# CRMWeaver: Building Powerful Business Agent via Agentic RL and Shared Memories

**Yilong Lai[1,2]\***, **Yipin Yang[1]**, **Jialong Wu[2]**, **Fengran Mo[3]**, **Zhenglin Wang[2]**, **Ting Liang[1]**,
**Jianguo Lin[1]**, **Keping Yang[1]†**

[1] Taobao & Tmall Group of Alibaba     [2] Southeast University     [3] University of Montreal

laiyilong0@gmail.com, {uipin.yang, kuiyu.lt, tengxiao, shaoyao}@alibaba-inc.com

## Abstract

Recent years have witnessed the rapid development of LLM-based agents, which shed light on using language agents to solve complex real-world problems. A prominent application lies in business agents, which interact with databases and internal knowledge bases via tool calls to fulfill diverse user requirements. However, this domain is characterized by intricate data relationships and a wide range of heterogeneous tasks, from statistical data queries to knowledge-based question-answering.

To address these challenges, we propose CRMWeaver, a novel approach that enhances business agents in such complex settings. To acclimate the agentic model to intricate business environments, we employ a synthesis data generation and RL-based paradigm during training, which significantly improves the model's ability to handle complex data and varied tasks. During inference, a shared memories mechanism is introduced, prompting the agent to learn from task guidelines in similar problems, thereby further boosting its effectiveness and generalization, especially in unseen scenarios. We validate the efficacy of our approach on the CRMArena-Pro dataset, where our lightweight model achieves competitive results in both B2B and B2C business scenarios, underscoring its practical value for real-world applications.

## 1 Introduction

With the advanced reasoning and decision-making capabilities, agentic systems based on large language models (LLMs) (Weng, 2023; Yao et al., 2023) are poised to advance artificial intelligence toward more human-like intelligence by autonomously interacting with external environments to accomplish complex goals. A large number of tasks have proven to be effective using an agent-based framework, including coding (Jimenez et al.,
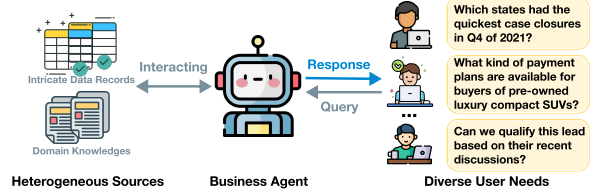


Figure 1: Challenges of business agents in handling diverse user needs and interacting with heterogeneous data sources.

2023; Wang et al., 2025a), web browsing (Zhou et al.; Deng et al., 2023), science discovery (Zhao et al., 2025; Xiang et al., 2025), and deep research (Wu et al., 2025a; Li et al., 2025).

Business agents constitute a prominent application scenario, where systems interact with enterprise databases and internal knowledge bases through tool calls to satisfy diverse user requirements. Unlike general-purpose agents, this domain is marked by highly intricate data dependencies (e.g., large-scale interconnected tables) and task heterogeneity that spans from statistical data queries to knowledge-based question answering. Figure 1 depicts the challenges of business agent. As a result, only proficient but prohibitively costly models have been able to achieve relatively satisfactory performance in these contexts.

To address these challenges, we propose **CRMWeaver**, a novel framework for building robust business agents through a two-fold optimization strategy. First, to adapt the model to the intricacies of business environments while mitigating the scarcity of high-quality training data, we employ large language models to synthesize complex and diverse training data. Based on the synthesis data, we introduce a two-stage training paradigm. In the first stage, we utilize rejection sampling and trajectory distillation to bootstrap the model, enabling it to solve the queries through multi-turn interaction. In the second stage, we employ reinforcement

learning with Decoupled Clip and Dynamic Sampling Policy Optimization(DAPO) (Yu et al., 2025) to further enhance the model's generalization in complex business environments[1]. Second, to cope with the diversity of user queries and the continual evolution of business schemas, we introduce a long-term memory module to improve the agent's stability across multiple tasks, particularly for previously unseen tasks. Specifically, we construct an index of task-specific guidelines built from successful training trajectories. At inference time, when similar tasks are identified, relevant memories are retrieved and injected into the current context, enabling the agent to leverage prior knowledge and improve task performance.

Our contributions are threefold:

1. We address the scarcity of high-quality training data and the need to adapt to complex business environments by combining synthetic data generation with a two-stage training pipeline.

2. We introduce a long-term memory module that enhances robustness and generalization by efficiently reusing solution guidelines from a stronger reasoning model, enhancing the agent's capabilities on diverse, complex tasks.

3. A lightweight model built on Qwen3-4B achieves performance comparable to Qwen3-235B-A22B-Instruct and Gemini 2.5-Pro on the business-agent benchmark CRMArena-Pro across B2B and B2C scenarios, demonstrating the effectiveness of CRMWeaver.

## 2 Preliminaries

To demonstrate the effectiveness of CRMWeaver, we choose to train and test on CRMArena (Huang et al., 2025a) and CRMArena-Pro (Huang et al., 2025b). Both provide realistic sandbox environments for evaluating LLM agents in a business environment, with a particular emphasis on complex multi-table dependencies and diverse enterprise tasks. CRMArena models 16 interconnected Salesforce objects with latent variables that capture hidden causal relations across thousands of records, enabling nine expert-validated customer service tasks covering service managers, analysts, and

agents. Building upon this foundation, CRMArena-Pro substantially expands both the structural and functional scope: it integrates 25 Salesforce objects across Service, Sales, and Configure-Price-Quote schemas, resulting in over 80,000 records across B2B and B2C organizations, and introduces nineteen tasks spanning four core business skills (database querying, textual reasoning, workflow execution, and policy compliance). These benchmarks collectively stress-test agents' abilities to navigate highly interconnected enterprise data, reason across both structured and unstructured records, and handle complex tasks through multi-turn interactions, thereby offering a uniquely challenging and diverse evaluation ground for advancing business-oriented LLM agents.

The task-solving agent in business environments follows the standard agent-environment interaction paradigm. At each time step $t$, the agent receives an observation $o_t \in \mathcal{O}$ from the business environment and selects an action $a_t \in \mathcal{A}$ based on the context $c_t = (o_1, a_1, \cdots, o_{t-1}, a_{t-1}, o_t)$, where $\mathcal{O}$ and $\mathcal{A}$ denote the observation and action spaces, respectively. This interaction process can be naturally formulated as a Partially Observable Markov Decision Process, characterized by the tuple $\langle \mathcal{U}, \mathcal{S}, \mathcal{A}, \mathcal{O}, \mathcal{T}, \mathcal{R} \rangle$. Here, $\mathcal{U}$ denotes the space of initial user queries, $\mathcal{S}$ the latent state space, $\mathcal{T} : \mathcal{S} \times \mathcal{A} \to \mathcal{S}$ the state transition function, and $\mathcal{R}$ denotes the reward function, which indicates task completion. The action space $\mathcal{A}$ corresponds to the set of tools the agent uses to interact with the business environments. In this work, we define three actions: 1) **Execute**, which employs SQL or Salesforce Object Search Language (SOSL) to retrieve relevant records or articles for problem solving. 2) **Date Calculation**, which computes precise temporal information given the current date and an interval count. 3) **Answer**, which allows the agent to finalize and submit a response once it has identified the solution.

The above setup defines the environments, tasks, and interaction paradigm that ground our work. We next introduce CRMWeaver, outlining its methodology for enhancing the problem-solving business-oriented tasks.

## 3 Methodology

In this section, we outline the methodology of CRMWeaver in four parts. We first introduce our data synthesis approach in § 3.1, designed to sup-

---

[1]To rigorously assess generalizability while mitigating data leakage, we conduct model training on CRMArena and evaluation on CRMArena-Pro, which encompass entirely distinct environments.
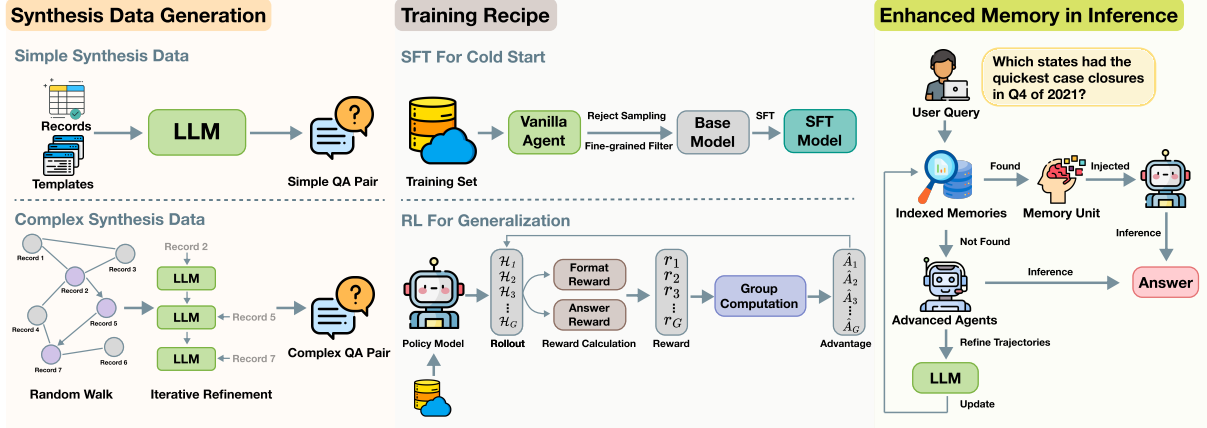
Figure 2: The overall pipeline of CRMWeaver. We first construct both simple and complex synthetic data by leveraging correlations among data records to support model training (§ 3.1). The training recipe adopts a two-stage paradigm, including SFT for model initialization and reinforcement learning for generalization (§ 3.3). Finally, the inference process incorporates a memory enhancement mechanism (§ 3.4).

port cold-start training by synthesizing challenging query and QA tasks that foster the model's ability to reason over complex and heterogeneous data. Then, we introduce the overall paradigm and setup of CRMWeaver in § 3.2. § 3.3 presents a two-stage training pipeline, consisting of supervised fine-tuning (SFT) for initialization, followed by a reinforcement learning phase with DAPO to enhance generalization in complex environments. Lastly, § 3.4 describes the integration of long-term memory at inference, which further strengthens the model's robustness and adaptability across diverse tasks, particularly when addressing previously unseen scenarios. Figure 2 provide an overview of our proposed approach.

### 3.1 Synthesis Data Construction

Given that CRMArena provides only a limited set of task-specific test queries targeting particular objects, we aim to enhance the model's generalization and reasoning capabilities over complex, multi-relational, and heterogeneous data. Inspired by recent work leveraging synthetic complex QA datasets to improve deep research agents (Wu et al., 2025a; Li et al., 2025; Shi et al., 2025), we construct a graph $G$, where nodes correspond to records and edges represent inter-table relationships, thereby capturing the relational structure inherent in CRMArena. For each synthesis data sample $X = (Q, A)$, we first extract a path $\Psi$ through a random walk on the graph $G$:

$$\Psi = (e_1, e_2, \cdots, e_k) \sim \text{RandomWalk}(G) \quad (1)$$

where $Q$ is the synthesis query, $A$ is the corresponding answer, $e_i$ is the sampled record in the path $\Psi$, and $k$ is the length of the sample path $\Psi$.

For each pair $\{e_{i-1}, e_i\}$ along the path $\Psi$, the connection is established through a specific field both in the records $e_{i-1}$ and $e_i$. Starting from the initial node $e_1$, we take the unique identifier of $e_1$ as the final answer $A$. We then define $q_i$ as a synthetic sub-query, derived from the partial path $(e_1, e_2, \cdots, e_i)$. To construct these sub-queries, we employ a prompt-based LLM strategy, where the LLM generates a more complex query $q_i$ given the preceding query $q_{i-1}$ and the corresponding record pair $\{e_{i-1}, e_i\}$:

$$q_i = \text{LLM}(\{e_{i-1}, e_i\} | q_{i-1}, A) \quad (2)$$

After $k$ iterations, we obtain a high-quality synthetic data $X = (Q, A)$. In addition, to broaden the scope and diversity of training tasks, we further construct synthetic datasets targeting statistical analysis and knowledge-base question answering, by exploiting tabular statistical features and internal knowledge resources. The detailed prompts used for constructing synthesis data are provided in the Appendix A.

Overall, our training corpus encompasses three categories:

- *Complex synthetic data*: Crafted using the method in Eq. (1) and Eq. (2) to strengthen the model's ability to search, integrate, and reason over multi-relational queries across interconnected tables. We provide some examples of complex synthesis queries in Appendix A.2.

3

- *Simple synthetic data*: Designed to instill basic business QA and data analysis skills, serving as a foundation for business problem-solving. (*e.g.* How long after purchase can Shoes & Clothings customers claim store credit for an overcharged order?)

- *Task-specific data*: Consisting of nine well-defined tasks from CRMArena. (*e.g.* Which states had the quickest case closures in Q4 of 2021?)

## 3.2 Agent Setup of CRMWeaver

We build CRMWeaver on top of the ReAct framework (Yao et al., 2023), which integrates reasoning and acting through the Thought–Action–Observation paradigm.

Unlike purely reactive agents that rely solely on observed states, or purely deliberative agents that operate only through internal reasoning, ReAct interleaves natural language reasoning traces with environment interactions. Formally, given a query $Q$, the ReAct agent addresses the task through iterative cycles conditioned on the action space $\mathcal{A}$ and observation space $\mathcal{S}$. At time step $t$, with historical trajectory $\mathcal{H}_t = (Q, \tau_1, \alpha_1, o_1, \cdots, \tau_{t-1}, \alpha_{t-1}, o_{t-1})$, the agent derives the current thought $\tau_t$ and action $\alpha_t$ according to the following equation:

$$(\tau_t, \alpha_t) = \pi_\theta(\tau_t, \alpha_t | \mathcal{H}_t) \qquad (3)$$

where $\pi_\theta$ is the policy of the model $\theta$, and $\alpha_t \in \mathcal{A}$. As mentioned in Section 2, the Action Space $\mathcal{A}$ comprises three tools: **execute**, **date calculation**, and **answer**, which enable the agent to interact with external business environments, perform precise temporal computations, and submit the final answers.

For practical supervised fine-tuning(SFT) and reinforcement learning described in the following section, for each time step $t$, we restrict the thought $\tau_t$, action $\alpha_t$ and observation $o_t$ must be enclosed in <think>, <tool_call> and <tool_response> tags, respectively, while the answer tool content is wrapped within <answer> tags. We provide a detailed trajectories example of CRMWeaver in Appendix C.

## 3.3 Two-stage Training Recipe

### 3.3.1 SFT for Agentic Model Cold Starts

To equip the model with the fundamental capability to address business tasks and ensure it follows the Thought-Action-Observation paradigm during problem-solving, we first leverage a powerful language model, GPT-4.1. Using the ReAct framework, we apply Reject Sampling (Yuan et al., 2023; et.al, 2024) to the training data mentioned in Section 3.1, generating multiple candidate execution trajectories. Subsequently, we filter these trajectories based on two primary criteria:

- Correctness of results: To ensure the accuracy of the trajectory data, we first exclude any trajectories where the final answer does not match the gold answer.

- Correctness of process: A portion of the training data has "None" as the answer. To ensure the model does not arrive at the result through simple guesswork, we discard trajectories with a very short number of execution turns.

This process yields a high-quality training dataset, $\mathcal{D}_{\text{SFT}}$, which contains a large number of execution trajectories. Each sample in the dataset is a high-quality trajectory $\mathcal{H} = (Q, \tau_1, \alpha_1, o_1, \cdots, \tau_t, \alpha_t, o_t)$. Following this, we train the model $\theta$ using supervised fine-tuning (SFT). It is important to note that, consistent with prior work (Chen et al., 2024; Wu et al., 2025a), we mask the loss for the observation part of the trajectories:

$$\mathcal{L}_{\text{SFT}} = -\frac{1}{\sum_{i=1}^{|\mathcal{H}|} \mathbb{I}[x_i \neq o]} \sum_{i=1}^{|\mathcal{H}|} \mathbb{I}[x_i \neq o] \cdot \log p_\theta(x_i \mid \mathcal{S}, x_{<i}; \theta) \qquad (4)$$

where $\mathbb{I}(\cdot)$ is the indicator function, $x_i$ is the i-th token of trajectory sample $\mathcal{H}$, $p_\theta$ is the output probability of model $\theta$ and $\mathcal{S}$ is the system description for solving business tasks.

### 3.3.2 Reinforcement Learning to Generalization

To further enhance the agentic model's capabilities, particularly its generalization in complex business environments, we employed a training strategy based on Reinforcement Learning, which allows the model to refine its policy through multi-turn interactions with an external environment.

In details, for each sample $(Q, A)$ from our training dataset $\mathcal{D}$, we rollout a group of trajectories, denoted as $\{\mathcal{H}_i\}_{i=1}^G$. Subsequently, we optimize the

policy of model $\theta$ through Decouple Clip and Dynamic Sampling Policy (DAPO) (Yu et al., 2025):

$$\mathcal{J}_{\text{RL}}(\theta) = \mathbb{E}_{\substack{(Q,A)\sim\mathcal{D} \\ \{\mathcal{H}_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot|q)}} \left[ \frac{1}{\sum_{i=1}^G |\mathcal{H}_i|} \sum_{i=1}^G \sum_{t=1}^{|\mathcal{H}_i|} \right.$$

$$\left. \min\left( r_{i,t}(\theta)\,\hat{A}_{i,t},\ \text{clip}\left(r_{i,t}(\theta),\, 1-\varepsilon_{\text{low}},\, 1+\varepsilon_{\text{high}}\right)\hat{A}_{i,t} \right) \right]$$

$$\text{s.t.}\quad 0 < \left| \{\mathcal{H}_i \mid \textbf{is\_equivalent}(A, \mathcal{H}_i)\} \right| < G \tag{5}$$

where

$$r_{i,t}(\theta) = \frac{\pi_\theta(\mathcal{H}_{i,t} \mid Q, \mathcal{H}_{i,<t})}{\pi_{\theta_{\text{old}}}(\mathcal{H}_{i,t} \mid Q, \mathcal{H}_{i,<t})},$$

$$\hat{A}_{i,t} = \frac{R_i - \text{mean}\left(\{R_i\}_{i=1}^G\right)}{\text{std}\left(\{R_i\}_{i=1}^G\right)}. \tag{6}$$

We use a rule-based reward function $R$, defined as follows:

$$R(\hat{a}_i, a) = 0.1 \times \text{score}_{\text{format}} + 0.9 \times \text{score}_{\text{answer}} \tag{7}$$

where the $\text{score}_{\text{format}}$ is the format score, which is awarded to the model $\theta$ only when its output contains complete and valid pairs of `<think>`...`</think>`, `<tool_call>`...`</tool_call>` tags, and ended with a valid `<answer>`...`</answer>` tag. The $\text{score}_{\text{answer}}$ is the answer score. For data analysis tasks, this reward is given only if the final answer $\hat{a}_i$ exactly matches the gold answer $a$. For QA tasks, we use the F1 score as the answer reward.

### 3.4 Improving Reasoning with Long-term Memory

For business-oriented tasks in CRMArena and CRMArena-Pro, we observe that while the surface forms of many problems vary, they often share common underlying solution strategies (*e.g.* "Find the agent with the shortest handle time who managed more than one case in Winter 2021" and "In the past 3 months, which agent achieved the lowest average handle time while handling more than 2 cases?"). Moreover, the tasks in business settings are diverse and cover various tasks, often comprising novel, unseen combinations of tasks.

To improve model performance on such heterogeneous business tasks, particularly in generalizing to unseen task types, we draw inspiration from recent work on long-term memory (Tang et al., 2025; Fang et al., 2025). We propose augmenting the inference process of our agent with a memory module that stores and retrieves workflow guidelines

---

**Algorithm 1:** Procedural Steps of the Long-term Memory Module

**Input:** Query $Q$, Indexed Memories $\mathcal{M}$, Threshold $\Phi$, System Prompt $\mathcal{P}$
**Output:** Execution trajectory $S$
**function** INFERENCEWITHMEMORIES($Q$, $\mathcal{M}$)
　　$\mathcal{E}, \phi \leftarrow \texttt{Top1-Retrieval}(Q, \mathcal{M})$　▷ Obtain guideline information and corresponding similarity.
　　**if** $\phi \geq \Phi$ **then**
　　　$\mathcal{P} \leftarrow \texttt{Concat}(\mathcal{P}, \mathcal{E})$　▷ Append guideline in the system prompt.
　　**else**
　　　$\texttt{Update}(Q, \mathcal{M})$　▷ Offline update indexed memories using advanced reasoning model.
　　$S \leftarrow \texttt{Solve}(Q \mid \mathcal{P}; \theta)$　▷ Solve the query based on the agentic model $\theta$.
**return** $S$
**end function**

---

from similar past tasks. Specifically, our method employs an indexed memory of workflow information to guide inference toward empirically proven solution patterns for analogous problems. For tasks that do not have a direct precedent in the indexed memories, we employ an offline refinement step. In this step, a more powerful reasoning model attempts to solve the query and generates a detailed corresponding workflow guideline. This guideline is then added to the memory index, thereby enhancing the model's capacity to handle future, novel task types. Algorithm 1 details the specific procedure for solving tasks with the long-term memory module.

It is noteworthy that during the update operation, we employ an advanced reasoning model twice. The final output's consistency determines whether the indexed memories $\mathcal{M}$ are updated. When storing information in memories, we use an LLM to convert specific trajectories into high-level guidelines, which then constitute the memory unit $\mathcal{E}$:

$$\mathcal{E} = \text{LLM}(\mathcal{H}_{\text{adv}}, Q) \tag{8}$$

where the $\mathcal{H}_{\text{adv}}$ is the trajectories of advanced reasoning model. Table 5 in appendix shows an example of the memory unit $\mathcal{E}$, and we also provide the details for generating and indexing for the memory unit in Appendix B.

| Model | B2B | | | | | B2C | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Workflow | Policy | Text | Database | Avg | Workflow | Policy | Text | Database | Avg |
| GPT-4o-mini | 19.5 | 33.8 | 12.6 | 19.3 | 21.3 | 13.5 | 25.6 | 11.7 | 23.8 | 18.6 |
| GPT-4o | 26.0 | 27.5 | 22.1 | 31.3 | 26.7 | 29.0 | 32.5 | 22.2 | 33.1 | 29.2 |
| Gemini-2.5-Flash | 67.5 | 33.5 | 25.1 | 41.6 | 41.9 | 80.0 | 31.0 | 26.3 | 49.3 | 46.7 |
| OpenAI-o1 | 67.0 | 43.3 | 23.4 | 56.3 | 47.5 | 74.5 | 34.5 | 29.5 | 59.6 | 49.5 |
| Gemini-2.5-Pro | 83.0 | 41.0 | 34.7 | 57.6 | 54.1 | 90.0 | 42.0 | 36.2 | 64.9 | **58.3** |
| Qwen3-235B-A22B-Instruct | 90.0 | 37.7 | 30.1 | 49.8 | 51.9 | 87.5 | 42.0 | 33.5 | 50.0 | 53.3 |
| Kimi-K2 | 92.0 | 42.5 | 34.5 | 48.0 | 54.3 | 93.0 | 41.8 | 36.2 | 52.2 | 55.8 |
| **CRMWeaver** | 90.0 | 27.3 | 32.1 | 73.0 | **55.6** | 90.5 | 31.3 | 33.8 | 72.8 | <u>57.1</u> |
| w.o Shared Memory | 90.5 | 27.5 | 27.8 | 72.0 | <u>54.5</u> | 89.5 | 27.5 | 32.2 | 72.0 | 55.3 |

Table 1: Experiments result on CRMArena-Pro across the B2B and B2C domains. The best average results among all backbones are **bolded**, and the second-best results are <u>underlined</u>.

## 4 Experiments

### 4.1 Experiments Setup

**Implementation Details** The backbone for our model $\theta$ is `Qwen3-4B-Instruct-2507`. Separately, the retrieval model discussed in § 3.4 is `BGE-small-en-v1.5`. We used the llama-factory framework (Zheng et al., 2024) for Supervised Fine-Tuning (SFT) on nearly 3k trajectories. After this, we proceeded with reinforcement learning, employing the verl[2] framework and 3.5k question-answer pairs. In our RL training, we set the rollout size $G = 16$, with $\varepsilon_{\text{low}}$ at 0.2 and $\varepsilon_{\text{high}}$ at 0.28. Moreover, to support stable model training during RL, we employ the official, locally deployable SQLite database[3] for the SQL tool. In the case of the SOSL tool, we follow the original configuration, which relies on invoking the Salesforce API. The threshold $\Phi$ in Algorithm 1 is set to 0.7. During inference, we set the temperature to 0.2 and limited the maximum number of interaction turns to 20. For the subsequent main experiments, the results of the OpenAI and Gemini series models are taken directly from the original paper (Huang et al., 2025b). We have also included results from the Qwen3-235B-A22B-Instruct and Kimi-K2 model, using the ReAct paradigm, for comparison. All experiments are conducted on 2 nodes equipped with 8 x Nvidia H20.

**Datasets Details** As metioned in §2, we use CRMArena-Pro (Huang et al., 2025b) as the testbed of our proposed method, which composed of nineteen distinct, expert-validated tasks categorized into four core business skills:

- **Database Querying & Numerical Computation**: This skill assesses an agent's ability

to formulate structured queries and perform numerical calculations on the retrieved data.

- **Information Retrieval & Textual Reasoning**: This involves processing and reasoning over unstructured text from sources such as knowledge bases or call transcripts to extract insights.

- **Workflow Execution**: This evaluates the agent's capacity to follow pre-defined business processes and execute actions based on specific rules.

- **Policy Compliance**: This tests an agent's ability to verify whether business operations adhere to established company policies and regulations.

The dataset encompasses both Business-to-Business (B2B) and Business-to-Consumer (B2C) scenarios. We adopted the single-turn setting, *i.e.* the user provides the task at once, resulting in a total of 3.8k test data samples. Evaluation in CRMArena-Pro is multifaceted. For sub-tasks such as knowledge QA and sales insight mining, we employed the F1 score as the evaluation metric, following standard settings. For other tasks, we used exact match as the metric.

### 4.2 Main Experiments

As depicted in Table 1, we present a comprehensive evaluation of our proposed method, CRMWeaver, against a suite of robust baseline models. The performance assessment was conducted on the CRMArena-Pro dataset, encompassing both B2B and B2C domains.

Our model demonstrates superior performance, achieving an average score of 55.6 in the B2B domain and a highly competitive score of 57.1 in the B2C domain. Notably, CRMWeaver exhibits

---

[2] https://github.com/volcengine/verl
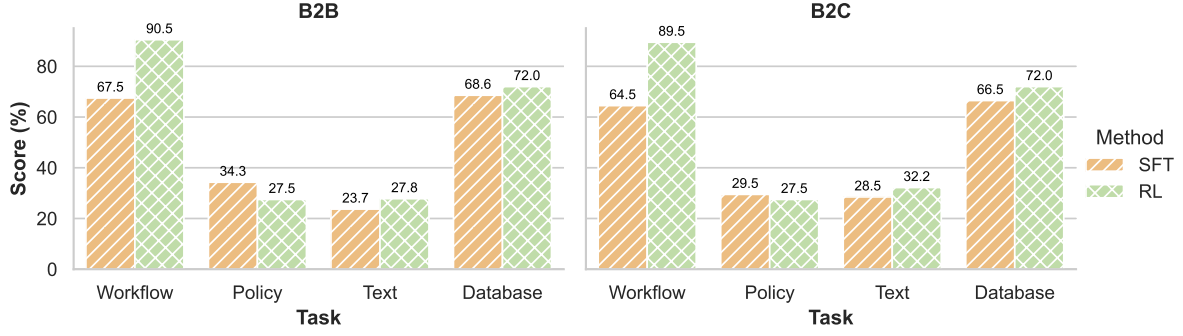[3] https://github.com/SalesforceAIResearch/CRMArena/tree/main/local_data

Figure 3: Performance comparison with SFT and RL on CRMArena-Pro across different tasks.

a pronounced advantage in database-related tasks, where it outperforms all baselines with scores of 73.0 (B2B) and 72.8 (B2C). These results underscore the model's enhanced capability in handling structured data analysis.

To further investigate the contribution of the shared memory component, we conducted an ablation study by removing this module from the agentic model, a variant denoted as "w.o. Shared Memory". The experimental results reveal a discernible degradation in performance, with the average scores declining to 54.5 in the B2B domain and 55.3 in the B2C domain. This outcome indicates that the shared memory module plays a crucial role in facilitating the model's overall performance. However, it is also important to note that even without the aid of shared memories, the agentic model still yielded competitive results, suggesting that the training-time optimization of CRMWeaver remains robust.

### 4.3 Impact of Reinforcement Learning with Agents

To study the impact of reinforcement learning on a model's capacity, we compare a model trained with only supervised fine-tuning (SFT) to one trained with both SFT and reinforcement learning. The experimental results are shown in Figure 3, which demonstrate that the application of RL consistently enhances model performance and generalization across diverse tasks. Notably, the RL-based model demonstrated superior outcomes in both B2B and B2C scenarios, suggesting that RL facilitates a more robust understanding of underlying task dynamics.

The most substantial gains are observed in the Workflow tasks, with the model's performance increasing from 67.5% to 90.5% in the B2B domain

and from 64.5% to 89.5% in the B2C domain. This improvement is attributed to the sequential nature of these tasks, as the RL approach effectively trained the model to follow a correct action sequence, thereby instilling greater confidence in its reasoning and reducing errors.

Conversely, a minor performance decrease was observed on Policy tasks, which was traced to the Solution Violation Identification sub-task. Further analysis revealed that the SFT model frequently outputs 'None' after only a few simple attempts, leading to quick, albeit unreasoned, correct predictions. In contrast, the RL-trained model pursued a more semantically grounded, step-by-step reasoning process. While this behavior demonstrates superior reasoning, it resulted in a lower score on a metric that inadvertently rewards such shortcuts. Potential solutions to these challenges could include synthesizing a more diverse set of similar tasks or designing more robust evaluation mechanisms that reward nuanced reasoning over quick, shortcut-based predictions.

### 4.4 The Frequency of Updates to Long-Term Memory

Our analysis indicates that updates to the Long-term Memory module are relatively infrequent, whereas its indexed memories exhibit consistently high utilization across a broad range of tasks.

In the CRMArena-Pro B2B scenario, the average update frequency is close to 3%, indicating that only 3 out of 100 business tasks, on average, trigger an indexed memories update. We do observe that for a small subset of tasks, such as Knowledge QA, the low similarity among questions results in a low probability of hitting the Long-term Memory. This can lead to the issue of frequent updates. To mitigate this, we employ a few-shot prompt-

7

ing approach with an LLM to classify whether a given query constitutes a knowledge-based question. This mechanism effectively prevents memory updates for QA tasks, thereby avoiding an excessive update frequency.

## 5 Related Works

**Business Agents and Benchmarks.** The evaluation of business agents has progressed from synthetic benchmarks to realistic enterprise environments. Early efforts like Webshop (Yao et al., 2022) relied on simplified webpages, while later benchmarks such as WorkArena (Drouin et al., 2024) and WorkBench (Styles et al., 2024) leveraged enterprise platforms (e.g., ServiceNow) and outcome-oriented tasks like email drafting or meeting scheduling. Tau-bench (Yao et al., 2024) further simulated user–agent interactions in domains such as retail and flight booking. CRMArena (Huang et al., 2025a) advanced realism by building on a Salesforce CRM system with expert-verified tasks spanning 16 objects and diverse information sources. Its successor, CRMArena-Pro (Huang et al., 2025b), expanded to 25 objects and 19 sub-tasks across four categories(Workflow Execution, Policy Compliance, Database Querying, and Information Retrieval), offering the most comprehensive and business-aligned evaluation to date.

**Long-term Memory of Language Agent.** Long-term memory has become a key strategy for improving agent performance by leveraging historical trajectories and domain-specific experiences (He et al., 2024; Mo et al., 2025b). Recent works explore diverse mechanisms for selective retention and retrieval. MemoryBank (Zhong et al., 2024) implements human-inspired memory dynamics to enable personalized, context-aware interactions, while MemoryLLM (Wang et al., 2024) encodes knowledge in a latent space to support knowledge editing. Subsequent studies extend this line with structured and transferable memories: Agent Workflow Memory (Wang et al., 2025b) promotes sub-workflow reuse, and Agent KB (Tang et al., 2025) introduces a shared knowledge base with a teacher–student retrieval scheme for cross-task generalization. Mem0 (Chhikara et al., 2025) further enhances scalability by constructing graph-based dynamic memories that capture complex rela-

tional dependencies beyond fixed context windows.

**Agentic Reinforcement Learning.** Agentic RL aims to enhance model capabilities through multi-turn interactions within complex, real-world environments. With the advent of Deepseek-R1 (Guo et al., 2025), rule-based RL becomes a dominant paradigm. A representative study is Search-R1 (Jin et al., 2025), which leverages end-to-end optimization of an LLM to improve its performance in multi-turn Retrieval-Augmented Generation (RAG) systems (Mo et al., 2025a). Subsequently, Re-Tool (Feng et al., 2025) adopted a "Think-Code-Observation" approach, utilizing Python code to enable multi-turn interactions with external code runtimes, thereby enhancing the model's mathematical reasoning abilities. The efficacy of Agentic RL has been further validated by several landmark works, including WebDancer (Wu et al., 2025a), Websailor (Li et al., 2025), and WebShaper (Tao et al., 2025). These systems leverage data synthesis and end-to-end agentic RL to develop powerful deep research agents capable of multi-turn interactions with web pages and search engines, demonstrating the effectiveness of Agentic RL in complex real-world tasks.

## 6 Conclusion

We presented CRMWeaver, a framework for building powerful business agents capable of operating in complex business environments and handling diverse domain-specified tasks.

CRMWeaver addresses two core challenges: data scarcity and domain complexity, through a two-fold optimization strategy. First, we synthesize complex, diverse training data and adopt a two-stage learning paradigm: (i) rejection sampling and trajectories distillation to bootstrap problem-solving competency, followed by (ii) reinforcement learning based on DAPO to strengthen generalization and decision-making in diverse environments. Second, we introduce a long-term memory module that indexes task-specific guidelines distilled from successful trajectories; at inference time, retrieval-augmented conditioning enables the agent to reuse prior knowledge and maintain stability across unseen tasks and evolving schemas. Together, these components yield a practical recipe for deploying business agents that reduce reliance on prohibitively expensive models while improving reliability in real-world, complex business settings.

## Limitation

**Multi-turns Settings.** Due to constraints in both training frameworks and implementation techniques, our reinforcement learning setup primarily focuses on scenarios where the user provides a single query and the agent interacts with tools through multiple tool call/observation turns. We have not yet extended our framework to more complex multi-turn user interactions (Mo et al., 2023) as explored in environments such as CRMArena-Pro (Huang et al., 2025b) and Tau-bench (Yao et al., 2024).

**GPU Utilization During Rollout.** Although we adopt asynchronous rollout optimization, the rollout process remains time-consuming and often leads to underutilized GPU resources. As a result, we are forced to restrict both the context length and the number of tool calls. We plan to address this limitation in future work by leveraging fully asynchronous RL frameworks such as AReal (Fu et al., 2025).

**Context Management** To handle the common issue of context window limitations during training and inference, we impose a simple constraint on the number of records returned, which inevitably restricts the agent's reasoning capacity. We envision adopting more sophisticated methods such as ReSum (Wu et al., 2025b), to manage contextual information more effectively.

**Hardware** Due to hardware constraints, we are unable to conduct agent RL training on larger-scale models (e.g., 14B or 32B parameters).

## References

Zehui Chen, Kuikun Liu, Qiuchen Wang, Wenwei Zhang, Jiangning Liu, Dahua Lin, Kai Chen, and Feng Zhao. 2024. Agent-FLAN: Designing data and methods of effective agent tuning for large language models. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 9354–9366, Bangkok, Thailand. Association for Computational Linguistics.

Prateek Chhikara, Dev Khant, Saket Aryan, Taranjeet Singh, and Deshraj Yadav. 2025. Mem0: Building production-ready ai agents with scalable long-term memory. *arXiv preprint arXiv:2504.19413*.

Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Sam 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.

Alexandre Drouin, Maxime Gasse, Massimo Caccia, Issam H Laradji, Manuel Del Verme, Tom Marty, David Vazquez, Nicolas Chapados, and Alexandre Lacoste. 2024. Workarena: how capable are web agents at solving common knowledge work tasks? In *Proceedings of the 41st International Conference on Machine Learning*, pages 11642–11662.

Aaron Grattafiori et.al. 2024. The llama 3 herd of models. *Preprint*, arXiv:2407.21783.

Runnan Fang, Yuan Liang, Xiaobin Wang, Jialong Wu, Shuofei Qiao, Pengjun Xie, Fei Huang, Huajun Chen, and Ningyu Zhang. 2025. Memp: Exploring agent procedural memory. *arXiv preprint arXiv:2508.06433*.

Jiazhan Feng, Shijue Huang, Xingwei Qu, Ge Zhang, Yujia Qin, Baoquan Zhong, Chengquan Jiang, Jinxin Chi, and Wanjun Zhong. 2025. Retool: Reinforcement learning for strategic tool use in llms. *arXiv preprint arXiv:2504.11536*.

Wei Fu, Jiaxuan Gao, Xujie Shen, Chen Zhu, Zhiyu Mei, Chuyi He, Shusheng Xu, Guo Wei, Jun Mei, Jiashu Wang, Tongkai Yang, Binhang Yuan, and Yi Wu. 2025. Areal: A large-scale asynchronous reinforcement learning system for language reasoning. *Preprint*, arXiv:2505.24298.

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.

Zihong He, Weizhe Lin, Hao Zheng, Fan Zhang, Matt W Jones, Laurence Aitchison, Xuhai Xu, Miao Liu, Per Ola Kristensson, and Junxiao Shen. 2024. Human-inspired perspectives: A survey on ai long-term memory. *arXiv preprint arXiv:2411.00489*.

Kung-Hsiang Huang, Akshara Prabhakar, Sidharth Dhawan, Yixin Mao, Huan Wang, Silvio Savarese, Caiming Xiong, Philippe Laban, and Chien-Sheng Wu. 2025a. CRMArena: Understanding the capacity of LLM agents to perform professional CRM tasks in realistic environments. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 3830–3850, Albuquerque, New Mexico. Association for Computational Linguistics.

Kung-Hsiang Huang, Akshara Prabhakar, Onkar Thorat, Divyansh Agarwal, Prafulla Kumar Choubey, Yixin Mao, Silvio Savarese, Caiming Xiong, and Chien-Sheng Wu. 2025b. Crmarena-pro: Holistic assessment of llm agents across diverse business scenarios and interactions. *arXiv preprint arXiv:2505.18878*.

Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. 2023. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770*.

Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Sercan Arik, Dong Wang, Hamed Zamani, and Jiawei Han. 2025. Search-r1: Training llms to reason and leverage search engines with reinforcement learning. *arXiv preprint arXiv:2503.09516*.

Kuan Li, Zhongwang Zhang, Huifeng Yin, Liwen Zhang, Litu Ou, Jialong Wu, Wenbiao Yin, Baixuan Li, Zhengwei Tao, Xinyu Wang, et al. 2025. Websailor: Navigating super-human reasoning for web agent. *arXiv preprint arXiv:2507.02592*.

Fengran Mo, Yifan Gao, Chuan Meng, Xin Liu, Zhuofeng Wu, Kelong Mao, Zhengyang Wang, Pei Chen, Zheng Li, Xian Li, et al. 2025a. Uniconv: Unifying retrieval and response generation for large language models in conversations. *arXiv preprint arXiv:2507.07030*.

Fengran Mo, Kelong Mao, Ziliang Zhao, Hongjin Qian, Haonan Chen, Yiruo Cheng, Xiaoxi Li, Yutao Zhu, Zhicheng Dou, and Jian-Yun Nie. 2025b. A survey of conversational search. *ACM Transactions on Information Systems*, 43(6):1–50.

Fengran Mo, Kelong Mao, Yutao Zhu, Yihong Wu, Kaiyu Huang, and Jian-Yun Nie. 2023. Convgqr: Generative query reformulation for conversational search. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4998–5012.

Dingfeng Shi, Jingyi Cao, Qianben Chen, Weichen Sun, Weizhen Li, Hongxuan Lu, Fangchen Dong, Tianrui Qin, King Zhu, Minghao Liu, et al. 2025. Taskcraft: Automated generation of agentic tasks. *arXiv preprint arXiv:2506.10055*.

Olly Styles, Sam Miller, Patricio Cerda-Mardini, Tanaya Guha, Victor Sanchez, and Bertie Vidgen. 2024. Workbench: a benchmark dataset for agents in a realistic workplace setting. In *First Conference on Language Modeling*.

Xiangru Tang, Tianrui Qin, Tianhao Peng, Ziyang Zhou, Daniel Shao, Tingting Du, Xinming Wei, Peng Xia, Fang Wu, He Zhu, et al. 2025. Agent kb: Leveraging cross-domain experience for agentic problem solving. *arXiv preprint arXiv:2507.06229*.

Zhengwei Tao, Jialong Wu, Wenbiao Yin, Junkai Zhang, Baixuan Li, Haiyang Shen, Kuan Li, Liwen Zhang, Xinyu Wang, Yong Jiang, et al. 2025. Webshaper: Agentically data synthesizing via information-seeking formalization. *arXiv preprint arXiv:2507.15061*.

Xingyao Wang, Boxuan Li, Yufan Song, Frank F. Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi Song, Bowen Li, Jaskirat Singh, Hoang H. Tran, Fuqiang Li, Ren Ma, Mingzhang Zheng, Bill Qian, Yanjun Shao, Niklas Muennighoff, Yizhe Zhang, Binyuan Hui, Junyang Lin, Robert Brennan, Hao Peng, Heng Ji, and Graham Neubig. 2025a. Openhands: An open platform for AI software developers as generalist agents. In *The Thirteenth International Conference on Learning Representations*.

Yu Wang, Yifan Gao, Xiusi Chen, Haoming Jiang, Shiyang Li, Jingfeng Yang, Qingyu Yin, Zheng Li, Xian Li, Bing Yin, Jingbo Shang, and Julian J. McAuley. 2024. Memoryllm: Towards self-updatable large language models. In *ICML*.

Zora Zhiruo Wang, Jiayuan Mao, Daniel Fried, and Graham Neubig. 2025b. Agent workflow memory. In *Forty-second International Conference on Machine Learning*.

Lilian Weng. 2023. Llm-powered autonomous agents. *lilianweng.github.io*.

Jialong Wu, Baixuan Li, Runnan Fang, Wenbiao Yin, Liwen Zhang, Zhengwei Tao, Dingchu Zhang, Zekun Xi, Gang Fu, Yong Jiang, et al. 2025a. Webdancer: Towards autonomous information seeking agency. *arXiv preprint arXiv:2505.22648*.

Xixi Wu, Kuan Li, Yida Zhao, Liwen Zhang, Litu Ou, Huifeng Yin, Zhongwang Zhang, Yong Jiang, Pengjun Xie, Fei Huang, Minhao Cheng, Shuai Wang, Hong Cheng, and Jingren Zhou. 2025b. Resum: Unlocking long-horizon search intelligence via context summarization. *Preprint*, arXiv:2509.13313.

Yanzheng Xiang, Hanqi Yan, Shuyin Ouyang, Lin Gui, and Yulan He. 2025. Scireplicate-bench: Benchmarking llms in agent-driven algorithmic reproduction from research papers. *arXiv preprint arXiv:2504.00255*.

Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. 2022. Webshop: Towards scalable real-world web interaction with grounded language agents. *Advances in Neural Information Processing Systems*, 35:20744–20757.

Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan. 2024. $\tau$-bench: A benchmark for tool-agent-user interaction in real-world domains. *Preprint*, arXiv:2406.12045.

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023. React: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*.

Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, et al. 2025. Dapo: An open-source llm reinforcement learning system at scale. *arXiv preprint arXiv:2503.14476*.

Zheng Yuan, Hongyi Yuan, Chengpeng Li, Guanting Dong, Keming Lu, Chuanqi Tan, Chang Zhou, and Jingren Zhou. 2023. Scaling relationship on learning mathematical reasoning with large language models. *Preprint*, arXiv:2308.01825.

Yilun Zhao, Kaiyan Zhang, Tiansheng Hu, Sihong Wu, Ronan Le Bras, Taira Anderson, Jonathan Bragg, Joseph Chee Chang, Jesse Dodge, Matt Latzke, et al. 2025. Sciarena: An open evaluation platform for

foundation models in scientific literature tasks. *arXiv preprint arXiv:2507.01001*.

Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyan Luo, Zhangchi Feng, and Yongqiang Ma. 2024. Llamafactory: Unified efficient fine-tuning of 100+ language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, Bangkok, Thailand. Association for Computational Linguistics.

Wanjun Zhong, Lianghong Guo, Qiqi Gao, He Ye, and Yanlin Wang. 2024. Memorybank: Enhancing large language models with long-term memory. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 19724–19731.

Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, et al. Webarena: A realistic web environment for building autonomous agents. In *The Twelfth International Conference on Learning Representations*.

# Appendix

## A Synthesis Data Construction

### A.1 Prompt for Complex Queries Generation

We provide the prompt used to generate the complex queries in Table 2 and Table 3.

### A.2 Complex Synthesis Data Examples

Here are some examples of complex synthesis data:

1. Which live chat transcript documents a case where a Portland-based customer, previously known for resolving product and sizing complaints in 2022, reached out on November 30, 2023, to request an immediate exchange after receiving a blue hoodie instead of black running shoes, and successfully initiated the return process with photo documentation through interactions involving three different support agents? Return only the chat transcript's ID. If no such records exist, return 'None'.

2. During the second week of January 2022, a customer with an account connected to a contact named Isabella Adams, who has an email address of isabella.morgan@example.com and a shipping address in Austin, Texas, placed an order that was activated and used a price book featuring the promotion 'Start 2022 with unbeatable deals on sports gear.' Which price book was referenced in this transaction? Return only the price book's id. If no such records exist, return 'None'.

3. Identify the durable and waterproof item, ideal for harsh weather conditions, that was added as a new selection by the customer whose previous order item was related to outdoor hiking, after a warranty issue with defective equipment was resolved in Portland and resulted in a successful transaction activation on March 19, 2022. What is the unique catalog reference for that item within the product listings? Return only the product's ID. If no such records exist, return 'None'.

## B Long Term Memory Indexing

To generate the long-term memory unit $\mathcal{E}$, we use o3-mini as the advanced reasoning model. After obtain the trajectories $\mathcal{H}_{adv}$, we use GPT4.1 with prompt described in Table 4 to generate the guideline for the given task.

## C Case Study of CRMWeaver

### C.1 Long-term Memory Example

Table 5 provide a example of memory unit $\mathcal{E}$.

### C.2 Full Trajectories Example

Table 6 is a full trajectories example of our proposed agent.

**Prompt for Initialize Seed Query**

You are a synthesis data generator. I will give you two schema and two corresponding data records. The two records are connected by a foreign key relationship. Your task is to generate a complex natural language question that includes the information from the records implicitly. Apart from this, the question's answer should be the id I provided with you. The question should be complex and not directly related to the id.

Input/Output Specifications:
Input:
– Source Schema: A description of the source schema.
– Source Record: A record from the source schema.
– Target Schema: A description of the target schema.
– Target Record: A record from the target schema.
– Connect Key: The connect relationship between the two records.
– Answer: Known correct answer.

Output:
– Must be in strict JSON format: {"Q": "generated question"}
– No explanations or extra fields allowed

Q must satisfy:
1. Be a complete natural language question
2. Allow deriving answer A by using the source and target records and connect key

Question Generation Principles:
1. Exact correspondence – Each question must fully base on the original conclusion, with the answer being its core content.
2. Derivability – The original conclusion must be directly derivable from the question and be the only correct answer.
3. Self–containment – Questions must be complete and independent, not relying on external references or unspecified context.
4. Information hiding – Do not reveal specific sources or data paths, but can include search hints.

5. Specificity and clarity – Questions should include details like specific times to ensure unique answers.
6. Single question – Generate only one question per conclusion.
7. If the conclusion can only be obtained from input content, include hints via data source identifiers in the question.
8. Language consistency – The language of each question must be the same as the conclusion's language.

Table 2: Prompt for Generating Seed Query

## Prompt for Initialize Seed Query

You are a synthesis data generator. I will give you an existing question and its question. Apart from this, I will also give you a schema and a record. The schema and record includes a foreign key exactly the same as the answer of the giving question. Your task is to generate a more complex question that merges the information from the record and the current given record implicitly. Apart from this, the question's answer should be the id I provided with you. The question should be complex and not directly related to the ID. And the question should not include explicit information such as id of the record or the name of the user.

Input/Output Specifications:
Input:
− Existing Question: A complex natural language question that has been previously generated.
− Existing Question's Answer: The answer to the existing question.
− Target Schema: A description of the target schema.
− Target Record: A record from the target schema.
− Connect Key: The connect relationship between the two records.
− Answer: Known correct answer.

Output:
− Must be in strict JSON format: {"Q": "generated question"}
− No explanations or extra fields allowed

Q must satisfy:
1. Be a complete natural language question
2. Allow deriving answer A by using the source and target records and the connect key

Question Generation Principles:
1. Exact correspondence − Each question must be fully based on the original conclusion, with the answer being its core content.
2. Derivability − The original conclusion must be directly derivable from the question and be the only correct answer.
3. Self−containment − Questions must be complete and independent, not relying on external references or unspecified context.
4. Information hiding − Do not reveal specific sources or data paths, but can include search hints.

5. Specificity and clarity − Questions should include details like specific times to ensure unique answers.
6. Single question − Generate only one question per conclusion.
7. If the conclusion can only be obtained from input content, include hints via data source identifiers in the question.
8. Language consistency − The language of each question must be the same as the conclusion's language.

Table 3: Prompt for Generating Deeper Challenge Queries

## Prompt for Generating Guideline

Please provide a concise guideline for the following execution log.
Break down the guideline into individual steps, and for each step, briefly describe what action was taken (note that it must be report the tool's name explicitly) and its outcome.
Ensure the summary is clear, structured, and captures all key actions in the order they occurred.

## Key Requirements
1. Focus exclusively on behavioral improvements derived from similar task patterns and experience.
2. Format output strictly as:
    1. [Specific suggestion 1]
    2. [Specific suggestion 2]
    ...
3. No headings, explanations, or markdown.

## Detailed Schema Information
{schema_information}

## Execution logs:
{successful_trajectories}

Table 4: The prompt used to generating the guideline of memory unit $\mathcal{E}$.

---

**Key:** Can we qualify this lead based on their recent discussions?...
**Value:**

1. Always use the execute tool to pull every relevant transcript and pricing/knowledge record before judging any BANT factor, ensuring no critical data is missed.
2. After data retrieval, perform an explicit cost–vs–budget calculation inside the analysis step to avoid incorrect 'Budget' conclusions.
3. When Authority is not clearly confirmed in the transcript, trigger an additional execute search (e.g., manager titles, purchasing roles) before deciding it is fulfilled or unmet.
4. Parse Timeline details quantitatively (days, dates) rather than qualitatively; if uncertain, run another execute query for scheduling constraints to improve accuracy.
5. Maintain a consistent workflow order –> execute –> analyze –> respond to increase reliability and reproducibility across similar tasks.

Table 5: An example of Memory Unit $\mathcal{E}$, where the **Key** is the indexed query used to match similar queries, and **Value** is the corresponding guideline.

**Question:**

Which agent ought to be designated to handle this lead? Return only the ID of the agent.
## Lead Routing Policy
The lead routing policy determines the best agent to assign the given new lead based on the following criteria:
− Territory Match: First, find the territory that matches the lead's region by comparing it with the territory description (comma−separated list of states/regions).
− Quote Success: Among agents associated with the matching territory, select those with the highest number of quotes with status "Accepted" (find the agent ID through the associated opportunity ID).
− Workload Balance: If there's still a tie, select the agent with the fewest unconverted (open) leads
## Lead's region: WI

**Step 1**

**<think>**
I need to find the territory that matches the lead's region (WI) by checking the Territory2 object. I will look for a territory with "WI" in its Description field.
**</think>**

**<tool_call>**
{"name": "execute", "arguments": {"query": "SELECT Id, Name FROM Territory2 WHERE Description LIKE '%WI%"}}
**</tool_call>**

**<tool_response>**
Observation: [('0MIWt0000007xzROAQ', 'US-Great Lakes')]
**</tool_response>**

**Step 2**

**<think>**
I have found the territory that matches the lead's region (WI), which is "US-Great Lakes" with ID "0MIWt0000007xzROAQ". Now, I need to find the agents associated with this territory by checking the UserTerritory2Association object.
**</think>**

**<tool_call>**
{"name": "execute", "arguments": {"query": "SELECT UserId FROM UserTerritory2Association WHERE Territory2Id = '0MIWt0000007xzROAQ'"}}
**</tool_call>**

**<tool_response>**
Observation: [('005Wt000003NHpdIAG',), ...]
**</tool_response>**

**Step 3**

**<think>**

I have found the list of agents associated with the territory "US-Great Lakes" (WI). Now, I need to determine which of these agents has the highest number of quotes with status "Accepted" by checking the Quote and QuoteLineItem objects.
**</think>**

**<tool_call>**
```
{"name": "execute", "arguments": {"query": "SELECT OwnerId, COUNT(*)
AS accepted_quotes FROM Quote WHERE Status = 'Accepted' AND OwnerId IN
('005Wt000003NHpdIAG', ···, '005Wt000003NJmcIAG') GROUP BY OwnerId ORDER BY
accepted_quotes DESC"}}
```
**</tool_call>**

**<tool_response>**
Observation: Error: no such column: OwnerId
**</tool_response>**

> **Step 4**

**<think>**
The error occurred because the OwnerId field does not exist directly on the Quote object. Instead, the agent (OwnerId) is associated with the Opportunity object, and the Quote is linked to that Opportunity via OpportunityId. I need to check the Opportunity object to find the OwnerId (agent) associated with each Quote.
**</think>**

**<tool_call>**
```
{"name": "execute", "arguments": {"query": "SELECT o.OwnerId, COUNT(*) AS
accepted_quotes FROM Quote q JOIN Opportunity o ON q.OpportunityId = o.Id
WHERE q.Status = 'Accepted' AND o.OwnerId IN ('005Wt000003NHpdIAG', ···,
'005Wt000003NJmcIAG') GROUP BY o.OwnerId ORDER BY accepted_quotes DESC"}}
```
**</tool_call>**

**<tool_response>**
```
Observation:  [('005Wt000003NIowIAG', 5), ('005Wt000003NJmcIAG', 4),···,
('005Wt000003NDXaIAO', 1)]
```
**</tool_response>**

> **Step 5**

**<think>**
The agent with ID '005Wt000003NIowIAG' has the highest number of accepted quotes (5), which is more than any other agent in the list. Since this agent has the highest quote success, they should be allocated to the lead according to the routing policy.
**</think>**

**<answer>**
```
005Wt000003NIowIAG
```
**</answer>**

Table 6: A Trajectories Example on CRMArena-Pro.