

Self-Generated In-Context Examples Improve LLM Agents for Sequential Decision-Making Tasks

Vishnu Sarukkai, Zhiqiang Xie, Kayvon Fatahalian
Stanford University

Abstract

Improving Large Language Model (LLM) agents for sequential decision-making tasks typically requires extensive task-specific knowledge engineering—custom prompts, curated examples, and specialized observation/action spaces. We investigate a different approach where agents automatically improve by learning from their own successful experiences without human intervention. Our method constructs and refines a database of self-generated trajectories that serve as in-context examples for future tasks. Even naive accumulation of successful trajectories yields substantial performance gains across three diverse benchmarks: ALFWorld (73% to 89%), Wordcraft (55% to 64%), and InterCode-SQL (75% to 79%). These improvements exceed those achieved by upgrading from gpt-4o-mini to gpt-4o and match the performance of allowing multiple attempts per task. We further enhance this approach with two innovations: database-level curation using population-based training to propagate high-performing example collections, and exemplar-level curation that selectively retains trajectories based on their empirical utility as in-context examples. With these enhancements, our method achieves 93% success on ALFWorld—surpassing approaches that use more powerful LLMs and hand-crafted components. Our trajectory bootstrapping technique demonstrates that agents can autonomously improve through experience, offering a scalable alternative to labor-intensive knowledge engineering.

1 Introduction

When creating LLM agents for sequential decision-making tasks, practitioners often improve agent performance by investing in task-specific knowledge engineering (through tedious prompt tuning [1], human-crafted in-context examples [2, 3] or custom observation and action spaces [4, 5]). Using these techniques, scaling agent performance comes from scaling human effort.

In this paper, we investigate an alternative path: enabling LLM agents to autonomously bootstrap their own performance by leveraging their own successful experiences via in-context learning. The efficacy of in-context learning depends critically on both the quality of the examples [2, 3] and their relevance to the current decision point [6, 7, 8]. This insight provides a natural direction for automated agent self-improvement: accumulating successful self-generated trajectories and estimating the most relevant and effective prior experiences to use as in-context examples for each action.

Our work assumes a ReAct-style agent [9] that retrieves different examples for each decision point based on their relevance to the current situation [10, 11]. We build on this foundation by focusing specifically on how to construct and refine the underlying database of self-generated examples. How can we identify which trajectories enhance performance on new tasks versus those that hinder performance? This database construction problem requires addressing both the collection of high-quality trajectories and the strategic curation of the most valuable ones for future retrieval at each decision point in the agent’s reasoning and acting loop.

We demonstrate that even naive database accumulation improves test-set performance from 73% to 89% on ALFWorld, 55% to 64% on Wordcraft, and 75% to 79% on InterCode-SQL. (Equivalent to what a baseline agent would achieve if it were allowed two to three attempts per task.) We further propose two database construction enhancements: (1) database-level curation that identifies and propagates high-performing example databases, and (2) exemplar-level curation that identifies helpful trajectories based on their empirical utility as in-context examples. These approaches do not require task-specific prompt engineering [12, 4] or custom observation/action space design [11, 4], but improve success rates on ALFWorld to 93%—surpassing approaches like AutoManual [4] that use more powerful LLMs and hand-crafted observation and action spaces, as well as hierarchical approaches like Autoguide [12]. The success rate improvement on ALFWorld exceeds the boost obtained from upgrading the agent’s underlying LLM from gpt-4o-mini to gpt-4o. Our results highlight the practical value of trajectory bootstrapping as a dimension for scaling test-time compute.

2 Preliminaries

2.1 Sequential Decision-Making Tasks

We focus on multi-step sequential decision-making tasks where agents must produce a series of actions over time based on observations of the environment. The sequential nature of these tasks introduces unique challenges for LLM agents, as they must interpret intermediate environmental feedback, maintain coherent reasoning across multiple steps, and adapt their strategy based on the evolving task state. This contrasts with one-shot generation tasks (e.g., solving math problems [13], one-shot code generation [14]) where feedback is only available after the complete solution is provided. Our example-driven learning strategy is potentially also suitable to single-step decision-making tasks, but we focus on the multi-step setting due to its applicability to a number of agentic tasks in real-world settings (embodied agents [15], browser-based tasks [16], etc.).

Formally, these tasks can be modeled as Partially Observable Markov Decision Processes (POMDPs), represented by the tuple $(\mathcal{S}, \mathcal{O}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \gamma)$, where \mathcal{S} denotes the underlying state space, \mathcal{O} the observation space, \mathcal{A} the action space, $\mathcal{T} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ defines the deterministic transition function, $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function, and $\gamma \in [0, 1]$ is the discount factor. The partial observability reflects that agents don’t have direct access to the full environment state but rather receive observations that provide limited information.

Given a task goal g , an episode consists of the agent interacting with the environment for a maximum of T timesteps. At each timestep t , the agent receives an observation $o_t \in \mathcal{O}$ of the current state, takes an action $a_t \in \mathcal{A}$, and the environment transitions to the next state according to the transition function \mathcal{T} . In our setting, we specifically consider sparse-reward environments where success is only determined at the end of an episode—the agent receives $\mathcal{R} = 1$ for successful task completion and $\mathcal{R} = 0$ otherwise. This is a standard setting in prior agentic work [9, 17, 12, 4].

2.2 ReAct-style Agent Loop

Our work assumes a ReAct-style [9] agent architecture that employs recent best practices for in-context retrieval [10, 11]. The agent operates through a three-phase approach (planning, reasoning, and acting) as formalized in Algorithm 1. Two key components differentiate our implementation from basic ReAct:

- To support complex, long-horizon tasks, we incorporate an initial planning step where the agent generates a high-level plan for the entire task before execution begins (1, line 3). Full-task planning has been shown to boost agent performance in prior work [10, 15]. This modification is fairly standard—either built directly into the algorithm [10], or included in the first reasoning step of human-crafted in-context examples [17, 4].

Algorithm 1 ReAct-style Agent Loop

```
1: function AGENT( $g, \mathcal{D}, \mathcal{E}, T$ )
2:    $C_p \leftarrow \text{Retrieve}(\mathcal{D}, \text{keys} = [g])$  ▷ Retrieve for plan
3:    $p \leftarrow \text{LLM}_{\text{plan}}(g, C_p)$  ▷ Generate initial plan
4:   Initialize  $\tau \leftarrow (g, p, \{\}, -)$ 
5:    $o_1 \leftarrow \mathcal{E}.\text{obs}()$ 
6:    $C_1 \leftarrow \text{Retrieve}(\mathcal{D}, \text{keys} = [g, p, o_1])$  ▷ Retrieve for current observation
7:   for  $t = 1$  to  $T$  do
8:      $r_t \leftarrow \text{LLM}_{\text{reason}}(\tau, o_t, C_t)$  ▷ Generate reasoning
9:      $C_{t+1} \leftarrow \text{Retrieve}(\mathcal{D}, \text{keys} = [g, p, r_t])$  ▷ Retrieve for current reasoning
10:     $a_t \leftarrow \text{LLM}_{\text{act}}(\tau, o_t, r_t, C_{t+1})$  ▷ Decide action
11:     $o_{t+1}, \text{done}, s \leftarrow \mathcal{E}.\text{step}(a_t)$  ▷ Execute action in environment
12:     $\tau \leftarrow \tau \cup (o_t, r_t, a_t)$ 
13:    if done then
14:      return  $(g, p, \{(o_i, r_i, a_i)\}_{i=1}^t, s)$ 
15:  return  $(g, p, \{(o_i, r_i, a_i)\}_{i=1}^T, 0)$  ▷ Failed due to timeout
```

- Rather than using the same per-task examples throughout an episode [9, 17, 12], we follow retrieve different trajectory segments for each decision point, ensuring the agent has access to the most relevant information at each step [10, 11]. See Appendix A for details.

The agent operates through three key LLM-based functions:

1. LLM_{plan} generates a high-level plan p for achieving the goal
2. $\text{LLM}_{\text{reason}}$ processes the current observation o_t to produce reasoning r_t
3. LLM_{act} determines the appropriate action a_t based on the reasoning

The $\text{Retrieve}()$ function selects the k most relevant examples from database \mathcal{D} based on the average cosine distance from the provided lookup keys to the corresponding examples in \mathcal{D} —see Appendix A and Algorithm 4 for details. The environment \mathcal{E} provides observations and processes actions, returning the next observation, a termination signal, and the success indicator when the episode ends.

This dynamic retrieval approach is critical for sequential decision-making tasks, as it allows the agent to access specialized knowledge relevant to each unique situation encountered during task execution. Our contribution builds upon this architecture by focusing specifically on how to construct and refine the underlying trajectory database that powers this retrieval mechanism.

Note that Algorithm 1 avoids strategies that employ task-specific prompting, observation spaces [11] or action spaces [5, 4]. The only task-specific knowledge is encapsulated in the content of the trajectory database \mathcal{D} . For simplicity, we eschew other techniques, like hierarchical learning [17, 12, 4], that are also task-agnostic, but add additional agent complexity. We view the benefits of hierarchical learning as orthogonal and complimentary to our database construction focus.

3 Problem Statement

Given the ReAct-style agent described in Section 2, our goal is to construct a trajectory database that maximizes LLM agent performance across sequential decision-making tasks. We focus specifically on how to build and refine the database of examples that the agent retrieves from at each decision point.

Formally, we aim to construct a database \mathcal{D} of trajectories, where each trajectory $\tau \in \mathcal{D}$ captures a complete task attempt:

$$\tau = (g, p, \{(o_t, r_t, a_t)\}_{t=1}^T, s)$$

We aim to maximize the agent’s expected performance across a distribution of tasks \mathcal{T} :

$$\mathcal{D}^* = \arg \max_{\mathcal{D}} \mathbb{E}_{g \sim \mathcal{T}} [\text{Success}(\text{Agent}(g, \mathcal{D}, \mathcal{E}, T))]$$

where $\text{Success}()$ returns the binary outcome s of the agent’s execution.

We assume that we are given: (1) \mathcal{D} initialized with a small number of human-generated trajectories, (2) a descriptor of the action space, and (3) access to a set of training tasks drawn from \mathcal{T} that the agent can attempt. All three assumptions are typical in ReAct-based agentic setups [9, 10, 17, 12, 4]. Given this setup, our focus is on the agent self-generating trajectories, then choosing the trajectories that should be added to \mathcal{D} to maximize the agent’s expected performance on novel tasks.

This problem presents two key questions:

- **Scaling Trajectory Collection:** As we gather more successful trajectories, does a larger database provide better guidance for future tasks?
- **Intelligent Trajectory Curation:** How can we determine which subset of collected trajectories will most effectively support the agent in solving new tasks?

4 Related Work

In-context learning for agent improvement Despite the current popularity of reinforcement learning-based approaches for improving agent capabilities [18, 19, 20, 21], in-context learning offers distinct scientific and practical advantages. These benefits include model-agnostic portability across different LLMs, efficiency in low-sample regimes [3, 22], and accessibility when implementation barriers exist for weight modification methods. Both empirical and theoretical work has established that in-context performance can scale effectively with additional examples [6, 7, 22, 8], suggesting that strategic example accumulation should lead to significant performance improvements. We focus on maximizing the value of limited examples through in-context methods, while hypothesizing that database quality, not just quantity, critically influences performance scaling. For completeness, we offer a preliminary investigation in App. C of how our collected trajectories could potentially serve as training data for fine-tuning approaches.

Automatic in-context examples Recent work has demonstrated the effectiveness of optimizing both instructional content and example curation in prompts. DSPy [23] introduced a framework for optimizing multi-step pipelines through instruction tuning and strategic example curation. Self-generated examples containing reasoning traces can eliminate the need for human-written examples, and these self-generated examples often contribute more to performance than optimized instructions alone [24]. These approaches typically select fixed exemplars for all task instances, whereas our method enables the dynamic selection of different in-context examples for each decision.

In-context self-improvement of LLM Agents Self-improvement methods for LLM agents either aim to solve one task (performing search/optimization) or transfer knowledge from prior tasks to novel ones (generalization) (see App. F.2 for further discussion). Approaches to solve a single task scale the number of sampled solutions at inference time [25, 26, 27] or incorporate feedback from failed attempts [28]. Knowledge transfer approaches include abstraction-based methods like ExpeL [17] and AutoGuide [12], while others employ task-specific information in their design—RAP [10] uses task-specific prompts and AutoManual [4] constructs task-specific state and action spaces (see App. F.1). Other dimensions of

self-improvement include hierarchical execution [29] and optimization techniques for multi-stage systems [30, 31, 32]—techniques complementary to our approach. Rather than developing complex architectures or leveraging task-specific information, we focus on identifying which trajectories most contribute to successful outcomes as in-context examples.

5 Methods

In this section, we discuss three algorithms for constructing database \mathcal{D} using a continual collection approach.

5.1 *Traj-Bootstrap*: Constructing a Database of Previously-Solved Tasks

Our trajectory-bootstrapping algorithm *Traj-Bootstrap* constructs a trajectory database \mathcal{D} by collecting successful agent experiences. As outlined in Section 3, we start with a minimal set of human-provided exemplars (which could be empty), then grow the database as the agent successfully completes training tasks. This process creates a positive feedback loop where successful examples help the agent solve new tasks, generating more successful examples.

Traj-Bootstrap operates on principles similar to reward-weighted regression in reinforcement learning [33], where only successful trajectories ($s = 1$) are stored in the database. This filtering mechanism ensures the agent learns from positive examples while avoiding potentially misleading failed attempts. Successful trajectories can be leveraged by asking the agent to imitate the successful patterns in these trajectories. However, failed trajectories are more challenging to operationalize due to the credit attribution problem: it is necessary to identify the ‘good’ vs ‘bad’ parts of the trajectory, before we can even guide the agent to imitate the good parts and avoid the mistakes made in the bad parts. Failed trajectories do offer the opportunity to guide exploration [28]; we leave this direction to future work.

5.2 *+DB-Curation*: Database-Level Data Curation

The *Traj-Bootstrap* algorithm displays unpredictable performance variation across training trials, even when following identical collection procedures. Figure 1 illustrates this variation across five trials on the InterCode-SQL benchmark (a benchmark we use for evaluation in Section 6). The variance arises from two factors: (1) the stochasticity of LLM outputs creating different initial trajectories, and (2) an amplification effect where early differences in collected examples lead to wide performance variation.

This observation motivates a data curation strategy inspired by population-based training in reinforcement learning [34]. Figure 1 shows that some databases lead to better task performance than others—so we identify the underperforming databases periodically during training and remove them, continuing growth from top-performing ones. We introduce *+DB-Curation*, a population-based training algorithm (Algorithm 2) to identify and propagate the most effective databases during the bootstrapping process.

We maintain N database instances initialized with identical human-provided exemplars. Each instance is used by a separate agent, accumulating successful trajectories independently. We trigger curation events when the number of tasks attempted reaches size thresholds (starting at size 10 and doubling thereafter: 10, 20, 40, 80, etc.). At each threshold, we evaluate database performance based on the agent’s success rate on all training tasks since the last threshold, and we replace the worst-performing database with a copy of the top-performing database.

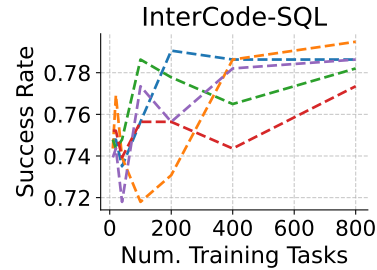


Figure 1: ***Traj-Bootstrap* leads to variance in test-time success rate.** Individual trials (5) shown as dashed lines, results on Intercode-SQL benchmark. There is noticeable variability in performance across trials.

Algorithm 2 Database Curation Logic for +DB-Curation

```
1: procedure OPTIMIZEDATABASES( $\{\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_N\}, interval$ )
2:   Initialize performance metrics  $\{m_1, m_2, \dots, m_N\}$  for each database
3:   for  $t = 1, 2, \dots, T_{train}$  do
4:     for  $i = 1, 2, \dots, N$  in parallel do
5:       Execute task  $t$  using database  $\mathcal{D}_i$ 
6:       If successful, add trajectory to  $\mathcal{D}_i$ 
7:       Update rolling performance metric  $m_i$  on recent tasks
8:     if  $t = 10 \times 2^j$  for any  $j \in \mathbb{N}$  then
9:       Sort databases by rolling performance on recent tasks
10:      Replace worst database with copies of best
```

The key insight of this approach is that database quality emerges from collective properties—like coverage, diversity, and complementarity across examples—not just individual trajectory quality. Moreover, a single trajectory collected early in training can influence many future trajectories by guiding the agent toward particular solution strategies, creating cascading database-level effects. By selecting and propagating entire databases, we preserve these beneficial emergent properties while using a simple, computationally efficient evaluation metric based on recent performance.

5.3 +Exemplar-Curation: Exemplar-Level Data Curation

While the database-level curation performed by +DB-Curation identifies entire sets of complementary trajectories, discarding whole databases can eliminate valuable trajectories. We find that even poor-performing databases contain individual high-quality trajectories that yield better outcomes when used as examples. Conversely, some trajectories marked as successful lead to failures when retrieved—because some trajectories lead to success in spite of some bad decisions that would be bad to repeat. Explicit identification of high-performing trajectories—those that exemplify generalizable reasoning patterns, as opposed to trajectories containing incorrect reasoning or actions (see Appendix for further analysis)—can improve efficiency compared to wholesale database removal.

This observation motivates *Exemplar-Curation*: identifying and selecting individual high-quality exemplars across multiple database instances based on their empirical utility as in-context examples. This approach parallels value-function learning in reinforcement learning [35], where we estimate the ‘value’ of each trajectory based on its contribution to successful outcomes.

We introduce a retrieval-weighted quality metric analogous to a value function to quantify each trajectory’s contribution to successful outcomes:

$$Q(\tau) = \frac{\sum_{i \in \mathcal{R}(\tau)} o_i \cdot f_i(\tau)}{\sum_{i \in \mathcal{R}(\tau)} f_i(\tau)} \quad (1)$$

where $\mathcal{R}(\tau)$ is the set of tasks for which trajectory τ was retrieved, o_i is the binary outcome of task i , and $f_i(\tau)$ is the retrieval frequency during task i .

This value metric measures how often a trajectory is associated with successful outcomes when retrieved as an in-context example. It weights outcomes by retrieval frequency, prioritizing trajectories frequently retrieved during successful completions while penalizing those associated with failures.

Algorithm 3 outlines our exemplar-level curation approach. For each task in the training set, we identify all successful trajectories across database instances and select only the exemplar with the highest value according to the metric. This approach constructs a composite database containing only the most effective exemplars as measured by their empirical contribution to successful outcomes on subsequent tasks.

Algorithm 3 Database Construction from Top Exemplars for +Exemplar-Curation

```
1: procedure SELECTEXEMPLARS( $\{\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_N\}, T_{train}$ )
2:    $\mathcal{D}_{composite} \leftarrow \emptyset$ 
3:   Compute quality metric  $Q(\tau)$  for each trajectory  $\tau \in \bigcup_{i=1}^N \mathcal{D}_i$ 
4:   for each task  $t \in T_{train}$  do
5:      $T_t \leftarrow \{\text{successful trajectories for task } t \text{ across all databases}\}$ 
6:     if  $T_t$  is not empty then
7:       Select top-1 trajectory from  $T_t$  by quality metric  $Q$ 
8:       Add selected trajectory to  $\mathcal{D}_{composite}$ 
9:   return  $\mathcal{D}_{composite}$ 
```

5.4 Train-Time vs Test-Time LLM Costs

The curation methods (+DB-Curation and +Exemplar-Curation) require N times more LLM inference calls during training compared to Traj-Bootstrap, as they maintain N parallel database instances. However, they do not require additional training tasks - all methods use the same task distribution and quantity. At test time, all three methods have identical computational costs—Algorithm 1 is simply provided with a different database \mathcal{D} for each method. This contrasts with approaches that scale the number of LLM calls per test-time task in order to improve performance [25, 26, 27, 28]. Our methods shift computational burden to training while maintaining efficient inference, a property our in-context methods share with fine-tuning methods.

6 Experiments

We evaluate our database construction methods through experiments addressing three key questions:

- Database scaling: How does task success rate scale with increasing database size?
- Improving database construction: How much do population-based training and exemplar-level curation improve task success rate?
- Overall effectiveness: How do our approaches compare to alternative approaches leveraging task-specific domain knowledge or hierarchical algorithms?

6.1 Experimental Setup

6.1.1 Benchmark Tasks

We evaluate our methods on three benchmarks:

- **ALFWorld** [36]: A text-based environment for navigation and object manipulation through textual commands, requiring exploration and sequential reasoning.
- **InterCode-SQL** [37]: An interactive coding environment where agents generate SQL queries to answer a user question, requiring understanding of database structures and query semantics.
- **Wordcraft** [38]: A simplified adaptation of Little Alchemy, where agents combine elements to create new ones through multi-step processes, requiring compositional reasoning.

These benchmarks were selected because they: (1) provide large enough task pools to support meaningful train/test splits, (2) represent diverse reasoning challenges relevant to sequential decision-making, and (3) have been used in prior work, enabling direct comparisons with existing methods.

Method	ALFWorld	InterCode-SQL	Wordcraft
Fixed-DB	0.73 \pm 0.02	0.75 \pm 0.01	0.55 \pm 0.03
Traj-Bootstrap	0.89 \pm 0.01	0.79 \pm 0.01	0.64 \pm 0.03
+DB-Curation	0.91 \pm 0.01	0.78 \pm 0.01	0.64 \pm 0.01
+Exemplar-Curation	0.90 \pm 0.02	0.81 \pm 0.01	0.72 \pm 0.02
+DB+Exemplar-Curation	0.93 \pm 0.03	0.82 \pm 0.01	0.69 \pm 0.01

Table 1: **Average success rate of our methods: self-collected trajectories provide the largest boosts in task success rate.** Traj-Bootstrap outperforms Fixed-DB across all three benchmarks. The combination of +DB-Curation and +Exemplar-Curation provides the best performance on both ALFWorld and InterCode-SQL. +Exemplar-Curation provides the best performance on Wordcraft.

6.1.2 Methods Compared

Our methods include:

- **Fixed-DB:** The baseline agent as described in Sec. 2, with a fixed database of human-provided initial examples and no database growth.
- **Traj-Bootstrap:** The simple progressive accumulation approach from Sec. 5.1.
- **Traj-Bootstrap+DB-Curation:** Our database-level trajectory curation from Alg. 2.
- **Traj-Bootstrap+Exemplar-Curation:** Our exemplar-level trajectory curation from Alg. 3.
- **Traj-Bootstrap+DB+Exemplar-Curation:** Applying both our database-level trajectory curation and propagation and our exemplar-level trajectory curation.

We compare these methods to two more advanced hierarchical designs. **Autoguide** [12] converts successful trajectories into explicit rules and retrieves the most contextually relevant rules, alongside low-level trajectories, at inference time. **AutoManual** [4] leverages hand-crafted task-specific observation and action spaces—see Appendix F.1 for details.

6.1.3 Implementation Details

Unless otherwise specified, we use GPT-4o-mini as our base LLM (temperature 0.1). For Fixed-DB and all Traj-Bootstrap agents, we retrieve the top- k most similar trajectories at each decision step ($k = 6$ for ALFWorld and InterCode-SQL, 10 for Wordcraft). We initialize each database with a small human-provided example set (18 for ALFWorld, 10 for InterCode-SQL, 4 for Wordcraft). With +DB-Curation, we maintain five database instances with curation every time the database size is doubled, starting with a minimum size of ten trajectories. We report success rates averaged over five random seeds. By default, we report success rate given the database at the end of the training process.

6.2 Traj-Bootstrap Results

Traj-Bootstrap performance improves with more training tasks Tab. 1 presents the final success rate metrics for our database construction methods. The performance of Traj-Bootstrap generally improves with increases in the number of training tasks attempted (Fig. 2) Performance continues to improve with more training tasks across all benchmarks, but exhibits diminishing returns—most gains occur within the first 25% of added training tasks. This efficiency decline occurs because each new example is retrieved less frequently as the database grows, influencing fewer generations, a pattern consistent with findings from Bertsch et al. [22] and Agarwal et al. [8]. As mentioned in Sec. 5, we observe performance variability across trials and within individual trials. Cross-trial variance indicates that some trials produce higher-performing databases when solving identical tasks. Within-trial fluctuations show that certain added trajectories can degrade performance.

Method	LLM(s)	Num Training Tasks	ALFWorld
Autoguide [12]	gpt-3.5-turbo + gpt-4-turbo	100	0.79*
Automanual [4]	gpt-4o-mini gpt-4-turbo + gpt-4o-mini	36	0.72 \pm 0.01
		36	0.91 \pm 0.01
Fixed-DB	gpt-4o-mini gpt-4o	0	0.73 \pm 0.05
		0	0.88 \pm 0.02
Traj-Bootstrap	gpt-4o-mini gpt-4o-mini	100	0.84 \pm 0.04
		3500	0.89 \pm 0.01
+DB-Curation	gpt-4o-mini gpt-4o-mini	100	0.86 \pm 0.02
		3500	0.91 \pm 0.01
+Exemplar-Curation	gpt-4o-mini gpt-4o-mini	100	0.86 \pm 0.03
		3500	0.90 \pm 0.02
+DB-Curation +Exemplar-Curation	gpt-4o-mini gpt-4o-mini	100	0.81 \pm 0.02
		3500	0.93\pm0.03

Table 2: **Comparison of agent success rates on ALFWorld: contextualizing the performance of Traj-Bootstrap.** The 15-point boost in average success rate from database construction via Traj-Bootstrap is similar to that achieved from upgrading Fixed-DB from gpt-4o-mini to gpt-4o. The performance of Traj-Bootstrap+DB+Exemplar-Curation exceeds Automanual [4], even though Automanual utilizes hand-designed observation and action spaces and a better LLM (gpt-4-turbo+gpt-4o-mini). * indicates results reported from original papers.

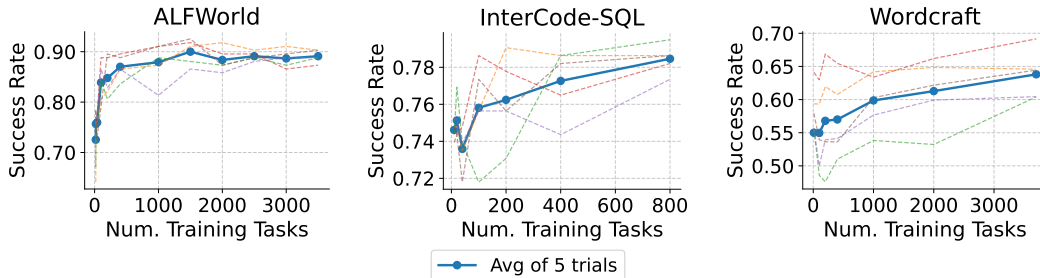


Figure 2: **Traj-Bootstrap results: success rate improves with increasing training tasks on all three benchmarks.** Individual trials (5) shown as dashed lines. All benchmarks exhibit diminishing returns as the database size increases. Trials show substantial performance variability, both within individual trials and across different trials.

+DB-Curation boosts performance on ALFWorld Fig. 3 illustrates how +DB-Curation can improve upon Traj-Bootstrap’s performance, despite exhibiting occasional performance dips at smaller database sizes. These dips result from inaccurate estimates (due to low sample count) of database quality early in the process—introducing noise into the curation process.

+Exemplar-Curation boosts performance on InterCode-SQL and Wordcraft As seen in Tab. 1, +Exemplar-Curation yields improvements in final task success rates on InterCode-SQL and Wordcraft, and also boosts success rate at intermediate database sizes for both InterCode-SQL and Wordcraft (Fig.3). To further highlight the impact of our exemplar-level curation metric, Fig. 4 compares databases built from the ‘best’ trajectories that are the most empirically effective in-context examples versus the least effective trajectories, as determined by Equation 1 in Sec. 5.3. The ‘best’ curve is identical to +Exemplar-Curation, while the ‘worst’ curve selects the bottom-1 trajectory instead of top-1 in Alg. 3, line 7. Using the database of high-quality examples yields a higher success rate across all database sizes for ALFWorld and Wordcraft, and for smaller database sizes for InterCode-SQL.

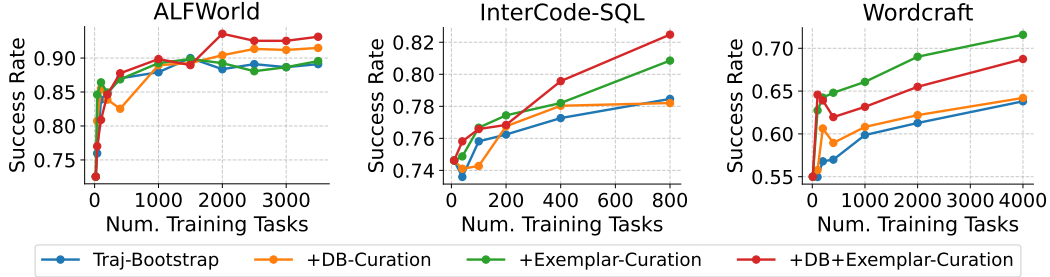


Figure 3: **Success rate comparison for Traj-Bootstrap and its variants (+DB-Curation, +Exemplar-Curation, +DB+Exemplar-Curation).** +DB-Curation enhances final success rate only on ALFWorld, but improves success rate for smaller DB sizes on all benchmarks. +Exemplar-Curation delivers success rate gains on both InterCode-SQL and Wordcraft. The combination of both enhancements delivers the largest gains on both ALFWorld and InterCode-SQL.

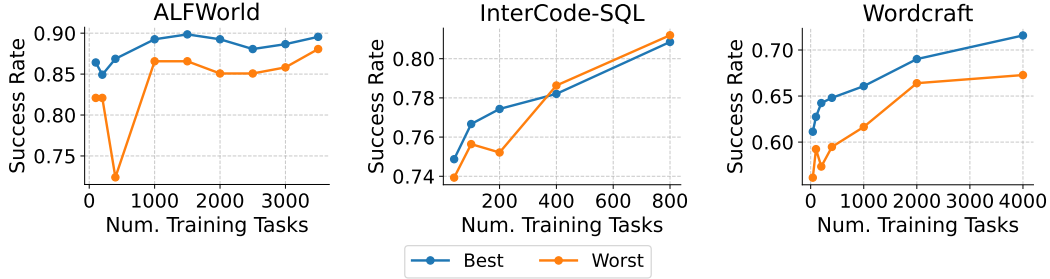


Figure 4: **The ‘best’ bootstrapped trajectories compared to the ‘worst’.** Databases constructed from the highest-quality successful trajectory per task, as measured by Eq. 1, outperform databases built from the lowest-quality successful trajectories on both ALFWorld and Wordcraft. The ‘best’ curve is identical to +Exemplar-Curation, while the ‘worst’ curve selects the bottom-1 trajectory instead of top-1 in Alg. 3, line 7.

+DB+Exemplar-Curation achieves best performance on ALFWorld and InterCode-SQL Fig. 3 and Tab. 1 highlight that +DB-Curation and +Exemplar-Curation can be complementary, as the combined +DB+Exemplar-Curation achieves the best final task success rates on both ALFWorld and InterCode-SQL (0.93 and 0.82 respectively). Note that on Wordcraft, +DB-Curation fails to provide boosts whether or not +Exemplar-Curation is used—Traj-Bootstrap and +DB-Curation perform identically (0.64), and +Exemplar-Curation (0.71) outperforms +DB+Exemplar-Curation (0.67).

6.3 Contextualizing performance boosts from Traj-Bootstrap

To contextualize the improvements achieved by Traj-Bootstrap, we compare with several alternative strategies: test-time sampling, using a better LLM, task-specific strategies, and hierarchical strategies.

Comparison with test-time scaling Our trajectory bootstrapping approaches achieve success rate improvements equivalent to scaling test-time compute by making multiple task attempts—an advantage in scenarios where multiple attempts are impractical or when success verification is unavailable at test time. Furthermore, our approaches provide these benefits without requiring any modifications to the test-time inference process. To demonstrate the magnitude of this benefit, we compare our methods to the alternative strategy of making multiple attempts at each test task with the Fixed-DB baseline and selecting the best outcome [25, 26, 27]. Tab. 3 reports the pass@k metrics for Fixed-DB across all three benchmarks, representing the probability of at least one successful attempt when making k independent attempts at each task. Using only a single attempt per task, Traj-Bootstrap approach achieves success rate comparable to Fixed-DB pass@2 or pass@3 on all three benchmarks. +DB-Curation and/or +Exemplar-Curation perform nearly on

Method	ALFWorld	Intercode-SQL	Wordcraft
Traj-Bootstrap	0.89 \pm 0.01	0.79 \pm 0.01	0.64 \pm 0.03
+DB+Exemplar-Curation	0.93 \pm 0.03	0.82 \pm 0.01	0.69 \pm 0.01
Fixed-DB@1	0.73 \pm 0.03	0.75 \pm 0.01	0.55 \pm 0.03
Fixed-DB@2	0.87 \pm 0.02	0.78 \pm 0.03	0.62 \pm 0.02
Fixed-DB@3	0.92 \pm 0.02	0.80 \pm 0.03	0.64 \pm 0.02
Fixed-DB@4	0.94 \pm 0.02	0.81 \pm 0.02	0.66 \pm 0.02
Fixed-DB@5	0.96 \pm 0.02	0.82 \pm 0.03	0.72 \pm 0.02

Table 3: **Pass@k of Fixed-DB on all benchmarks.** On all benchmarks, using only a single test-time attempt per task, Traj-Bootstrap achieves success rates between that of Fixed-DB at pass@2 and at pass@3. +DB+Exemplar-Curation achieves success rates between pass@3 and pass@5.

the level of Fixed-DB pass@4 on ALFWorld, pass@5 on InterCode-SQL, and pass@5 for Wordcraft.

Comparison with model upgrades On ALFWorld, after 3500 training tasks Traj-Bootstrap yields a 20-point success rate boost over Fixed-DB, significantly outperforming the 15-point improvement gained by upgrading Fixed-DB to a more powerful LLM.

Comparison with task-specific strategies Tab. 2 shows that Traj-Bootstrap+DB+Exemplar-Curation using GPT-4o-mini achieves a success rate exceeds (0.93) that exceeds that of Automanual [4] configured to use a combination of GPT-4-turbo and GPT-4o-mini (0.91). Thus, our methods outperform an approach that uses a more powerful LLM and customized observation and action spaces. See App. G for a comparison to hand-crafted approaches on InterCode-SQL.

Comparison with hierarchical algorithms Given 100 training tasks, Autoguide, a hierarchical rule-learning approach, achieves a 0.79 success rate (using a combination of gpt-3.5-turbo + gpt-4-turbo). Given the same number of training tasks our best approach achieves significantly greater success rate (0.86) with gpt-4o-mini (Tab. 2). While this comparison employs different algorithms and LLMs, the performance of Traj-Bootstrap suggests that self-constructed databases of low-level trajectories can be competitive with hierarchical approaches.

6.4 Extending Traj-Bootstrap

Can we predict agent success? Beyond improving agent performance, we can also utilize our self-collected examples to implement useful agent diagnostics, such as predicting an agent’s success on novel tasks. On InterCode-SQL and Wordcraft, we train a calibrated Random Forest classifier of agent success based on task goal and initial observation embeddings. Classifier quality (measured via AUROC) improves with database size, reaching 0.77 for InterCode-SQL and 0.71 for Wordcraft with our largest databases. The predicted probabilities also closely match empirical success rates, indicating well-calibrated predictions. See App. D for details.

Can we use our self-collected databases for fine-tuning? We fine-tune GPT-4o-mini using trajectories from our best-performing Traj-Bootstrap+DB+Exemplar-Curation database for each benchmark. The resulting fine-tuned agents (ReAct-Finetune) outperform our in-context approach on ALFWorld (23-point vs. 20-point boost) and Wordcraft (19-point vs. 14-point), but perform worse on InterCode-SQL (4-point vs. 7-point). This suggests our self-collected examples are effective not only for in-context learning but also for creating fine-tuned agents. See App. C for details.

7 Discussion

The success of our approach reveals performance gains that stem primarily from accumulating successful examples, establishing a foundation for agent self-improvement where the

quantity and quality of accessible data rivals the importance of architectural complexity. This parallels trends in traditional deep learning, where data curation often yields substantial improvements. Our findings point to promising research directions that approach LLM agent enhancement from a data-centric perspective—advancing both strategic data collection methods (balancing exploration versus exploitation across diverse tasks) and refined filtering techniques to maximize performance.

Acknowledgments Thank you to Brennan Shacklett, Purvi Goel, Zander Majercik, William Wang, Bradley Brown, Jon Saad-Falcon, and William Mark for valuable discussions and feedback. Support for this project was provided by Roblox and Meta, and API credits were provided by OpenAI and together.ai.

References

- [1] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- [2] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [3] Jerry Wei, Jason Wei, Yi Tay, Dustin Tran, Albert Webson, Yifeng Lu, Xinyun Chen, Hanxiao Liu, Da Huang, Denny Zhou, et al. Larger language models do in-context learning differently. *arXiv preprint arXiv:2303.03846*, 2023.
- [4] Minghao Chen, Yihang Li, Yanting Yang, Shiyu Yu, Binbin Lin, and Xiaofei He. Automanual: Generating instruction manuals by llm agents via interactive environmental learning. *arXiv preprint arXiv:2405.16247*, 2024.
- [5] Ke Yang, Yao Liu, Sapana Chaudhary, Rasool Fakoor, Pratik Chaudhari, George Karypis, and Huzefa Rangwala. Agentoccam: A simple yet strong baseline for llm-based web agents. *arXiv preprint arXiv:2410.13825*, 2024.
- [6] Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? investigations with linear models. *arXiv preprint arXiv:2211.15661*, 2022.
- [7] Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In *International Conference on Machine Learning*, pages 35151–35174. PMLR, 2023.
- [8] Rishabh Agarwal, Avi Singh, Lei Zhang, Bernd Bohnet, Luis Rosias, Stephanie Chan, Biao Zhang, Ankesh Anand, Zaheer Abbas, Azade Nova, et al. Many-shot in-context learning. *Advances in Neural Information Processing Systems*, 37:76930–76966, 2024.
- [9] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*, 2023.
- [10] Tomoyuki Kagaya, Thong Jing Yuan, Yuxuan Lou, Jayashree Karlekar, Sugiri Pranata, Akira Kinose, Koki Oguri, Felix Wick, and Yang You. Rap: Retrieval-augmented planning with contextual memory for multimodal llm agents. *arXiv preprint arXiv:2402.03610*, 2024.
- [11] Ruiwen Zhou, Yingxuan Yang, Muning Wen, Ying Wen, Wenhao Wang, Chunling Xi, Guoqiang Xu, Yong Yu, and Weinan Zhang. Trad: Enhancing llm agents with step-wise thought retrieval and aligned decision. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 3–13, 2024.

-
- [12] Yao Fu, Dong-Ki Kim, Jaekyeom Kim, Sungryull Sohn, Lajanugen Logeswaran, Kyunghoon Bae, and Honglak Lee. Autoguide: Automated generation and selection of context-aware guidelines for large language model agents. *arXiv preprint arXiv:2403.08978*, 2024.
 - [13] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*, 2021.
 - [14] Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770*, 2023.
 - [15] Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M Sadler, Wei-Lun Chao, and Yu Su. Llm-planner: Few-shot grounded planning for embodied agents with large language models. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 2998–3009, 2023.
 - [16] Hongliang He, Wenlin Yao, Kaixin Ma, Wenhao Yu, Yong Dai, Hongming Zhang, Zhenzhong Lan, and Dong Yu. Webvoyager: Building an end-to-end web agent with large multimodal models. *arXiv preprint arXiv:2401.13919*, 2024.
 - [17] Andrew Zhao, Daniel Huang, Quentin Xu, Matthieu Lin, Yong-Jin Liu, and Gao Huang. Expel: Llm agents are experiential learners. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 19632–19642, 2024.
 - [18] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
 - [19] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36:53728–53741, 2023.
 - [20] Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024.
 - [21] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
 - [22] Amanda Bertsch, Maor Ivgi, Uri Alon, Jonathan Berant, Matthew R Gormley, and Graham Neubig. In-context learning with long-context models: An in-depth exploration. *arXiv preprint arXiv:2405.00200*, 2024.
 - [23] Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vardhamanan, Saiful Haq, Ashutosh Sharma, Thomas T Joshi, Hanna Moazam, et al. Dspy: Compiling declarative language model calls into self-improving pipelines. *arXiv preprint arXiv:2310.03714*, 2023.
 - [24] Krista Opsahl-Ong, Michael J Ryan, Josh Purtell, David Broman, Christopher Potts, Matei Zaharia, and Omar Khattab. Optimizing instructions and demonstrations for multi-stage language model programs. *arXiv preprint arXiv:2406.11695*, 2024.
 - [25] Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V Le, Christopher Ré, and Azalia Mirhoseini. Large language monkeys: Scaling inference compute with repeated sampling. *arXiv preprint arXiv:2407.21787*, 2024.
 - [26] Evan Wang, Federico Cassano, Catherine Wu, Yunfeng Bai, Will Song, Vaskar Nath, Ziwen Han, Sean Hendryx, Summer Yue, and Hugh Zhang. Planning in natural language improves llm search for code generation. *arXiv preprint arXiv:2409.03733*, 2024.

-
- [27] Xiyao Wang, Zhengyuan Yang, Linjie Li, Hongjin Lu, Yuancheng Xu, Chung-Ching Lin, Kevin Lin, Furong Huang, and Lijuan Wang. Scaling inference-time search with vision value model for improved visual comprehension. *arXiv preprint arXiv:2412.03704*, 2024.
- [28] Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 36:8634–8652, 2023.
- [29] Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291*, 2023.
- [30] Jon Saad-Falcon, Adrian Gamarra Lafuente, Shlok Natarajan, Nahum Maru, Hristo Todorov, Etash Guha, E Kelly Buchanan, Mayee Chen, Neel Guha, Christopher Ré, et al. Archon: An architecture search framework for inference-time techniques. *arXiv preprint arXiv:2409.15254*, 2024.
- [31] Shengran Hu, Cong Lu, and Jeff Clune. Automated design of agentic systems. *arXiv preprint arXiv:2408.08435*, 2024.
- [32] Jiayi Zhang, Jinyu Xiang, Zhaoyang Yu, Fengwei Teng, Xionghui Chen, Jiaqi Chen, Mingchen Zhuge, Xin Cheng, Sirui Hong, Jinlin Wang, et al. Aflow: Automating agentic workflow generation. *arXiv preprint arXiv:2410.10762*, 2024.
- [33] Jan Peters and Stefan Schaal. Reinforcement learning by reward-weighted regression for operational space control. In *Proceedings of the 24th international conference on Machine learning*, pages 745–750, 2007.
- [34] Max Jaderberg, Valentin Dalibard, Simon Osindero, Wojciech M Czarnecki, Jeff Donahue, Ali Razavi, Oriol Vinyals, Tim Green, Iain Dunning, Karen Simonyan, et al. Population based training of neural networks. *arXiv preprint arXiv:1711.09846*, 2017.
- [35] Andrew G Barto. Reinforcement learning: An introduction. by richard’s sutton. *SIAM Rev*, 6(2):423, 2021.
- [36] Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew Hausknecht. Alfworld: Aligning text and embodied environments for interactive learning. *arXiv preprint arXiv:2010.03768*, 2020.
- [37] John Yang, Akshara Prabhakar, Karthik Narasimhan, and Shunyu Yao. Intercode: Standardizing and benchmarking interactive coding with execution feedback. *Advances in Neural Information Processing Systems*, 36:23826–23854, 2023.
- [38] Minqi Jiang, Jelena Luketina, Nantas Nardelli, Pasquale Minervini, Philip HS Torr, Shimon Whiteson, and Tim Rocktäschel. Wordcraft: An environment for benchmarking commonsense agents. *arXiv preprint arXiv:2007.09185*, 2020.
- [39] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 11 2019. URL <http://arxiv.org/abs/1908.10084>.
- [40] Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. The faiss library. *arXiv preprint arXiv:2401.08281*, 2024.
- [41] Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable real-world web interaction with grounded language agents. *Advances in Neural Information Processing Systems*, 35:20744–20757, 2022.
- [42] Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W Cohen, Ruslan Salakhutdinov, and Christopher D Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. *arXiv preprint arXiv:1809.09600*, 2018.

A Key Agent Details

In Sec. 2, we establish an agent design that enables it to learn in-context from its own self-collected experiences. Here, we elaborate on a few key design decisions in our agent design:

- **Standardized prompts:** we use the same simple, task-agnostic prompt templates for all tasks, rather than writing new prompts per task. These prompts are in App. B. Alternate approaches incorporate domain-specific information into their prompts—we discuss these approaches in Appendices F.1 and G.
- **Two-level retrieval:** We retrieve trajectories at both trajectory level (for planning) and state level (for reasoning and action selection), enabling the agent to leverage both strategic patterns and situation-specific techniques. Database \mathcal{D} contains self-collected trajectories, and retrieval is performed at the trajectory level for the initial plan p , and in the state-level observation-reasoning-action loop for both r_t and a_t .
- **Multi-key retrieval:** All retrieval is performed by KNN, with similarity metric defined as the average of cosine similarities across the specified ‘key’ variables. For instance, in Line 3 of Alg. 1, we retrieve from \mathcal{D} using two keys: goal g and plan p . We return similar trajectories based off the average of the cosine similarities of goals and plans when comparing each trajectory to the current trajectory. When doing state-level retrieval (Lines 7 and 10), we additionally find the most similar states within the selected trajectories via state-level key o_t or r_t , then return a window of states around the most similar state. This is similar to the retrieval scheme in [10]. See detailed pseudocode for retrieval in Alg. 4.
- **Thought-based retrieval:** For the first step of a trajectory, we retrieve using the trajectory-level keys (g, p) as well as the current observation o_1 (Alg. 1, line 6)—but for all subsequent steps we use reasoning r_t as a key instead of observation o_t (Alg. 1, line 9). This approach, inspired by Zhou et al. [11], enables generalization across trajectories with similar reasoning, and similarity across natural-language r_t can be handled by generic embedding functions more easily than potentially bespoke observations o_t . By retrieving at every step, we aim to retrieve the most relevant trajectories for each decision.
- **Generic embedding mechanism:** Since g , p , and r_t are all natural-language strings, we employ standard embeddings (all-MiniLM-L6-v2 [39]) that generalize across domains without task-specific engineering.

B Additional Implementation Details

B.1 Prompt Templates

Across all benchmarks, we use standardized prompt templates for the core components of our retrieval-based ReAct agent. The same templates were used across all benchmarks with no task-specific modifications. These templates are intentionally minimalist, focusing on providing the necessary context and retrieved examples while avoiding task-specific prompt engineering.

The templates are included below. Across all templates, the in-context examples follow the format specified in the prompt itself (for plan, the in-context examples are of form “goal,plan”, etc):

Plan:

```
1 system_prompt: 'You are an expert at generating high-level plans of actions to
  ↳ achieve a goal.\n Here is your action space: {action_space}.\n Here are some
  ↳ examples of goal,plan from episodes that successfully achieved similar goals:
  ↳ {examples}'
2 user_prompt: 'goal: {goal}\n plan: '
```

Algorithm 4 Multi-key Retrieval

```
1: procedure MULTIKEYRETRIEVAL( $\mathcal{D}$ , traj_keys, state_key, query,  $k$ , window_size)
2:   similarities  $\leftarrow []$ 
3:   for each trajectory  $\tau$  in  $\mathcal{D}$  do
4:     sim  $\leftarrow 0$ 
5:     for each key in traj_keys do
6:       sim  $\leftarrow$  sim + CosineSimilarity(query[key],  $\tau$ [key])
7:     sim  $\leftarrow$  sim / |traj_keys| ▷ Average similarity across trajectory keys
8:     similarities.append(sim,  $\tau$ )
9:   similar_trajectories  $\leftarrow$  TopK(similarities,  $k$ )
10:  if state_key is not None then ▷ State-level retrieval with window
11:    windowed_results  $\leftarrow []$ 
12:    for each trajectory  $\tau$  in similar_trajectories do
13:      state_similarities  $\leftarrow []$ 
14:      for each state  $s$  in  $\tau$ .states do
15:        state_sim  $\leftarrow$  CosineSimilarity(query[state_key],  $s$ [state_key])
16:        state_similarities.append(state_sim,  $s$ , index( $s$ ))
17:       $\_,$  most_similar_state, idx  $\leftarrow$  Max(state_similarities)
18:      window_start  $\leftarrow$  max(0, idx -  $\lfloor$  window_size / 2  $\rfloor$ )
19:      window_end  $\leftarrow$  min(| $\tau$ .states|, idx +  $\lceil$  window_size / 2  $\rceil$ )
20:      windowed_results.append( $\tau$ .states[window_start : window_end])
21:  return windowed_results
22: else
23:  return similar_trajectories
```

Reason:

```
1 system_prompt: 'You are an expert at reasoning about the most appropriate action to
  ↳ take towards achieving a goal.\n Here is your action space: {action_space}.\n
  ↳ Here are some examples of goal,plan,observation,reasoning,action from episodes
  ↳ that successfully achieved similar goals: {examples}'
2 user_prompt: 'goal: {goal}\n plan: {plan}\n trajectory: {trajectory}\n reasoning: '
```

Act:

```
1 system_prompt: 'You are an agent in an environment. Given the current observation,
  ↳ you must select an action to take towards achieving the goal: {self.goal}.\n
  ↳ Here is your action space: {action_space}.\n Here are some examples of
  ↳ goal,plan,observation,reasoning,action from episodes that successfully achieved
  ↳ similar goals: {examples}'
2 user_prompt: 'goal: {goal}\n plan: {plan}\n trajectory: {trajectory}\n action: '
```

B.2 Retrieval Implementation

For all retrieval steps, we implement hybrid search across all the desired retrieval keys—ex. goal, plan, observation, reasoning. We return the top- k examples by averaged distance across each of the keys. We implement a sliding window approach for state-level retrieval to enhance contextual relevance—we include the surrounding context (preceding and following states) up to a window of 5 steps to provide coherent episode fragments.

The retrieval mechanism is implemented using FAISS [40] for efficient similarity search as the database grows. We use exact nearest neighbor search.

B.3 Population-Based Training Details

Our database-level curation approach maintains a population of 5 database instances. Each instance is initialized with the same set of human-provided exemplars. The population

undergoes curation every time the database size doubles, and performance is evaluated on the tasks attempted since the previous doubling.

The replacement strategy follows standard population-based training practices: the bottom 20% of databases (based on validation performance) are replaced with copies of the top 20%.

B.4 Quality Metric Computation

For exemplar-level curation, we track the retrieval patterns of each trajectory throughout the training process. For each task, we record: 1. Which trajectories were retrieved 2. How many times each trajectory was retrieved during the solution process 3. Whether the task was successfully completed

After completing all training tasks, we compute the quality metric $Q(\tau)$ for each trajectory τ as:

$$Q(\tau) = \frac{\sum_{i \in \mathcal{R}(\tau)} o_i \cdot f_i(\tau)}{\sum_{i \in \mathcal{R}(\tau)} f_i(\tau)} \quad (2)$$

where $\mathcal{R}(\tau)$ is the set of tasks for which trajectory τ was retrieved, $o_i \in \{0, 1\}$ is the outcome of task i , and $f_i(\tau)$ is the retrieval frequency of trajectory τ during task i .

To ensure statistical significance, we only compute the quality metric for trajectories that were retrieved for at least 3 different tasks. For trajectories with insufficient retrieval data, we assign a neutral quality score equal to the average success rate across all tasks.

B.5 Note on planning step

Following the convention from RAP [10], we omit the planning step on benchmarks with short trajectory length (Intercode-SQL, Wordcraft). This planning step is valuable for maintaining long-horizon coherence on the ALFWorld benchmarks (30 steps), and is standard in prior ReAct-based agentic work [10, 17, 12], whether the planning step is explicitly separate from reasoning, or incorporated into the first reasoning step.

C Can Self-Collected Examples Improve a Fine-Tuned LLM Agent?

We have shown that self-collected databases improve the performance of in-context LLM agents. Here, we test whether the same data can also benefit fine-tuning.

Using the OpenAI fine-tuning API, we fine-tune GPT-4o-mini on each benchmark using data from our best-performing database construction method: Traj-Bootstrap+DB+Exemplar-Curation, collected over the full training set. We fine-tune using a simple ReAct-format prompt:

```

1 {
2   'system': 'You are a ReAct agent that helps users accomplish tasks. Given a goal,
   ↳ you will receive observations about the environment and respond with your
   ↳ reasoning and actions. For each observation, first think through the problem
   ↳ step by step (Thought), then decide on an action (Action). Your actions should
   ↳ be clear, concise, and directly executable in the environment.',
3   'user': 'Goal: {goal} \n Initial observation: {observations[0]}',
4   'assistant': 'Thought: {reasoning[i]}\nAction: {action[i]}',
5   'user': 'Observation: {observations[i+1]}',
6   ...
7 }
```

We refer to the resulting fine-tuned model as ReAct-Finetune. To run the agent, we prompt it with a goal and initial observation, then alternate assistant messages (for reasoning and action) with user messages (for new observations).

Method	ALFWorld	InterCode-SQL	Wordcraft
Traj-Bootstrap+DB+Exemplar-Curation	0.93 \pm 0.03	0.82\pm0.01	0.69 \pm 0.01
ReAct-Finetune	0.96\pm0.01	0.79 \pm 0.01	0.74\pm0.01

Table 4: **Trained on the same data, fine-tuned agent ReAct-Finetune is competitive with our best in-context approach.** This suggests that our self-collected data is effective not only for in-context prompting but also for creating fine-tuned agents. All values are averages over 5 trials.

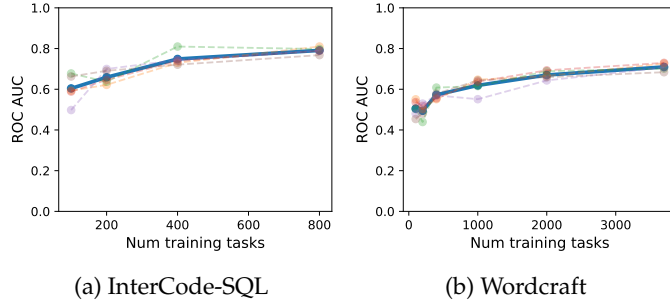


Figure 5: **AUROC of success prediction improves with more self-collected examples.** Performance continues to rise with increasing database size.

Tab. 4 shows that ReAct-Finetune slightly outperforms the in-context agent on ALFWorld (0.96 vs. 0.93) and Wordcraft (0.74 vs. 0.69), while performing slightly worse on InterCode-SQL (0.79 vs. 0.82). These results suggest that self-collected examples are effective not only for in-context prompting but also for creating competitive fine-tuned agents.

D Can We Predict Agent Success Rates

We have shown that increasing the number of self-collected examples improves agent performance. Here, we test whether the same examples can also predict performance on new tasks.

On InterCode-SQL and Wordcraft, task difficulty is partly observable from the goal g and initial observation o_1 . For InterCode-SQL, g is a natural-language query. For Wordcraft, g is the desired element and o_1 specifies available crafting elements. In contrast, ALFWorld task difficulty depends heavily on scene layout, which g and o_1 do not reveal. We therefore exclude ALFWorld from this analysis.

We use the same embedding model as in retrieval (all-MiniLM-L6-v2 [39]) to encode the concatenated string $[g; o_1]$. We train a calibrated Random Forest classifier to predict task success/failure, calibrating its outputs via 5-fold cross-validation with a learned sigmoid function. For each of 5 independent Traj-Bootstrap trials, we evaluate (1) the classifier’s AUROC on held-out tasks, and (2) its calibration.

As shown in Fig. 5, prediction performance improves as the database grows. For InterCode-SQL, AUROC rises from 0.60 (100 tasks) to 0.77 (800 tasks). Wordcraft shows a similar trend, improving from 0.50 (100 tasks) to 0.71 (4000 tasks). In both cases, predictive accuracy increases alongside task performance.

Fig. 6 shows the calibration of the final classifiers (trained on all available training tasks). Predicted success probabilities closely match observed success rates, indicating well-calibrated models.

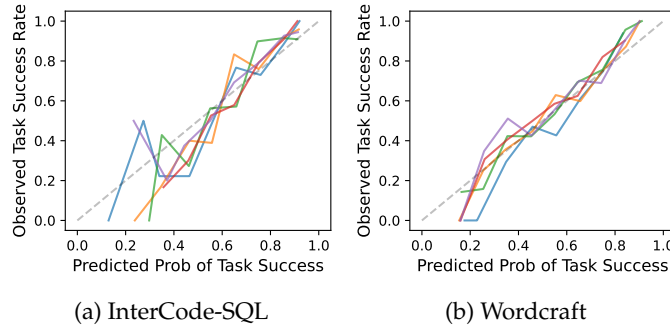


Figure 6: **Predicted probabilities are well-calibrated.** For both benchmarks, predicted and empirical success rates generally align.

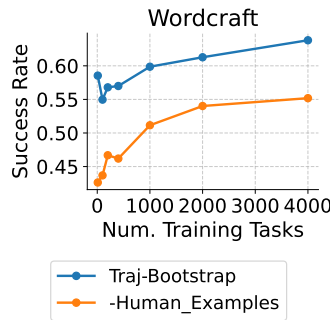


Figure 7: **Ablating the value of initial human-provided examples, Wordcraft.** Traj-Bootstrap, initialized by default with a database of 5 human-provided trajectories for Wordcraft, achieves better performance with these starting examples than when initialized from an empty database (-Human-Examples). Performance still scales with database size for -Human-Examples—but in this case fails to reach the performance achieved via 5 human-provided examples, even after self-collecting trajectories on 4000 training tasks.

E Is it Possible to Bootstrap an Agent Without Initial Hand-Crafted Examples?

Providing a small number of hand-crafted in-context examples is standard practice in the LLM agent literature [9, 17, 12, 4, 10]. However, what if we initialized Traj-Bootstrap with an empty database? In order to understand the value of the initial human-provided examples, we test Traj-Bootstrap with and without the initial human-provided examples on Wordcraft. We refer to the variant initialized with an empty database as -Human-Examples. Traj-Bootstrap, initialized by default with a database of 5 human-provided trajectories for Wordcraft, achieves better performance with these starting examples than when initialized from an empty database (-Human-Examples). Performance still scales with database size for -Human-Examples—but in this case fails to reach the performance achieved via 5 human-provided examples, even after self-collecting trajectories on 4000 training tasks. On at least this one task, the initial human-provided trajectories shaped the reasoning and action patterns of the agent in a way that boosted the continual database construction process. We leave exploration of hand-crafting these in-context examples to future work.

F Key Details of Prior Agentic Approaches

F.1 How does Automanual Leverage Hand-Crafted Information

Rather than learning from self-collected examples, an alternate approach to agent construction is to leverage practitioner domain knowledge. Beyond implementing both a

hierarchical learning system and code-based action spaces, Automanual [4] incorporates domain knowledge about the ALFWorld task into multiple components of the algorithm. In this section we include some code from the official Automanual GitHub to illustrate.

Observation spaces : Automanual uses a modified observation space that enhances the ALFWorld string by adding two critical pieces of information: 1) The current location of the agent, 2) What the agent is currently holding. Both of these pieces of information typically have to be deduced from the trajectory of previous observations and actions, but Automanual tracks them explicitly:

```
1 if "Nothing happens" not in observation:
2     self.last_obs = observation
3     if "go to" in script:
4         self.cur_loc = re.search(r'go to (\S+)', script).group(1)
5         self.cur_loc_info = observation
6     if "open" in script or "close" in script:
7         self.cur_loc_info = observation
8     if "take" in script:
9         self.holding = re.search(r"(?<=take\s)(.*?)(?=\sfrom)", script).group(1)
10        self.cur_loc_info = ""
11    if "put" in script:
12        self.holding = "nothing"
13        self.cur_loc_info = ""
14 elif "go to" in script:
15     loc = re.search(r'go to (\S+)', script).group(1)
16     if loc == self.cur_loc:
17         observation = self.cur_loc_info
18 observation += f" You are at {self.cur_loc} and holding {self.holding}."
```

Action spaces : in ALFWorld, any task typically involves three main components: 1) Searching for an object, 2) Performing an action with the object (heating, cooling, cleaning, etc.), 3) Placing the object somewhere. Automanual significantly simplifies both the search and placement operations by providing multi-action helper functions within its code-based action space:

```
1 # Define a helper method to find object that is needed
2 def find_object(agent, recep_to_check, object_name):
3     for receptacle in recep_to_check:
4         observation = agent.go_to(receptacle)
5         # Check if we need to open the receptacle. If we do, open it.
6         if 'closed' in observation:
7             observation = agent.open(receptacle)
8         # Check if the object is in/on the receptacle.
9         if object_name in observation:
10            object_ids = get_object_with_id(observation, object_name)
11            return object_ids, receptacle
12    return None, None
13
14 # Define a helper method to put object in/on the target receptacle
15 def go_to_put_object(agent, target_receptacle, object_id):
16     observation = agent.go_to(target_receptacle)
17     # check if target_receptacle is closed. If so, open it.
18     if 'closed' in observation:
19         observation = agent.open(target_receptacle)
20     observation = agent.put_in_or_on(object_id, target_receptacle)
21     return observation
```


Method	Intercode-SQL Success Rate
GameSQL	0.73
GameSQL+Cheat	0.84
Fixed-DB	0.74
Traj-Bootstrap	0.79
+DB-Curation	0.78
+Exemplar-Curation	0.81
+DB+Exemplar-Curation	0.82

Table 5: **Comparison of agent success rates on InterCode-SQL: contextualizing the performance of Traj-Bootstrap.** Without cheats, the hand-crafted GameSQL agent (0.73) performs comparably to Fixed-DB (0.74). With handicap access to the database schema, GameSQL+Cheat (0.84) slightly outperforms +DB+Exemplar-Curation (0.82). The boost from our database-construction techniques nearly matches the boost from providing the GameSQL agent with access to privileged database schema information.

F.2 A note on training and test sets

The distinction between how different techniques leverage data is crucial in understanding the generalization capabilities of LLM agents. We can categorize existing approaches based on how they treat training and test data:

Single-Task Optimization : Some approaches focus exclusively on improving performance on a single task instance without concern for generalization. For example, Shinn et al. [28] leverages feedback from failed attempts to incrementally improve performance on the same task. Similarly, search methods [25, 27] expand the solution space for a specific problem instance. While these approaches can solve individual tasks, they don’t transfer knowledge across different problems, essentially ‘overfitting’ to a single instance.

Mixed Train-Test Evaluation : Some recent work blurs training and test boundaries. For instance, RAP [10] makes multiple passes over the same dataset, allowing the system to ‘learn’ from some test examples before evaluating on others within the same set. This approach does not assess true generalization capability, as the model has indirect exposure to the test distribution during its learning phase.

Full Train-Test Separation : Several papers maintain a clear separation between training and test data: 1) ExpeL [17] extracts general rules from a training set of trajectories and applies them to entirely separate test tasks, 2) AutoGuide [12] generates contextual guidelines from training experiences that are evaluated on distinct test scenarios. 3) AutoManual [4] constructs hierarchical ‘manuals’ from training interactions that are then applied to novel test tasks.

Our approach similarly ensures that trajectories used for database construction come exclusively from designated training tasks, with evaluation conducted on a separate set of test tasks never seen during the database construction phase. This separation is essential for validating that the knowledge captured by the agent generalizes to new problems rather than memorizing specific solutions.

G Comparison to Hand-Crafted InterCode-SQL Agent

We further contextualize the performance of Traj-Bootstrap by comparing to two hand-crafted agents on InterCode-SQL. The Intercode-SQL paper [37] provides a hand-crafted agent, GameSQL to solve the task, and optionally provides the agent with a ‘handicap’—giving the agent information on all relevant parts of the database schema. We denote the assisted version as GameSQL+Cheat. Neither agent provides in-context examples, and both share a bespoke, hand-crafted prompt (see App. J).

As seen in Tab. 5, Fixed-DB performs similarly to GameSQL (0.74 vs 0.73), and the performance of our best method, +DB+Exemplar-Curation, approaches the performance of GameSQL+Cheat (0.82 vs 0.84). Therefore, our database-construction techniques lift the performance of a generic ReAct-style agent nearly as much as the lift provided to the hand-crafted agent via ‘handicap’ access to the database schema.

H Benchmark Details

H.1 ALFWorld

ALFWorld [36] is a text-based environment that aligns with embodied tasks, allowing agents to navigate and manipulate objects through textual commands. We use the standard ALFWorld benchmark consisting of 3500 training tasks and 134 out-of-distribution test tasks across 6 task categories:

- Pick & Place: Find and move objects to specified locations
- Clean & Place: Find, clean, and place objects
- Heat & Place: Find, heat, and place objects
- Cool & Place: Find, cool, and place objects
- Pick Two & Place: Find and move two objects to a specified location
- Look at Object: Find an object and examine it under light

Following [10], for the ALFWorld benchmark we perform similarity search over task categories in addition to the other retrieval keys (goal, plan, observation, action). We do this to follow the convention in this prior work.

For our initial human-provided exemplars, we used the 18 successful trajectories (3 per task category) provided by Zhao et al. [17]. These trajectories were used to initialize all database instances.

The success criteria for ALFWorld tasks are defined by the environment and require the agent to satisfy all conditions specified in the goal. For example, in a ‘Heat & Place’ task, the agent must find the target object, place it in the microwave, turn on the microwave, and finally place the heated object at the specified destination. Both Autoguide and Automanual allow 50 actions for task completion—but choosing to employ “reasoning” counts as an action. Since we force our agent to reason at every step, we allow our agents (Fixed-DB, Traj-Bootstrap and variants) only 30 steps for task completion (on Autoguide and Automanual, the agent does not reason in practice at most steps ex. in a search procedure).

For this benchmark, we do not provide an action space string to the LLM, relying purely on the in-context examples to communicate the action space.

H.2 InterCode-SQL

InterCode-SQL [37] is an interactive coding environment for evaluating language agents’ SQL programming abilities. We use a subset of the InterCode benchmark focusing on SQL query generation, built upon the Spider SQL dataset. Of the 1034 tasks in the dataset, we randomly assign 800 tasks to train and the remaining 234 tasks to test.

Each task provides a natural language query request. The agent must generate a syntactically correct SQL query that retrieves the requested information. The agent must first execute queries to understand the database schema. The environment provides feedback on syntax errors and execution results, but the agent is only allowed to submit a solution once.

The success criteria for InterCode-SQL tasks require the agent to submit a solution query within 10 steps. The environment executes the query and compares the results against a ground-truth reference.

For our initial human-provided exemplars, we collected 10 human-created trajectories for 10 randomly-selected training tasks. These trajectories were used to initialize all database

instances. For all solved trajectories, we append the solution query to the goal string—since the goal of the task is to ‘discover’ this query through interacting with the SQL database.

We used the following action space string for InterCode-SQL:

```
Your action space is outputting valid mysql commands to solve the sql task.
You will be evaluated on the Latest Standard Output.
If you believe the latest observation is the final answer, you can complete the task by
    running 'submit' by itself.
You have 10 iterations to solve the task.
Follow the syntax and logical flow from the provided examples exactly.
```

H.3 Wordcraft

Wordcraft [38] is a simplified adaptation of the game Little Alchemy, where agents must combine elements to create new elements through multi-step processes. We randomly select 4000 training tasks and 500 test tasks from the subset of tasks requiring up to 2 steps to solve, with the train-test split separating the tasks into disjoint sets of goal elements.

The agent starts with a set of elements and must discover combinations that creates a particular target element specified in the goal. The environment provides feedback on successful combinations and updates the available elements accordingly.

The success criteria for Wordcraft tasks require the agent to create the target element within 4 steps, while the minimum solution length is up to 2 steps.

For our initial human-provided exemplars, we collected 4 human-annotated trajectories from randomly-selecting training tasks. These trajectories were used to initialize all database instances. We collected fewer initial trajectories for Wordcraft than for InterCode-SQL (4 vs 10) since Wordcraft is a slightly simpler task, requiring up to 4 steps for task completion while InterCode-SQL requires up to 10.

We used the following action space string for Wordcraft:

```
Output strings with the names of the two entities we would like to combine in this step.
```

H.4 Note on Benchmark Selection

We selected three sequential decision-making benchmarks that cover different reasoning challenges—**ALFWorld** [36] tests text-based navigation and object manipulation, **InterCode-SQL** [37] tests interactive code generation, and **Wordcraft** [?] tests compositional reasoning.

While prior works [17, 12] test on WebShop [41], we encountered bugs in generating achievable goals on the full benchmark (confirmed by <https://github.com/princeton-nlp/WebShop/issues/43>) and identified tasks with incorrect rewards.

We excluded QA benchmarks (HotPotQA [42], etc.) because performance depends on information retriever quality and LLM self-evaluation efficacy, two factors that would confound our study of LLM Self-Improvement. We plan to test our algorithms on QA benchmarks in future work.

I Computational Resources

All experiments were conducted using the following computational resources:

- 1 NVIDIA A5000 GPU (24GB memory) for embedding computation
- 64GB RAM

The majority of computation was spent on OpenAI API calls for the LLM-based decision-making. Database operations including embedding computation, storage, and retrieval accounted for less than 5% of the total computation time.

For embedding computations, we used all-MiniLM-L6-v2 [39].

For LLM inference, we used the OpenAI API for GPT-4o-mini, which required approximately:

- 2,000,000 API calls for ALFWorld
- 200,000 API calls for InterCode-SQL
- 500,000 API calls for Wordcraft

The total cost of API usage was approximately \$3,000 USD.

J GameSQL Prompt

Yang et al. [37] write this hand-crafted prompt for the GameSQL agent:

```
1 '''
2 {self.language}Env` is a multi-turn game that tests your ability to write
3 a {self.language} command that produces an output corresponding to a natural
4   ↳ language query.
5
6 ## GAME DESCRIPTION
7 At the start of this game, you are given a natural language query describing some
8 desired output (i.e. "Find the first name of a student who have both cat and dog
9   ↳ pets").
10 Aside from the natural language query, you have no information about the tables you
11   ↳ have access to.
12
13 The game will be played in a series of turns. Each turn, you can submit a
14   ↳ {self.language} command.
15 You will then get a response detailing the output of your {self.language} query
16   ↳ along with a reward
17 that tells you how close your {self.language} command is to the correct answer.
18
19 The goal of this game is to write a {self.language} command that gets a reward of 1.
20   ↳ The game will automatically
21 terminate once you get a reward of 1.
22
23 ## INPUT DESCRIPTION
24 Each turn, you can submit a {self.language} command. Your {self.language} command
25   ↳ should be formatted as follows:
26
27 ```{self.language}
28 Your {self.language} code here
29 ```
30
31 Your {self.language} command can help you do one of two things:
32 1. Learn more about the tables you have access to
33 2. Execute {self.language} commands based on these tables to generate the correct
34   ↳ output.
35
36 ## OUTPUT DESCRIPTION
37 Given your {self.language} command input, `{self.language}Env` will then give back
38   ↳ output formatted as follows:
39
40 Output: <string>
41 Reward: <decimal value between 0 and 1>
42
43 The output is a string displaying the result from executing your {self.language}
44   ↳ query.
45 The reward is a decimal value between 0 and 1.
46
47 ## REWARD DESCRIPTION
```

```
38 The reward should be interpreted as a ratio. It tells you how many rows your
   ↳ {self.language}
39 command outputted correctly compared to the correct answer.
40
41 ## RULES
42 1. Do NOT ask questions. Your commands are fed directly into a SQL compiler.
43
44 ## STRATEGY
45 You are free to play as many turns of the game as you'd like to inspect tables
46 and develop your {self.language} command.
47
48 The best strategy for this game is to first write {self.language} commands that
   ↳ help you learn
49 about the tables that you have access to. For instance, in a SQL environment, you
   ↳ might use `SHOW TABLES`
50 and `DESC <table name>` to learn more about the tables you have access to.
51
52 Once you have a good understanding of the tables, you should then write
   ↳ {self.language} commands
53 that would answer the natural language query using the tables you have access to.
54 '''
```