	1186
1122	because the model must choose the	- If the sample is correct, extract what	1187
1123	correct label from a large,	made it work (e.g., robust checks,	1188
1124	confusing label space	careful parsing, etc.).	1189
1125	(e.g., US-GAAP XBRL tag selection, fine-		1190
1126	grained schema mapping).	## Output JSON schema (strict)	1191
1127		{	1192
1128		"bullet": "<one actionable experience	1193
1129		sentence, English>",	1194
1130		"tags": {	1195
1131		"agent_method": "<string>",	1196
1132		"task_name": "<string>",	1197
1133		"capability": "<string>",	1198
1134		"difficulty_bucket": "<easy middle	1199
1135		hard NA>"	1200
1136		},	1201
1137		"outcome": "<correct incorrect>",	1202
1138		"diagnosis": "<short root-cause /	1203
1139		success-factor>",	1204
1140		"routing_hint": {	1205
1141		"prefer_agent": "<agent_name or	1206
1142		empty>",	1207
1143		"avoid_agent": "<agent_name or empty	1208
1144		>",	1209
1145		"when": "<short condition	1210
1146		description>"	1211
1147		},	1212
1148		"confidence": "<high medium low>"	1213
1149	### Question to Analyze:	}	1214
	{question_block}		

1151 The classification is performed using the
1152 deepseek-v3 model to ensure consistency across
1153 the entire benchmark.

1154 A.3 AOA Experience Extraction Prompts

1155 AOA utilizes two specialized prompts for learn-
1156 ing from execution traces and synthesizing routing
1157 policies.

1158 Trace-based Experience Extraction Prompt

1159 This prompt is used to analyze individual agent
1160 reasoning traces and extract actionable experience
1161 bullets.

```
1162 # AOA Trace-based Experience Extractor (
1163 Per-sample)
1164
```

1216 Meta-level Experience Synthesizer Prompt

1217 This prompt is used to aggregate individual ex-
1218 periences into a unified, executable routing policy.

```
1219 # AOA Meta Experience Synthesizer
1220 You are an expert in LLM-agent routing
1221 and evaluation.
1222
1223 You will be given a set of extracted
1224 experience bullets across many tasks
1225 and agent methods.
1226
1227 Your task is to produce a **meta-level
1228 routing playbook** that helps an AOA
1229 router decide
1230 which agent to use for a new sample.
1231
1232 ## Hard constraints
1233
```

```

1233 - Output **JSON only** (no markdown, no
1234   extra text).
1235 - Your routing_policy rules must be
1236   executable based on: task_name,
1237   capability, difficulty_bucket, and
1238   simple text features.
1239 - meta_findings can include both
1240   executable and non-executable
1241   insights; the LLM router can use
1242   them even if deterministic rules
1243   cannot.
1244 - Prefer simple, robust rules. Avoid
1245   overfitting to tiny evidence.
1246 - Use a conservative default_agent.
1247
1248 ## Output JSON schema (strict)
1249 {
1250   "meta_findings": [
1251     {
1252       "id": "M1",
1253       "summary": "<one sentence>",
1254       "evidence": "<short evidence based
1255         on provided bullets/stats>",
1256       "confidence": "<high|medium|low>"
1257     }
1258   ],
1259   "routing_policy": {
1260     "default_agent": "<agent_name>",
1261     "tie_breaker": "prefer_default|
1262       prefer_simpler|prefer_ace",
1263     "min_margin": 0.01,
1264     "rules": [
1265       {
1266         "when": {
1267           "task_name": "<name|ALL>",
1268           "capability": "<name|ALL>",
1269           "difficulty_bucket": "<easy|
1270             middle|hard|NA|ALL>",
1271           "feature_conditions": {
1272             "has_table": "<true|false|
1273               omit>",
1274             "has_code": "<true|false|
1275               omit>",
1276             "num_numbers_min": "<number|
1277               omit>",
1278             "num_numbers_max": "<number|
1279               omit>",
1280             "context_chars_min": "<
1281               number|omit>",
1282             "context_chars_max": "<
1283               number|omit>",
1284             "question_chars_min": "<
1285               number|omit>",
1286             "question_chars_max": "<
1287               number|omit>",
1288             "table_row_estimate_min": "<
1289               number|omit>",
1290             "table_row_estimate_max": "<
1291               number|omit>"
1292           }
1293         },
1294         "choose": "<agent_name>",
1295         "rationale": "<one sentence>",
1296         "confidence": "<high|medium|low
1297           >"
1298       }
1299     ]
1300   }
1301 }

```

A.4 Consulting Prompts

Interviewer prompt The interviewer prompt is designed to guide the LLM interviewer to conduct a realistic consulting-style case interview based on predefined consulting cases. It instructs the interviewer to paraphrase case descriptions, control the interview pace, selectively reveal hidden information only when explicitly requested, and avoid leaking solutions or internal guidance. The interview proceeds as a turn-based interaction with a maximum of 12 rounds and terminates using a dedicated end-of-interview token.

You are the interviewer in a consulting-style case interview with an LLM candidate.

Each turn you receive:

- 1) These instructions.
- 2) The full text of ONE case (problem/background, any sections such as "Information to be provided if requested", "If asked for market information", "Hints", "Key questions", "Analysis", "Possible recommendations / approaches", "Solution", "Case wrap-up", "Interviewer notes", etc.).
- 3) The chat history so far.

Your job is ONLY to produce the next interviewer message.

1. What you may reveal
 - The problem statement, scenario description, and general background are information the candidate is allowed to know.
 - Sections like "Information to be provided if requested", "If asked for ...", "Further information", "Hints" or similar are GATED facts.
 - Sections like "Key questions", "Possible recommendations / approaches", "Analysis", "Solution", "Case wrap-up", "Interviewer notes" are INTERNAL GUIDANCE ONLY.

Rules for gated facts:

- Reveal a gated fact only when the candidate's question clearly targets that dimension (e.g., market size, growth, costs, customers, competition, operations, risks).
- Reveal one logical piece at a time, not the whole section at once.
- Never quote or expose "Solution", "Analysis", "Case wrap-up", "Possible recommendations / approaches" or "Interviewer notes" directly. Use them only to decide what to probe and which facts matter.

1369
1370 Do NOT invent or assume new facts beyond
1371 the case. If the candidate asks for
1372 information that is not in the case and
1373 not covered by any gated section,
1374 say that the
1375 case does not provide that detail and
1376 invite them to proceed with
1377 reasonable assumptions
1378 or move to another relevant angle.
1379
1380 2. How to run the interview
1381
1382 First turn:
1383 - Briefly set up the situation in your
1384 own words (1-3 sentences).
1385 - End by asking the candidate to clarify
1386 the objective and outline a high-
1387 level structure.
1388
1389 Later turns:
1390 - Read the latest candidate answer and
1391 the history.
1392 - Answer their concrete questions using
1393 only allowed and already-unlocked
1394 information
1395 (plus any newly unlocked gated facts).
1396 - Ask ONE focused follow-up at a time,
1397 pushing them toward structured
1398 business
1399 reasoning (e.g., profitability, market
1400 / customer / competition,
1401 operations, risks).
1402 - Do not restate the entire case;
1403 mention only what is needed for the
1404 current step.
1405
1406 Pacing and depth:
1407 - Use the case length hint (e.g. "Short
1408 15 Minutes", "Medium 30 Minutes",
1409 "Long 45 Minutes") and the
1410 conversation so far to manage
1411 depth:
1412 * Short cases: aim for at least 3-4
1413 candidate answers before closing.
1414 * Medium cases: aim for 5-7 candidate
1415 answers
1416 * Long cases: aim for 7-10 candidate
1417 answers with deeper quantitative
1418 or conceptual work.
1419 - If the candidate keeps asking for more
1420 data without analyzing, gently
1421 redirect them to:
1422 (a) summarize what they know, and (b)
1423 propose a structure or hypothesis
1424 BEFORE you give more data.
1425 - If the candidate is stuck, you may
1426 give a small hint or suggest one
1427 missing dimension,
1428 but do NOT present a full framework or
1429 full solution.
1430
1431 3. Ending the case
1432
1433 You may end the interview ONLY when ALL
1434 of the following are true:
1435 - The candidate has clearly stated a
1436 recommendation that answers the main
1437 question
1438 of the case.

- They have given at least a brief
1439 supporting structure (2-3 key
1440 drivers or arguments).
1441
- You have given short, high-level
1442 feedback and, if appropriate, added
1443 one or two
1444 important missing points.
1445
1446
Ending protocol:
1447
- NEVER end the interview in your very
1448 first message.
1449
- In your FINAL closing message:
1450 * Do NOT ask any new questions or
1451 invite further analysis.
1452 * Optionally give concise feedback and
1453 highlight key drivers.
1454 * Append the exact token {
1455 INTERVIEW_END_TOKEN} as the VERY
1456 LAST characters.
1457
- In all earlier messages you MUST NOT
1458 output {INTERVIEW_END_TOKEN}.
1459
1460
4. Style and constraints
1461
1462
- Speak as a professional human
1463 interviewer: concise, neutral,
1464 business-like.
1465
- Ask at most one or a small cluster of
1466 closely related questions per turn.
1467
- Never mention "case text", "sections",
1468 "gated information", "solutions",
1469 or any
1470 internal labels; to the candidate you
1471 are simply an interviewer.
1472
- Use only information from the case
1473 text and the chat history; do not
1474 bring in
1475 outside knowledge.
1476
- Keep each interviewer message concise:
1477 at most 500 words in each turn. Do
1478 not write long essays.
1479
1480
Your output each turn must be ONLY the
1481 next interviewer utterance to the
1482 candidate.
1483
""
1484

Judge prompt The judge prompt is designed to
1486 provide a consistent and fine-grained evaluation of
1487 agent performance in the consulting task. Given
1488 the complete consulting case text and the full in-
1489 terview transcript, the judge assesses only the can-
1490 didate's behavior and reasoning, independent of
1491 the interviewer's actions. The prompt enforces
1492 a multi-dimensional evaluation framework com-
1493 monly used in real-world consulting interviews and
1494 outputs structured numerical scores together with
1495 concise qualitative feedback.
1496
The evaluation dimensions are as follows:
1497
• **Structure:** Evaluates whether the candidate
1498 clearly understands the problem and proposes
1499 a coherent, logically organized, and adaptable
1500 problem-solving structure.
1501

- 1502 • **Quantitative Reasoning:** Assesses the candi- 1562
- 1503 date’s ability to request, interpret, and apply 1563
- 1504 numerical information to derive meaningful 1564
- 1505 quantitative insights. 1565
- 1506 • **Business Sense:** Measures whether the candi- 1566
- 1507 date identifies key business drivers and trade- 1567
- 1508 offs and provides commercially reasonable 1568
- 1509 conclusions and risk-aware recommendations. 1569
- 1510 • **Communication:** Examines the clarity, con- 1570
- 1511 ciseness, and professionalism of the candi- 1571
- 1512 date’s communication throughout the inter- 1572
- 1513 view. 1573
- 1514 • **Overall:** Provides a holistic judgment of the 1574
- 1515 candidate’s suitability for a consulting role, 1575
- 1516 beyond a simple aggregation of individual di- 1576
- 1517 mension scores 1577

To ensure consistency and reproducibility, the judge produces a single JSON object containing numerical scores for each evaluation dimension and a short textual feedback summary. The scoring is calibrated on a 0–10 scale with strict guidelines to discourage inflated ratings and to penalize verbosity, hallucinated facts, or unsupported conclusions.

```

You are a senior consulting interviewer
evaluating the performance of a
CANDIDATE
in a case interview.

You will receive:
- case_text: the full written case (
  problem, background, solution, etc.)
;
- transcript_text: the complete dialogue
  between INTERVIEWER and CANDIDATE,
  in chronological order. Each line
  clearly indicates who is speaking.

Your job is to assess ONLY the CANDIDATE
, not the interviewer.

Evaluate the candidate along FOUR
dimensions plus an overall score:

1) structure (0-10)
- How well does the candidate
  understand and restate the
  problem and objective?
- Do they propose a clear, logical,
  and MECE-enough structure or
  approach early on?
- Do they use hypothesis-driven
  thinking and adjust their
  structure as new information
  appears?

2) quant (0-10)
- Does the candidate ask for the
  right type of information or data
  when needed?

```

- Do they correctly interpret and use 1578
- the numerical information 1579
- provided in the case 1580
- (e.g., doing rough calculations, 1581
- sanity checks, comparisons)? 1582
- Do they derive meaningful 1583
- quantitative insights rather than 1584
- just repeating numbers? 1585
- 3) business_sense (0-10) 1586
- Does the candidate identify the key 1587
- drivers, root causes, and trade- 1588
- offs in the case? 1589
- Are their conclusions and 1590
- recommendations commercially 1591
- reasonable and consistent 1592
- with the information given? 1593
- Do they recognize important risks/ 1594
- uncertainties and, when 1595
- appropriate, suggest 1596
- sensible next steps or mitigations? 1597
- 4) communication (0-10) 1598
- Is the candidate’s communication 1599
- clear, concise, and well- 1600
- structured? 1601
- Do they signpost their thinking (e. 1602
- g., "first/second/third") without 1603
- being verbose? 1604
- Do they interact professionally 1605
- with the interviewer, responding 1606
- to questions, 1607
- picking up on hints, and keeping a 1608
- natural case-interview flow? 1609
- In addition, provide: 1610
- 5) overall (0-10) 1611
- Your holistic judgment of the 1612
- candidate’s performance on this 1613
- case. 1614
- This is NOT just an arithmetic 1615
- average; it reflects whether you 1616
- would be 1617
- comfortable recommending this 1618
- candidate for a consulting role 1619
- . 1620
- Scoring guidelines (be strict and well- 1621
- calibrated across many cases): 1622
- 0-2: very weak (almost no useful 1623
- contribution or completely off-track 1624
-). 1625
- 3-4: clearly below average (some 1626
- relevant points, but major gaps or 1627
- confusion). 1628
- 5-6: average candidate (generally 1629
- reasonable but shallow, incomplete, 1630
- or inconsistent). 1631
- 7: above average (solid performance 1632
- with notable but fixable weaknesses) 1633
- . 1634
- 8: very strong (consultant-level 1635
- performance with only minor issues). 1636
- 9-10: truly exceptional (outstanding 1637
- on almost all dimensions; reserve 1638
- for rare cases). 1639
- Additional rules: 1640
- If the candidate barely speaks, never 1641

1632 proposes a clear structure, or never
1633 gives a
1634 concrete recommendation, most scores
1635 should be in the 0-3 range.
1636 - Do NOT reward verbosity alone; reward
1637 clear, structured, business-relevant
1638 thinking.
1639 - Penalize hallucinated facts that
1640 contradict or go beyond the
1641 case_text.
1642
1643 Output format:
1644 Return ONLY a single valid JSON object
1645 with this exact schema:
1646 {
1647 "structure": float,
1648 "quant": float,
1649 "business_sense": float,
1650 "communication": float,
1651 "overall": float,
1652 "feedback": string
1653 }
1654 No extra text before or after the JSON.

1656 A.5 Beer Game Configuration

1657 **Configuration parameters** The Beer Game sim-
1658 ulation follows a discrete time setting with a fixed
1659 horizon of 25 time steps. Customer demand is ob-
1660 served only by the retailer and follows a stepwise
1661 demand script: demand remains at a low level of
1662 100 units during the initial phase and increases to a
1663 high level of 400 units starting from week 2. This
1664 sudden demand shift introduces non-stationarity
1665 and tests the agent’s ability to adapt to changing
1666 market conditions.

1667 Information and material flows are subject to
1668 delays. Orders placed by downstream agents are
1669 transmitted upstream with a one-week information
1670 delay, while physical shipments experience a two-
1671 week delivery delay. The system is initialized in
1672 a steady state with a target inventory level of 400
1673 units to avoid transient effects at the beginning of
1674 the simulation.

1675 Inventory dynamics incur explicit economic
1676 costs. Each unit of inventory held generates a hold-
1677 ing cost of 0.5 per time step, while each unit of
1678 unmet demand (backorder) incurs a higher penalty
1679 of 1.0, reflecting the greater economic impact of
1680 stockouts. In addition, a minimum inventory cost
1681 of 200 is imposed to model fixed operational ex-
1682 penses independent of inventory fluctuations.

1683 In all experiments, the evaluated agent controls
1684 the retailer role, which is closest to customer de-
1685 mand and therefore most exposed to demand un-
1686 certainty. All other supply-chain roles (wholesaler,
1687 distributor, and factory) follow predefined equation-
1688 based policies.

1689 **Opponent policies** Non-controlled agents in the
1690 supply chain follow fixed equation-based ordering
1691 rules. Under the *typical* policy, the order quantity
1692 at time step t is determined by an inventory and
1693 backlog correction rule:

$$1694 q_t = d_t + \alpha(I^* - I_t) + \beta B_t, \quad (1)$$

1695 where d_t denotes the observed demand, I_t is the
1696 current inventory level, I^* is the target inventory,
1697 and B_t represents the backlog. The parameters
1698 α and β control the adjustment speed toward the
1699 target inventory and backlog compensation, respec-
1700 tively. And in our experiment we set $\alpha = 1$ and
1701 $\beta = 1$ as fix parameters.

1702 Under the *smoothing_4* policy, demand is first
1703 smoothed using a four-step moving average:

$$1704 \tilde{d}_t = \frac{1}{4} \sum_{k=0}^3 d_{t-k}, \quad (2)$$

1705 and the smoothed demand \tilde{d}_t is then substituted
1706 for d_t in the ordering rule:

$$1707 q_t = \tilde{d}_t + \alpha(I^* - I_t) + \beta B_t, \quad (3)$$

1708 A.6 Enterprise Digital Twin (EDT) 1709 Implementation Details

1710 This appendix provides low-level execution details
1711 of the Enterprise Digital Twin (EDT) task, comple-
1712 menting the main text description. We focus on (i)
1713 scenario specification, (ii) step-wise simulation dy-
1714 namics, (iii) project life-cycle and stochastic events
1715 (extension and follow-on), (iv) firm-level account-
1716 ing and metrics, and (v) the evaluation protocol
1717 used in our benchmark.

1718 Here is a brief version of the scenario applied in
1719 our task, and we use it as an example for illustrating
1720 the mechanism of EDT environment. The details
1721 of this task scenario can be found in our codes.

```
1722 "scenarios": {  

1723   "interactive": {  

1724     "runspecs": {  

1725       "starttime": 1,  

1726       "stoptime": 96,  

1727       "dt": 1  

1728     },  

1729     "properties": {  

1730       "revenue_risk_level": {  

1731         "type": "Double",  

1732         "value": 0.5  

1733       },  

1734       "fixed_cost": {  

1735         "type": "Double",  

1736       }  

1737     }  

1738   }  


```

```

1739     "value": 20000.0
1740   }
1741 },
1742
1743 "agents": [
1744   {
1745     "name": "consultant",
1746     "count": 1,
1747     "properties": {
1748       "name": {
1749         "type": "String",
1750         "value": "Consultant 1"
1751       },
1752       "salary": {
1753         "type": "Double",
1754         "value": 6000.0
1755       },
1756       "workplace_cost": {
1757         "type": "Double",
1758         "value": 2000.0
1759       }
1760     }
1761   },
1762   ... # 11 other consulatant agents
1763
1764   {
1765     "name": "project",
1766     "count": 1,
1767     "properties": {
1768       "name": {
1769         "type": "String",
1770         "value": "Project 1: Core
1771           Upgrade"
1772       },
1773       "contracted_effort": {
1774         "type": "Double",
1775         "value": 70.0
1776       },
1777       "contracted_probability": {
1778         "type": "Double",
1779         "value": 1.0
1780       },
1781       "extension_probability": {
1782         "type": "Double",
1783         "value": 0.25
1784       },
1785       "extension_effort": {
1786         "type": "Double",
1787         "value": 10.0
1788       },
1789       "follow_on_probability": {
1790         "type": "Double",
1791         "value": 0.1
1792       },
1793       "is_follow_on": {
1794         "type": "Boolean",
1795         "value": false
1796       },
1797       "deadline": {
1798         "type": "Double",
1799         "value": 30.0
1800       },
1801       "consultants": {
1802         "type": "Double",
1803         "value": 2.0
1804       },
1805       "start_time": {
1806         "type": "Double",
1807         "value": 1.0
1808

```

```

    },
    "billing_rate": {
      "type": "Double",
      "value": 16000.0
    }
  }
},
... # 9 other project agents
]

```

A.6.1 Scenario specification

An EDT episode is defined by a JSON scenario under a scenario manager (e.g., smEDT), consisting of three top-level blocks: runspecs, properties, and agents. The runspecs block defines the discrete simulation horizon via starttime, stoptime, and dt. In our benchmark implementation, the simulator advances by discrete step calls until termination, and stoptime acts as a step limit. The scalar dt is used as a per-step scaling factor for both work delivery and cost accrual.

The properties block contains firm-level parameters, most importantly the global fixed_cost and revenue_risk_level. The agents block instantiates a multiset of consultant agents and project agents, plus a controlling component that aggregates flows into evaluation metrics.

A tested agent does not act during the episode. Instead, it outputs a compact scenario-level decision schema $\{C, R, P\}$, where C is the retained number of consultants, R is the global revenue_risk_level, and P encodes per-project acceptance and start/deadline windows. This schema is applied as a constrained transformation to a template scenario to produce a materialized scenario JSON. Non-controllable template fields (e.g., salary, workplace cost, and baseline project parameters) remain unchanged. The simulator then executes the materialized scenario for the full horizon and returns step-wise and terminal outcomes.

A.6.2 Consultant dynamics and capacity constraints

Each consultant agent is characterized by two per-step cost parameters: salary and workplace_cost. These costs are counted every step regardless of utilization and same for every consultant in our experiment settings. Let N_c denote the retained number of consultants (set by the tested agent), s denote the salary and w denote the workplace cost. At each step t , the total consultant operating cost

contribution is

$$\text{Cost}_t^c = N_c \cdot (s + w) \cdot dt. \quad (4)$$

Consultants provide the only labor capacity for project delivery. At any step, a consultant can contribute to at most one project. Moreover, each project j has an integer staffing requirement req_j (scenario field consultants) that caps concurrent workers on that project: at any step, no more than req_j consultants can deliver effort to project j . The simulator enforces sticky assignment: once a consultant begins working on a project, it remains assigned to that project until the project finishes its current workload (base scope and any triggered extension), after which the consultant becomes available for reassignment. This stickiness induces non-trivial opportunity costs: starting a long project early can lock capacity and delay higher-margin projects.

Per-step work delivery is modeled in units of *effort*. For a working consultant, delivered effort per step is scaled by dt . Let $k_{j,t}$ be the number of consultants actually working on project j at step t , with $0 \leq k_{j,t} \leq \text{req}_j$. Then the delivered effort to project j at step t is upper bounded by

$$e_{j,t} \leq k_{j,t} \cdot dt. \quad (5)$$

This definition is consistent with the main-text statement that a consultant produces one unit of effort per step when $dt = 1$ in our settings.

A.6.3 Project state

Each project j is parameterized by:

contracted_effort: base scope $\text{Effort}_j^{\text{base}}$

billing_rate: per-unit revenue Rate_j

consultants: staffing cap req_j

start_time: start time start_j

deadline: deadline dead_j

contracted_probability: reliability π_j^{base}

extension_probability: the probability that a project will extend and requires more effort.

extension_effort the required effort in extension scope $\text{Effort}_j^{\text{ext}}$

follow_on_probability follow-on probability of a project π_j^{fo}

A project is active only within its permissible window. Work can accrue only when $t \geq \text{start}_j$ and $t \leq \text{dead}_j$ and the episode has not terminated. Projects maintain a remaining workload state $E_{j,t}$ (initialized to $\text{Effort}_j^{\text{base}}$). Given delivered effort $e_{j,t}$, the update is

$$E_{j,t+1} = \max(0, E_{j,t} - e_{j,t}). \quad (6)$$

Revenue is generated proportional to delivered effort. The instantaneous revenue contribution from project j at step t is

$$\text{Rev}_{j,t} = e_{j,t} \cdot \text{Rate}_j, \quad (7)$$

and total step revenue is

$$\text{Rev}_t = \sum_j \text{Rev}_{j,t}. \quad (8)$$

If the project reaches its deadline with $E_{j,t} > 0$, the unfinished portion is not deliverable after dead_j , meaning it cannot produce further revenue within that project instance.

A.6.4 Extensions and follow-on projects

EDT introduces two stochastic mechanisms that can create additional revenue opportunities while consuming capacity and increasing uncertainty: extensions (scope growth) and follow-on projects.

Extensions. If a project completes its base scope strictly before its deadline (i.e., $E_{j,t} = 0$ at some $t < \text{dead}_j$), an extension event may trigger. Let $g(R)$ be the risk gating function induced by the global revenue_risk_level $R \in [0, 1]$. Operationally, the extension triggers when a project-level draw passes a threshold that depends on both π_j^{ext} and R . When triggered, the project remaining workload is increased by $\text{Effort}_j^{\text{ext}}$:

$$E_{j,t} \leftarrow E_{j,t} + \text{Effort}_j^{\text{ext}}. \quad (9)$$

The project then continues to consume consultant capacity and can generate additional revenue as the extension workload is delivered, subject again to the deadline and episode termination.

Follow-on projects. At a project's deadline step $t = \text{dead}_j$, a follow-on opportunity may be instantiated. As with extensions, instantiation is gated by the global risk level R and the project parameter π_j^{fo} . A follow-on project is created as a new project agent with $\text{is_follow_on} = \text{true}$. It inherits the primary economic structure of its parent project (e.g.,

similar staffing requirement and billing rate) while suppressing further follow-on chaining (follow-on probability set to zero), preventing infinite cascades. The follow-on has its own start time after creation and competes for the same consultant pool.

These stochastic mechanisms create non-linear portfolio effects. High R increases the chance of extensions and follow-ons, potentially improving earnings but also locking capacity and raising uncertainty; low R yields more predictable revenue trajectories but less upside.

A.6.5 Firm-level accounting and returned metrics

At each step t , the firm accrues expenses and revenue, which are aggregated by the controlling component into cumulative metrics. Let $\text{Cost}_t^{\text{fix}} = \text{fixed_cost} \cdot dt$ denote step fixed cost, and let Cost_t^c denote consultant costs from Eq. (1). Then total step expenses are

$$\text{Exp}_t = \text{Cost}_t^{\text{fix}} + \text{Cost}_t^c. \quad (10)$$

Cumulative revenue and expenses are computed as running sums:

$$\text{AcRev}_T = \sum_{t=1}^T \text{Rev}_t, \quad (11)$$

$$\text{AcExp}_T = \sum_{t=1}^T \text{Exp}_t. \quad (12)$$

Cumulative earnings (profit) are then

$$\text{Earnings}_T = \text{AcRev}_T - \text{AcExp}_T. \quad (13)$$

Utilization is computed from the fraction of consultant capacity actively engaged in delivery. Let busy_t denote the number of consultants assigned to any project at step t . Then step utilization is

$$\text{Util}_t = \frac{\text{busy}_t}{N_c}, \quad (14)$$

with $\text{Util}_t = 0$ when $N_c = 0$. The simulator reports both per-step utilization and an overall average utilization computed across steps.

A.6.6 Evaluation protocol in the benchmark

For each tested model, evaluation proceeds over a fixed set of template scenarios. For each episode, the model is provided with a structured description of (i) horizon, (ii) cost parameters, and (iii) project

parameters. It outputs the schema $\{C, R, P\}$ subject to strict formatting constraints. The evaluator materializes a new scenario by applying the schema to the template (disabling projects, adjusting start/deadline windows, setting R , and selecting the first C consultants). The BPTK server is then launched to execute the scenario, and the evaluator reads step-wise outputs to compute and store metrics (including cumulative earnings, revenue, expenses, cash, utilization, and revenue risk). When multiple runs per scenario are enabled, the model may adapt its schema based on prior-run feedback, enabling learning-style search over scenario configurations under a fixed action space.

A.7 Dataset Statistics

We provide the detailed sample distribution of the FirmBench task corpus. Table 5 summarizes the number of samples in the training, validation, and testing sets for each component. Note that for our currently implemented *online* evaluation mode, only the *test* set is utilized for evaluation and sequential adaptation.

In addition, we introduce the Beer Game and Enterprise Digital Twin (EDT) as simulation-based serious game tasks, which differ fundamentally from the sample-based datasets.

Dataset	Train	Valid	Test	Total
<i>Structured Reasoning</i>				
CodeFinQA	4,409	200	788	5,397
CodeTAT-QA	2,654	200	288	3,142
ConvFinQA	133	-	132	265
FinCode	7	2	47	56
finer	1,000	500	441	1,941
FinKnow	100	50	589	739
formula	500	300	200	1,000
FormulaEval	50	-	50	100
SEC-NUM	6,646	200	2,000	8,846
TAT-QA	120	-	120	240
<i>Interactive Consulting</i>	12	-	48	60
Total	15,631	1,452	4,706	21,789

Table 5: Statistics of FirmBench tasks across training, validation, and testing splits.

A.8 AI Assistance Disclosure

AI-based tools were used during the preparation of this work to assist with code development and to correct grammatical and stylistic issues in the manuscript. All scientific content, experimental design, results, and conclusions were conceived, implemented, and verified by the authors.