# Aegis: Taxonomy and Optimizations for Overcoming Agent-Environment Failures in LLM Agents

Kevin Song
University of Toronto, Vector Institute
Toronto, Canada
xinyang.song@utoronto.ca

Anand Jayarajan
University of Toronto, Vector Institute
Toronto, Canada
anandj@cs.toronto.edu

Yaoyao Ding
University of Toronto, Vector Institute
Toronto, Canada
yaoyao@cs.toronto.edu

Qidong Su
University of Toronto, Vector Institute
Toronto, Canada
qdsu@cs.toronto.edu

Zhanda Zhu
University of Toronto, Vector Institute
Toronto, Canada
zhandazhu@gmail.com

Sihang Liu
University of Waterloo
Waterloo, Canada
sihangliu@uwaterloo.ca

Gennady Pekhimenko
University of Toronto, Vector Institute
Toronto, Canada
pekhimenko@cs.toronto.edu

## Abstract

Large Language Models (LLMs) agents augmented with domain tools promise to autonomously execute complex tasks requiring human-level intelligence, such as customer service and digital assistance. However, their practical deployment is often limited by their low success rates under complex real-world environments. To tackle this, prior research has primarily focused on improving the agents themselves, such as developing strong agentic LLMs, while overlooking the role of the system environment in which the agent operates.

In this paper, we study a complementary direction: improving agent success rates by optimizing the system environment in which the agent operates. We collect 142 agent traces (3,656 turns of agent-environment interactions) across 5 state-of-the-art agentic benchmarks. By analyzing these agent failures, we propose a taxonomy for agent-environment interaction failures that includes 6 failure modes. Guided by these findings, we design Aegis, a set of targeted environment optimizations: 1) environment observability enhancement, 2) common computation offloading, and 3) speculative agentic actions. These techniques improve agent success rates on average by $6.7 - 12.5\%$, without any modifications to the agent and underlying LLM.

## 1  Introduction

Large Language Models (LLMs) are increasingly being deployed as autonomous agents to perform complex, multi-turn real-world tasks. LLM agents are equipped with tools, usually in the form of functions, to interact with the outside world, such as database queries [20, 47, 54], file system operations [35, 40, 44], and web searches [32, 45, 53]. The combination of powerful foundation models and domain-specific tools enables autonomous systems that can perform complex real-world tasks, such as customer service [15, 46], business data management [12, 13], and digital assistant [11, 21].

Despite this promise, practical agent deployments suffer from low task success rates [13, 19, 29, 38, 43, 52]. Agent failures occur when agents perform tasks incorrectly, producing misleading output or undesirable side effects, requiring costly human interventions [16, 37, 41]. To address this challenge, prior works focus on improving the agent/LLM, such as building foundation LLMs with stronger agentic abilities [10, 34] or fine-tuning LLMs for specific workloads [36, 48, 49]. We refer to this paradigm as *agent-for-system*. However, while state-of-the-art LLMs demonstrate strong intrinsic abilities in mathematics and logic [6, 9, 31, 34], their performance on real-world agentic tasks remains limited. OpenAI o3 [34], a model optimized for agentic use cases, improves success rates on customer service workloads [46] by only 2-3% [34] over its predecessor o1, illustrating that strengthening models alone has limited effectiveness.

To understand why agents fail despite strong underlying LLMs, we conducted a preliminary analysis of failure traces from state-of-the-art agentic benchmarks [35, 46]. We found that failures stem not only from flawed LLM reasoning, but also from the agent's interaction with its surrounding environment, such as inability to discover information, follow workload constraints, and process large volumes of data [3, 13, 20]. This led us to ask: *How significant is the system environment for agent reliability, and how much improvement can environment optimizations achieve?* In this paper, we investigate this complementary but largely overlooked direction: optimizing the system environment to improve agent reliability, an approach which we term *system-for-agent*.

However, designing effective environments for agents is non-trivial, as real-world system environments are complex

1

**Table 1.** Taxonomy of agent system-interaction failures and corresponding environment optimizations.

| Category | Subcategory | Description | Optimization |
|---|---|---|---|
| Exploration Failures (§ 4.1) | State-space Navigation Failure | Agent fails to navigate the environment to retrieve all necessary data required to complete the task. | Environment lookahead (§ 5.1) |
| | State Awareness Failure | Agent has incorrect understanding about its current position within the environment. | Explicit agent state changes (§ 5.1) |
| Exploitation Failures (§ 4.2) | Tool Output Processing Failure | Agent makes computational errors (e.g., comparisons, ranking) when processing information gathered from tool outputs. | Offload common computations (§ 5.2) |
| | Domain Rule Violation | Agent fails to follow domain rule by either performs forbidden actions or incorrectly blocks valid actions. | Offload domain rule validations (§ 5.2) |
| | User Instruction Following Failure | Agent fails to follow user's specific instructions as requested. | No direct environment optimizations |
| Resource Exhaustion (§ 4.3) | – | Agent fails due to exceeding allocated maximum number of turns or tokens before task completion. | Speculative agentic actions (§ 5.3) |

and vary widely across different domains. To bridge this gap, we present an in-depth study of agent-environment interactions and a systematic methodology for optimizing system environments to improve agent reliability. We collect agent failure traces comprising 142 failed agent tasks and 3,656 turns of agent interactions. To precisely localize agent failures, we introduce a *subtask-based abstraction* inspired by Hierarchical Task Networks (HTNs) from the automated planning literature [8, 30]. Using this abstraction, we annotate the collected failure traces to identify the specific subtask where each failure first occurred. Based on these annotations, we propose a taxonomy of agent-environment interaction failures, summarized in Table 1. We categorize failures into three main types: 1) *Exploration failures*, which occur when agent fails to gather all information required to complete the task, 2) *Exploitation failures*, which arise when agent incorrectly processes information it has collected, and 3) *Resource exhaustion*, which accounts for failures due to exceeding the pre-allocated turns/token limit.

Based on our analysis, we propose Aegis, a set of system-level optimizations to address the identified failure modes, as summarized in Table 1. To mitigate exploration failures, we enhance environment observability through expanding the agent's observation window and explicitly communicating state changes. To reduce exploitation failures, we offload deterministic reasoning operations such as ranking and rule-checking from the agent to be performed in the environment instead. Lastly, we introduce speculative agentic actions to address resource exhaustion by preemptively bundling related tool calls, reducing turn count and token consumption.

We evaluate the effectiveness of our environment optimizations across 5 state-of-the-art agent benchmarks, achieving $6.7 - 12.5\%$ average success rate improvements on average without any modifications to the agent. For context,

this improvement is comparable to or greater than those typically seen between major LLM model generations, which are typically a few percent [1, 33, 34]. We further validate these improvements using a model fine-tuned for agentic workloads: xLAM-2 8B [36]. Our environment optimizations increase success rates by $3.3 - 8.8\%$ on average, demonstrating that system-for-agent complements model fine-tuning. Lastly, we observe an unexpected benefit of environment optimizations: the monetary cost of LLM inference API is reduced by $7.1 - 17.7\%$ as a result of more efficient agent-environment interactions.

In summary, we make the following contributions:

- We collect and analyze 142 failed tasks and 3,656 agent interactions across 5 diverse agentic workloads. Based on this analysis, we propose a taxonomy for categorizing agent-environment interaction failures, consisting of 6 distinct failure modes.

- We introduce Aegis, a set of targeted environment optimizations to address each type of failure in our taxonomy.

- We implement and evaluate system environment optimizations on agentic benchmarks, demonstrating 6.7-12.5% improvements in task success rates on 3 representative LLM models, GPT 4.1, GPT 4.1 mini, and o3, without modifying the agent or underlying LLM model.

## 2 Background

### 2.1 Large Language Model Agents

Large Language Model (LLM)-based agents augment LLMs with the ability to call external tools or functions, enabling them to complete complex tasks autonomously. Through tool use, agents can gather information and manipulate their surrounding system environment to complete a variety of objectives. A typical agent execution trace is illustrated in Figure 2. The agent is first presented with a system prompt

**Table 2.** Agent goals, environments, and key tools for each workload.

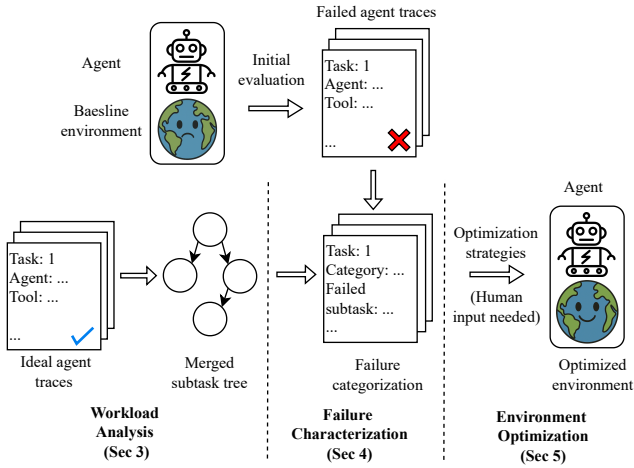| Domain | Agent Goal | Environment | Key Tools |
|---|---|---|---|
| Airline [46] | Book or modify flights | 3 DBs for user, reservations, and flights | `get_user_details,` `search_direct_flight, modify_flight` |
| Retail [46] | Return or exchange orders | 3 DBs for user, orders, and products | `get_user_details, get_order_details,` `return_item` |
| File system [35] | File system operations | Emulated Linux file system | `cd, ls, rm, cp` |
| CRM [12] | Customer Relationship Management (CRM) tasks | Realistic Salesforce platform (SQL DBs) | `get_cases, calculate_avg,` `get_shipping_state` |
| Medical [14] | Mmanaging medical records | Medical records DB over HTTP | `GET Patient, GET Record` |



**Figure 1.** Methodology overview.



System prompt:
- Available functions: [tool1, tool2, ...]
- Domain rules: <rule1, rule2, ...>
**User task:** "Help accomplish <goal>"

**Agent:** "Let me think step-by-step..." Invokes tool1(params)
**Tool:** tool1 response: <result>
**Agent:** Invokes tool2(params)
...
**Agent:** <terminate>

**Figure 2.** Typical agent execution flow.

## 2.2 The Accuracy Challenge of LLM Agents

Despite their promise, LLM agents often fail to complete complex, real-world tasks [3, 13, 41, 50]. This is because real-world workloads demand specialized domain knowledge, interaction with complex environments, and processing of large-scale data [3, 13, 20, 23, 38]. In practice, failures require costly fallback mechanisms: humans may need to take over incorrectly performed tasks, or the agent must re-execute the task from scratch, both of which reduce the efficiency gains promised by automation [4, 16, 18, 37, 41].

Current efforts to improve agent reliability primarily focus on enhancing the underlying LLM. One strategy is to improve the general capabilities of foundation LLM models. For instance, state-of-the-art LLMs such as GPT-4.1 and o3 are explicitly optimized for agentic use cases [33, 34]. Another approach is to fine-tune LLM for a specific workload. This requires building custom pipelines for data collection, model training, and evaluation [22, 36, 39, 48, 49]. The Salesforce xLAM-2 models, for example, are fine-tuned from Llama base models and demonstrated a higher task success rate than their base counterparts [49]. While these agent-for-system strategies are effective at improving agent capabilities, they require a costly and complex multi-stage pipeline of data preparation, model training, and evaluation [36, 48, 49].

In this work, we study a complementary but overlooked direction: enhancing agent reliability by optimizing its surrounding environment. Our methodology, illustrated in Figure 1, consists of three steps. Workload Analysis (§ 3): we

that specifies the set of available functions, their arguments, descriptions, and return values. The system prompt may also define domain rules that constrain the agent's behavior, such as prohibiting certain actions. The agent is given a task, such as retrieving data, performing analysis, or modifying the environment. During execution, the agent proceeds in discrete steps. In each step, the agent may reason and optionally invoke tools via a function calls, and the environment returns a response. This loop continues until the agent decides to terminate and optionally returns a final answer.

To evaluate agent performance on real-world tasks, several benchmarks have been proposed. In this work, we study five state-of-the-art benchmarks (Table 2), each of which includes an environment, available tools, and a set of tasks. For each task, the benchmark defines an expected execution, including the correct sequence of tool calls and the final environment state. Agent correctness is measured by comparing its actual trace against these expectations (e.g., whether the correct tools were invoked in the right order, and whether the final environment state matches the expected outcome).
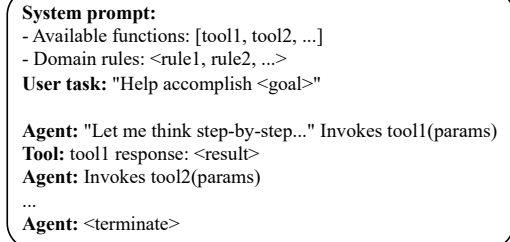
analyze the fundamental structure of each workload using a subtask abstraction to establish a baseline for correct agent behavior. Failure Taxonomy (§ 4): captures agent failures stemming from environment interactions and identifies recurring failure patterns. Environment Optimizations (§ 5): target common failure modes to improve reliability.

## 3 Agentic Workload Analysis

In this section, we study properties of agentic workloads and describe the subtask abstraction for agent-environment interaction analysis.

### 3.1 Methodology

We first establish the following definitions:

- Agent: LLM operating as an autonomous decision-maker that initiates actions to complete user objectives[1].
- System environment: The surrounding system with which an agent interacts, such as databases and file systems.
- Task: A problem assigned to an agent. In this paper, each task corresponds to one entry in an agent benchmark.
- Agent trace: A sequence of actions performed by an agent to fulfill a task, successfully or unsuccessfully.

We collect agent traces from the five benchmarks detailed in Table 2. From each benchmark, we randomly sample 50 tasks, for a total of 250 tasks. We partition the tasks in each workload into an *analysis set* (30 tasks) and an *evaluation set* (20 tasks). We perform failure analysis and develop environment optimizations on the analysis set (§ 3 and § 4) and then measure agent success rates on the unseen evaluation set to test the generalizability of our optimizations.

### 3.2 The Subtask Abstraction

We adapt the Hierarchical Task Network (HTN) framework from automated planning literature [30]. HTN decomposes complex tasks into manageable components with explicit dependencies. We use two key HTN concepts:

- Subtask: A discrete, logical unit of work that contributes to the completion of the overall task.
- Subtask graph: A partial ordering of subtasks, where nodes represent subtasks and edges represent dependencies between subtasks.

Consider an airline workload task: "Change my flight tomorrow to the cheapest option". This decomposes into four subtasks: 1) retrieve user profile, 2) get all user reservations, 3) find the cheapest alternative flight, and 4) modify the reservation. The dependencies are: subtask 4 depends on subtasks 2 and 3 (the agent needs both the correct reservation and cheapest flight), while subtask 2 depends on subtask 1 (user profile is required to retrieve reservations).

The key benefit of the subtask abstraction is that it enables us to localize failures to the precise point agent deviates from the correct solution. Prior agent failure studies [3, 7, 50] adopt a task-level failure attribution method, which assigns one failure label, such as "Incorrect Information Retrieval", for each failed agent task. We find this method insufficient for analyzing agent-environment failures, as it does not reveal which part of the system environment should be optimized. For instance, "Incorrect Information Retrieval" failures can be due to different underlying reasons, such as misinterpreting a specific tool response, failing to retrieve a particular piece of data, or misunderstanding the task.

We derive subtasks in two steps. First, each tool call is mapped to a subtask. In the airline example, this yields 4 subtasks: `get_user`, `get_reservations`, `find_flights`, and `modify_reservation`. Second, subtasks that encompass multiple logical operations are split. For instance, we decompose `modify_reservation` into `judge modify` (reasoning about modification eligibility according to airline domain rules) and `modify reservation` (executing the change). This decomposition is critical because it allows us to distinguish between different failure types, which require different solutions. In this case, a `modify_reservation` could stem from two distinct causes: an error in judging the modification's eligibility or an error in executing the change itself.

### 3.3 Exploration and Exploitation Subtasks

We divide subtasks into two categories: *exploration* and *exploitation*. Exploration subtasks gather new information from the environment (e.g., listing flights on a given date, retrieving patient records, browsing directory contents). Exploitation subtasks act on knowledge the agent already possesses (e.g., changing a reservation, deleting a file, ordering a medical procedure). We make this distinction as exploration and exploitation subtasks benefit from distinct environment optimizations. Exploration subtasks require efficient and accurate information collection, benefiting from environments with high observability. Exploitation subtasks require environments that reduce the likelihood of action errors when the agent processes complex information. This distinction drives our optimization strategy. Section 5 demonstrates how we design targeted environment improvements for exploration and exploitation failures.

### 3.4 Workload Analysis Insights

To analyze agent workload properties, we construct *merged subtask graphs* by first building individual subtask graphs for each task in each workload (30 tasks in the analysis set), then aggregating them into a single merged graph. The nodes and edges of the merged graph are the union of those in the individual graphs. A node's weight is the fraction of tasks in which the subtask appears, while an edge's weight is the fraction of tasks where the corresponding dependency exists. For example, a node appearing in 3 of 30 tasks has a weight of
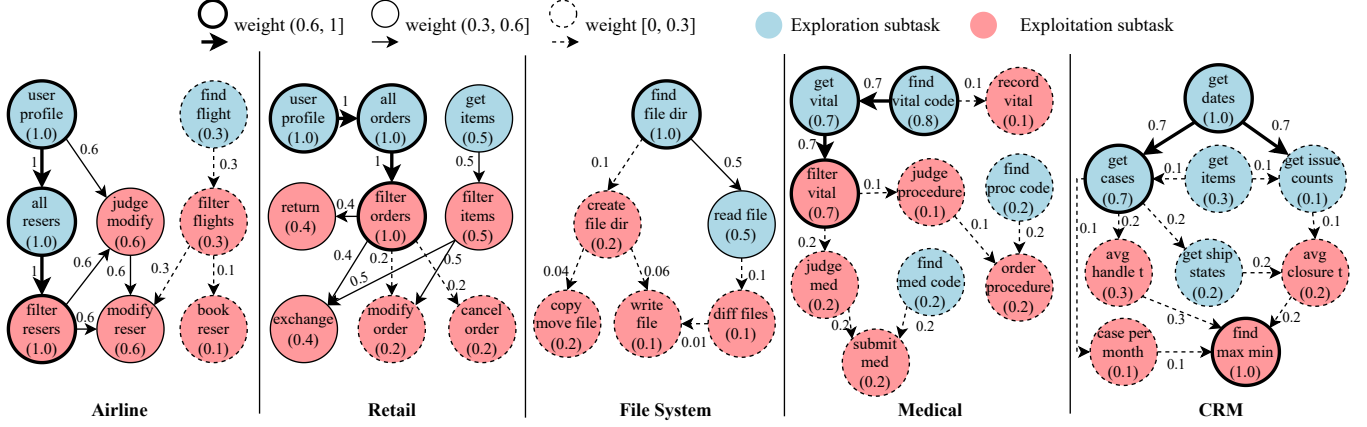
---

**Figure 3.** Merged subtask graphs. Node/edge thickness represent weight. Color encodes explore/exploit subtask.



**(a)** Node weight distributions.
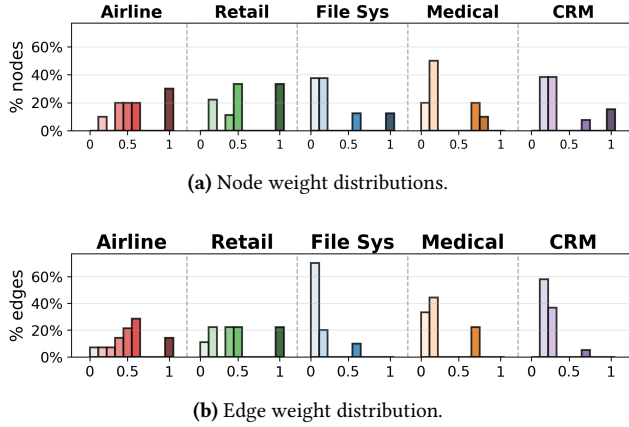


**(b)** Edge weight distribution.

**Figure 4.** Distribution of node and edge weights in each merged subtask graph. X-axis shows weight ranges for each workload from 0 to 1. Y-axis shows percent of nodes/edges that fall within each bucket.

0.1, and an edge appearing in 15 of 30 tasks has a weight of 0.5. Intuitively, node weights indicate how likely an ideal agent would perform each subtask, while edge weights indicate transition probabilities between subtasks.

Figure 3 shows the merged subtask graph for each workload, and Figure 4 plots the weight distributions of nodes and edges. Our analysis reveals two key properties. First, **agent workloads contain frequently recurring subtasks**, represented by high-weight nodes. As shown in Figure 4a, all workloads except the file system have subtasks that appear in over 70% of tasks. These common subtasks typically occur early in the workflow and involve *retrieving critical data objects* required by subsequent steps. For instance, in the airline workload, an agent must retrieve user and reservation details before it can modify a booking. This observation

motivates environment optimizations that improve the reliability of common subtasks, as detailed in Section 5. Second, **workloads exhibit common subtask transitions**, indicated by high-weight edges (Figure 4b). A high weight edge from subtask A to B means that performing A is very likely to be followed by performing B. For instance, in the retail workload, retrieving a user's profile is followed by retrieving their order history in all tasks. This motivates speculative agentic action optimization, as detailed in Section 5.

## 4  Failure Taxonomy

This section presents our taxonomy of agent environment interaction failures. We define a *failure* as the first unsuccessful subtask (as defined in Section 3.2) in an agent's execution trace. This focus on the first failure allows us to identify the root cause of errors, including cascading failures where an early mistake (e.g., retrieving the wrong data) makes all subsequent steps incorrect. We do not classify transient errors that the agent self-corrects without impacting the final outcome, such as a corrected command syntax, as failures. Applying this methodology, we analyzed 142 failed agent traces across 5 workloads and 3 models, containing 3,656 turns of agent-environment interactions. For each of the failed traces, we annotate the subtask where the failure occurred and the failure category. Table 1 presents our complete failure taxonomy, and Figure 5 shows the distribution of failure categories across workloads and models. We now explain each failure category in detail.

### 4.1  Exploration Failures

Exploration failures occur when an agent fails to gather information critical to completing its task. Exploration failures account for 30%, 27%, and 29% of all failures under GPT 4.1, GPT 4.1 mini, and o3, respectively. We conceptualize agent exploration as a search process where agents use tools to navigate an environment and collect information. For instance,
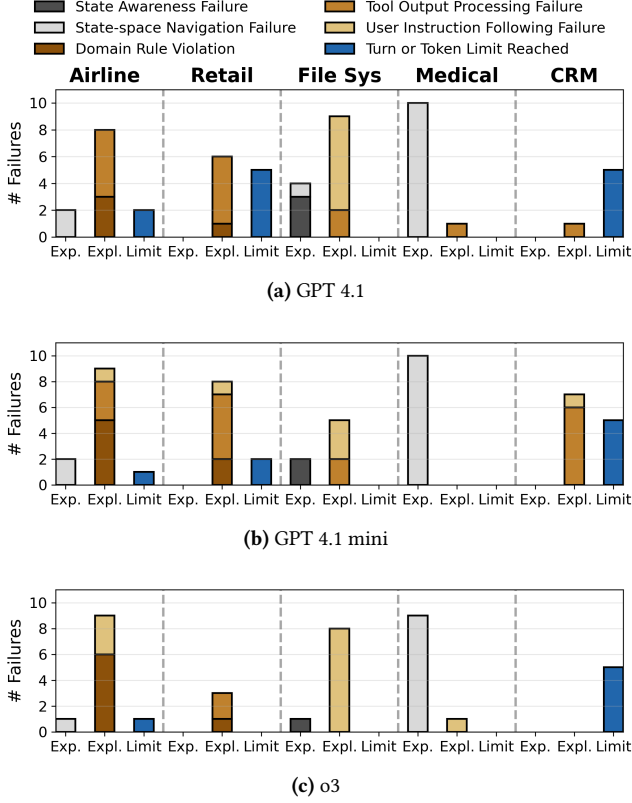
**Figure 5.** Failure category distribution. Exp: exploration. Expl: exploitation. Limit: resource exhaustion.

in a file system, the agent uses the `cd` tool to navigate to various directories to gather locations of different files. The search process can be multi-step. In the medical workload, to retrieve records associated with a specific patient, the agent must first invoke the `GET Patient` tool, then use information from the tool response to invoke `GET Record&patient_id`. The goal of exploration is to find the *critical context*: the minimal set of information required to complete the given task. For instance, to manipulate a specific file, the agent must at least know the file's location within the file system. Using this conceptualization of agent exploration, we identify two distinct subtypes of this failure:

**State-space Navigation Failure** occurs when the agent fails to navigate the environment to collect the critical context. This often happens when the agent's limited view of the environment causes incorrect navigation. This is illustrated in Figure 8b, where the critical context lies beyond the agent's observation window. The agent then explores the environment in a best-effort manner, hoping to reach the critical information. This failure mode is dominant in the medical workload, accounting for 91%, 100%, and 90% of failures under GPT 4.1, GPT 4.1 mini, and o3, respectively. A common type of failure occurs when a medical agent fails

to retrieve a patient's complete medical history. The `GET Record` tool returns partial patient records by default for better resource efficiency. However, the agent is often not aware of this output truncation, causing it to perform downstream tasks based on this incomplete information.

**State Awareness Failure** occurs when the agent loses track of its current position within the environment during navigation. This is distinct from a navigation failure, where the agent correctly identifies its position but does not know how to navigate to find information. In this case, the agent's internal representation of its current state is incorrect. For example, in the file system workload, agents often misidentified their current working directory, leading them to execute file operations in the wrong location. This failure accounted for 23% and 29% of failures in the file system workload under GPT 4.1 and GPT 4.1 mini, respectively. This type of failure is less common under o3, accounting for 10% of all failures.

## 4.2 Exploitation Failures

Exploitation failures occur when an agent possesses the necessary information gathered through exploration but fails to utilize this information correctly. Exploitation failures account for the majority of failures in nearly all workloads (Figure 5), comprising 75% of failures in airline, 88% in retail, 60% in file system, and 33% in CRM under GPT 4.1. We categorize exploitation failures into three subtypes based on the source of the information that the agent failed to process correctly: tool outputs, domain rules, and user instructions.

**Tool Output Processing Failure** occurs when the agent incorrectly utilizes information gained from tool outputs in previous turns. Tool output processing failure is particularly common in the airline, retail, and CRM workloads, accounting for 25%, 67%, and 33% of their respective failed tasks under GPT 4.1. For instance, a CRM agent correctly gathered all relevant customer service handling times but miscalculated the average case handling time. We observe that this failure commonly arises from incorrect logical or mathematical operations performed by the agent. We further classify them based on the operation which the agent performs incorrectly: comparison (e.g., all items with price less/greater than threshold), calculation (e.g., total price), retrieval (e.g., find the most expensive item), and sorting (e.g., earliest/latest arrival time), which comprises 33%, 25%, 25%, and 17% of all tool output processing failures, respectively.

Surprisingly, we observe that agent fails at relatively simple operations such as comparisons. We confirm this by prompting the agent with only the essential information for the failed operation. For instance, we prompt the agent to find the average time from a list by providing only the list of time intervals. We find that in this case the LLM can consistently perform the previously failed operation correctly. We conclude that the agent's failure to perform these operations is not an intrinsic limitation of the LLM, but a consequence of *context distractions* [38]. State-of-the-art LLMs are capable

of solving difficult math problems [6, 9, 31, 34]. However, in agentic workloads, the agent's context is crowded with additional information, such as system prompts, user instructions, and tool outputs that are not immediately relevant to the current operation. As a result of the fundamental nature of the attention mechanism behind LLMs, their reasoning quality degrades under this distraction [19, 38].

**Domain Rule Violation** occurs when an agent's action violates domain rules of the workload. This failure type is the most common in the airline workload, accounting for 21%, 42%, and 50% of all failures under GPT 4.1, GPT 4.1 mini, and o3, respectively. The domain rule most commonly violated is the modification policy of the airline, which restricts the modification of reservations based on the time of purchase, flight cabin, and user membership. Similar to tool output processing failures, we observe that while the LLM can correctly apply domain rules when presented in isolation, the addition of extra information makes it more difficult for the agent to adhere to domain rules correctly.

In our analysis, we further observe two subtypes of domain rule violations. The first is *invalid action*, where the agent performs an action that violates a domain rule (e.g., modifying a flight that is not allowed to be modified). The second type is a *lack of correct action*, where the agent incorrectly concludes that a valid action is forbidden (e.g., informing a user they cannot upgrade an eligible ticket). Addressing these two types of domain rule violations requires different optimizations, as we will detail in Section 5.2.

**User Instruction Following Failure** occurs when the agent fails to follow specific instructions from the user. For example, when asked to create a file with specific content, an agent in the file system workload created the file but populated it with its own interpretation of the content. This failure was most frequent in the file system workload, constituting 40%, 40%, and 89% of its failures under GPT 4.1, GPT 4.1 mini, and o3. Environment optimizations do not directly address this type of failure since we do not modify the user instructions. We still include this for completeness of our taxonomy to account for all agent failures we observe and to understand the limitations of system-level optimization.

### 4.3 Resource Exhaustion Failure

This failure occurs when an agent cannot complete a task before reaching the maximum number of turns or tokens allocated. Such limits are standard practice in agent benchmarks to control cost, bound runtime, and penalize inefficient reasoning [12, 14, 23, 24, 27, 46]. We set a limit of 20 turns and 20,000 tokens per task, consistent with prior work [12, 23, 27, 35], and have verified that all tasks are solvable within these constraints. As shown in Figure 5, resource exhaustion is a major source of failure in the retail and CRM workloads, accounting for 45% and 83% of failures respectively. Token exhaustion is the primary issue in CRM. As detailed in Figure 6a, subtasks such as `get_cases` can return
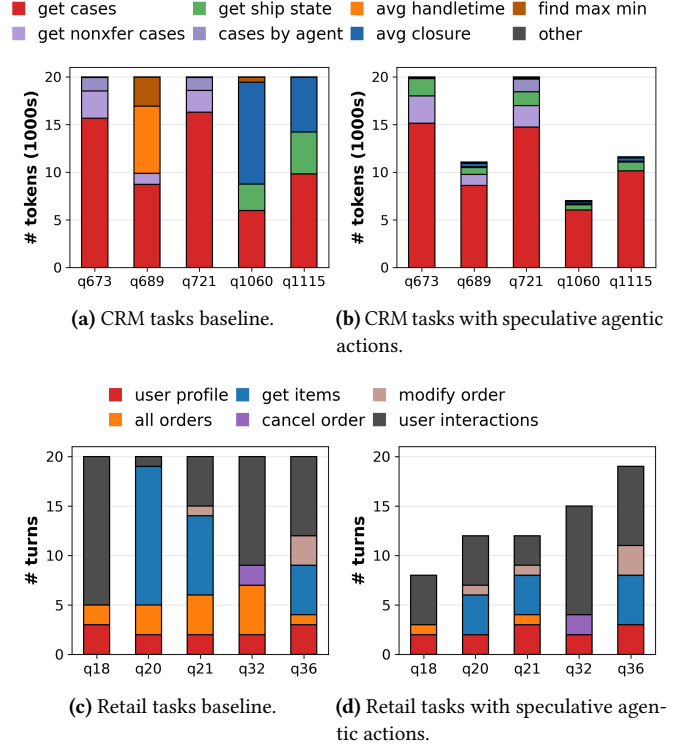


**(a)** CRM tasks baseline.

**(b)** CRM tasks with speculative agentic actions.

**(c)** Retail tasks baseline.

**(d)** Retail tasks with speculative agentic actions.

**Figure 6.** Token/turn breakdown on all GPT 4.1 tasks failed due to resource exhaustion on CRM and retail. Each bar represents a task failed due to resource exhaustion. Stacks show token/turn breakdown by subtasks.

large volumes of data (e.g., over 200 cases in task 721). Subsequent tools that require this data as input argument, such as `calc_avg`, force the agent to generate a correspondingly large argument list, which further exhausts the token budget. On the other hand, turn exhaustion is the dominant issue in retail. As shown in Figure 6b, agents often exceed the turn limit through lengthy `user interactions` or repeated calls to `get_items` to fetch data. We also observe that the agent consistently spends 4-7 turns to retrieve the user profile, then all user orders. These examples show that inefficient interactions, whether from verbose tool outputs or repetitive queries, are a primary cause of resource exhaustion.

### 4.4 Subtask Failure Analysis

We next link the failure categories back to the specific subtasks where they arise. For every subtask, we compute 1) failure rate: the number of times the subtask failed divided by the number of times the subtask appears 2) frequency: the number of times the subtask appears divided by the total number of subtasks. Figure 7 plots failure rate (y-axis) against frequency (x-axis) for all workloads. Subtasks in the upper-right quadrant are both common and brittle, and therefore represent the highest-leverage optimization targets. Across
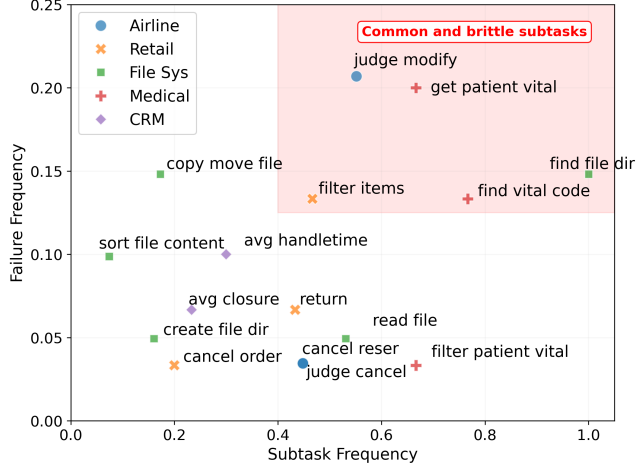
**Figure 7.** Subtask frequency vs. failure frequency (GPT 4.1).

workloads, we identify the following brittle subtasks, which are the primary candidates for environment improvements.

- Airline: `judge reservation modification`
- Retail: `filter product items`
- File system: `find file/directory`
- Medical: `find vital code`, `get patient vital data`
- CRM: `compute average handle time`

## 5 System Environment Optimizations

In this section, we utilize insights from Section 3 and our failure taxonomy to optimize the system environment. Our optimization focuses exclusively on the agent's environment. We do not alter the user's task, the agent's core logic (i.e., the underlying model and prompting techniques), or the task's correctness criteria. We do not add new tools, only augment existing ones. We propose Aegis, a set of targeted optimizations corresponding to each failure category, as shown in Table 1. We now describe each optimization in detail.

### 5.1 Enhancing Environment Observability

To address exploration failures, our principle is to improve the agent's observability of the environment.

**Environment Lookahead.** In Section 4.1, we observe that an agent's exploration is often inefficient because the critical context lies beyond its observation window. The agent then explores the environment in a best-effort manner, hoping to reach the desired information by chance. To address this type of failure, we introduce *environment lookahead*, a technique where exploration tools are augmented to provide a "peek" into adjacent locations in the environment with respect to the agent's current position. As illustrated in Figure 8, this expanded observation window gives the agent a broader view of its surroundings, increasing the probability of successful navigation. For instance, when the environment
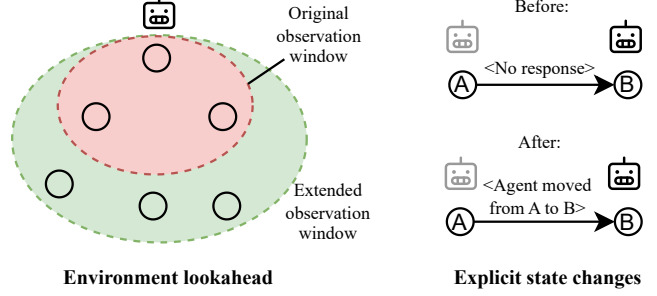


**Figure 8.** Illustration of exploration optimizations.

detects that the agent is retrieving paginated patient records, it can return a message indicating the total number of records and the number of records already retrieved. This hint effectively guides the agent towards the correct exploration actions and reduces the chance of missing critical context.

**Explicit State Changes.** This optimization targets state awareness failures, which occur when an agent's internal understanding of its current location within the environment is incorrect. This failure causes the agent to issue subsequent commands from an incorrect assumption about its current state. To mitigate this, we augment the environment to explicitly communicate the agent's state changes. Whenever a tool call modifies the agent's state, we augment the tool's response to include a confirmation of the new state, as shown in Figure 8. For example, the response from the `cd` command is augmented to confirm the new working directory and list its contents, directly reinforcing the agent's understanding of its location. From the agent's perspective, this optimization turns a reasoning task (tracking the current state through a potentially complex sequence of exploration) into a simple retrieval task (reading the state from the last tool response).

### 5.2 Offloading Common Computational Patterns

In Section 4.2, we observe that while LLMs are fundamentally capable of performing complex reasoning tasks, they still fail to perform simple mathematical operations and rule-based reasoning in agentic workloads due to context distractions. To address exploitation failures, our key principle is to *minimize the amount of reasoning* the agent has to perform to reduce the chance of information synthesis failures. To do this, we offload common computations/reasoning operations from the agent to the environment. Unlike the agent, we can perform operations such as sorting, calculation, and constraint checking reliably in the environment.

**Offload Tool Output Processing.** This optimization targets tool output processing failures. We offload common computations such as sorting (Section 4.2) by augmenting tools to return pre-computed results alongside their primary output. For instance, a tool that retrieves products for the retail agent is augmented to also return the products sorted by price. We note that this optimization is opportunistic, as

the agent may not always need the pre-computed results. However, since such computations are common, the benefits of having pre-computed results readily available outweigh the minor cost of slightly longer tool responses[2].

**Offload Domain Rule Validation.** To address domain rule violations, we offload rule validation from the agent to the environment. Our analysis in Section 4.2 identified two subtypes of these violations: performing invalid action and lack of correct action. We address invalid actions by implementing *environment guardrails*, which embed rule validation directly into tools. For instance, when the airline agent attempts to modify a reservation, the tool first validates the request against domain rules (e.g., time of purchase, flight cabin) and rejects the action if a rule is violated. For failures caused by a lack of correct action, guardrails are ineffective because the agent makes no attempt. Instead, we enrich tool outputs with *domain rule hints*. For example, when the agent retrieves the user's reservation information, the tool response is augmented to indicate whether each reservation is modifiable based on domain rules. This optimization transforms the agent's task from complex rule deduction into simple information retrieval, significantly reducing such errors.
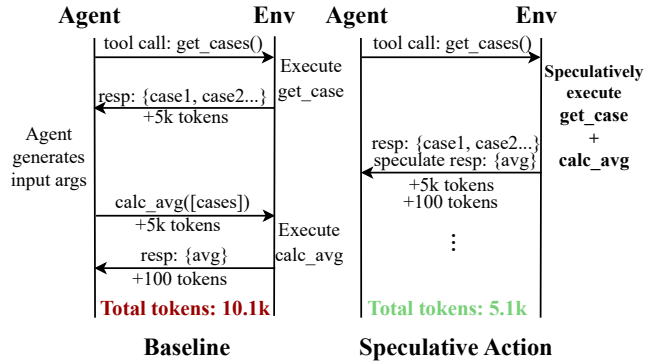


**Figure 9.** Speculative agentic actions example (CRM task).

### 5.3 Speculative Agentic Actions

This optimization targets resource exhaustion failures by speculatively executing likely subtask sequences in a single turn. Using our subtask transition analysis, when an agent calls a tool that typically precedes another, the environment preemptively executes both and bundles the results into a single response. We present an example in Figure 9. As we observed in section 4.3, in the CRM workload, `get_cases` often returns a large set of results. The subsequent `calc_avg` function call requires the list of cases as an input argument, requiring the agent to generate a substantial number of additional tokens. With speculation, when the agent invokes `get_cases`, the environment speculatively

---

[2]The alternative, creating dedicated tools (e.g., `get_cheapest_item()`), would increase overall cost and latency due to an additional tool-use turn.

also executes `calc_avg` using the output from `get_cases`. The response from `calc_avg` is bundled with `get_cases` output and returned to the agent. When this speculation matches the agent's next step (a "hit"), we save the tokens required to generate the input argument of `calc_avg`.

Speculative agentic actions also reduce turn count. In the same CRM example, a speculation hit eliminates an additional agent-environment round trip, as the average result is already available. Thus, we apply speculative agentic actions in retail workload to address tasks where the turn limit is reached. In the retail workload, we observe that agents almost always retrieve the user profile before fetching the user's orders (Figure 3). Thus, we execute `get_user_orders` within `get_user_profile` speculatively.

We note that this optimization is probabilistic. In situations where, e.g., the CRM agent calls `get_cases` without intending to execute `calc_avg`, the speculation results in a "miss," and the extra tokens consumed by `calc_avg` output are wasted (the "miss cost"). As with other speculative techniques (e.g., CPU cache prefetching), this optimization is most effective when the miss cost is low.

## 6 Evaluation

### 6.1 Success Rate

Figure 10 shows the success rate before (Baseline) and after (Optimized) system environment optimizations. We evaluate three models representing different capability levels: general-purpose GPT 4.1, the smaller, economical GPT 4.1 mini, and the state-of-the-art reasoning model o3. We show the results for the analysis set and evaluation set separately.

On average, the environment optimizations lead to 10.3%, 6.7%, 12.5%, 10.0%, and 7.5% improvement in success rate for the five workloads, respectively. To contextualize these gains, the improvements from our system-level approach are comparable to or greater than those typically seen between major LLM model generations, which are often a few percent [1, 33, 34]. For instance, the reported success rate improvement between o1 and o3 on two of our workloads, airline and retail, was 2-3% [34]. This result demonstrates that system-for-agent optimization is a powerful and complementary path to enhancing agent reliability. On average, our environment optimizations improve the success rate of the analysis set by 11.8% and the evaluation set by 7.0%. The higher improvement on the analysis set is expected, as our methodology directly targets failures observed within it. The improvement on the unseen evaluation set demonstrates that our optimizations generalize effectively to new tasks.

### 6.2 Failure Breakdown

We analyze the distribution of failure types before and after environment optimizations in Table 3. We observe that our targeted optimizations successfully reduce the most frequent failures. For instance, the compute offloading optimization
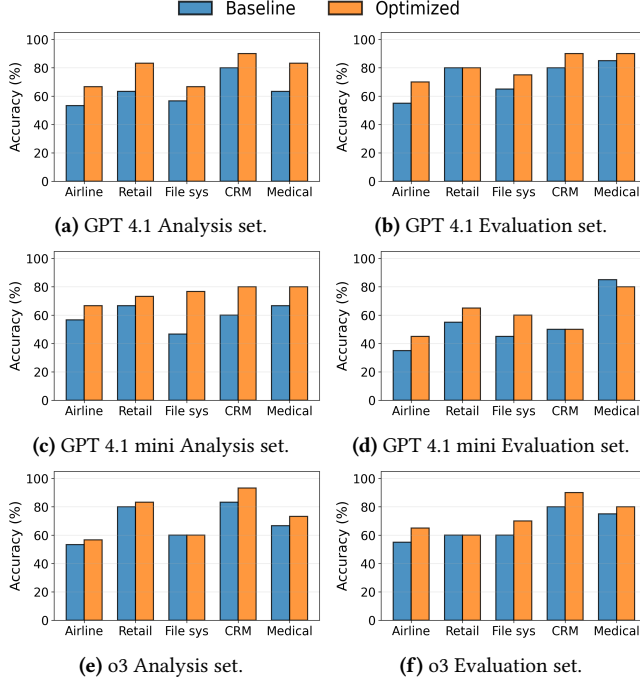
**(a)** GPT 4.1 Analysis set.　　**(b)** GPT 4.1 Evaluation set.

**(c)** GPT 4.1 mini Analysis set.　　**(d)** GPT 4.1 mini Evaluation set.

**(e)** o3 Analysis set.　　**(f)** o3 Evaluation set.

**Figure 10.** Overall task success rates.

**Table 3.** Failure breakdown (counts) before (Base) and after (Opt) system optimizations for GPT 4.1 in the analysis set (30 tasks total in each workload). Dashed lines represent 0 failures before and after optimizations. Underscores show cases where failures have been reduced.

| Failure category | Airline | | Retail | | File Sys | | CRM | | Medical | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Base | Opt | Base | Opt | Base | Opt | Base | Opt | Base | Opt |
| State Awareness | - | - | - | - | 3 | 1 | - | - | - | - |
| State-space Navigation | 2 | 0 | 0 | 1 | 1 | 1 | - | - | 10 | 0 |
| Domain Rule Violation | 3 | 3 | 1 | 0 | - | - | - | - | 0 | 1 |
| Tool Output Processing | 5 | 2 | 5 | 1 | 2 | 2 | 1 | 1 | 1 | 0 |
| User Inst. Following | 0 | 2 | 0 | 3 | 7 | 6 | - | - | - | - |
| Resource Exhaustion | 2 | 2 | 5 | 0 | - | - | 5 | 2 | 0 | 4 |



**(a)** Monetary cost breakdown.　　**(b)** Turns breakdown

**Figure 11.** GPT 4.1 monetary cost and turns before and after environment optimizations.

cuts *Tool Output Processing* failures in the airline workload by 60% (from 5 to 2). Similarly, exploration optimizations eliminated the majority of exploration failures in airline, file system, and medical workload.

A key observation is that mitigating one type of failure can reveal underlying, previously masked failures. For instance, after addressing navigation failures in the medical workload, agents gathered more extensive patient data, which in turn caused some tasks to exceed the token limit, triggering new *Resource Exhaustion Failure* failures. A similar effect occurred in the airline and retail workloads, where resolving exploitation failures revealed *User Instruction Following* failures not previously observed. This observation highlights the compounding nature of agent failures and the need for a holistic approach to improving agent success rate.

**6.2.1 Speculative Agentic Actions.** Figures 6c and d show the effects of speculative agentic actions on the same set of tasks analyzed in Section 4.3. Out of the 5 CRM tasks, Q689, Q1060, and Q1115 become correct by speculatively executing `calculate_avg` after `get_cases`, reducing the token cost of avg handletime and avg closure as described in Section 5.3 Q673 and Q721 remain unsolved as the `get_cases` function returns a large amount of data that still exceeds the token limit despite speculative action savings. Among the 5 retail tasks, Q21, Q32, and Q36 become correct. The turn savings come from speculatively retrieving all user orders within the `get_user_profile` function.

### 6.3 Cost Analysis

Environment optimizations reduce cost in addition to improving accuracy. Figure 11 reports the monetary cost and turns before and after our optimizations. We compute cost using OpenAI API pricing and decompose it into three categories according to the standard LLM inference pricing model: "Prompt" (input tokens that miss the prompt cache), "Cached" (input tokens that hit the prompt cache), and "Completion" (output tokens, including reasoning tokens for o3). On average, environment optimizations reduce monetary cost by 17.4%, 7.1%, and 17.7% for GPT 4.1, GPT 4.1 mini, and o3, respectively. These cost savings are driven by the reduction in the number of turns consumed to complete the task: 16.8%, 7.2%, and 16.6% for the same models. Turn reductions are a result of speculative agentic actions, as we demonstrated in Section 6.2.1.

One exception is the medical workload, where cost increases by 29-33% after optimization. This is because in the baseline environment, agents often fail due to insufficient exploration (Section 4), terminating early without retrieving all relevant data. The observability enhancement enables the agent to correctly gather more information, despite higher token usage. We consider this a favorable tradeoff.
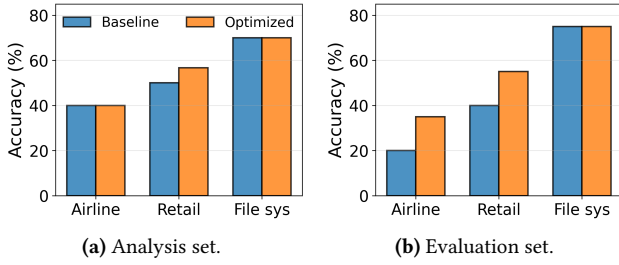
**(a)** Analysis set.      **(b)** Evaluation set.

**Figure 12.** Task success rates on xLAM-2 8B LLM model.

## 6.4 Fine-tuning

To evaluate the effect of our optimizations on fine-tuned models, we test them on xLAM-2 8B, a model fine-tuned from Llama 3 8B [28] that achieves state-of-the-art accuracy compared to models of similar size on $\tau$-bench (airline and retail) and BFCL (file system) workloads [27, 36]. We deploy this open-source model locally using an NVIDIA RTX A6000 GPU and evaluate its performance, shown in Figure 12 We observe that environment optimizations improve the average task success rate by 3.3% on the analysis set and 8.8% on the evaluation set. This result demonstrates that system-level optimizations are complementary to model fine-tuning; they can further enhance agent reliability beyond the benefits gained from fine-tuning. We notice that xLAM-2 8B performs exceptionally well on file system, and our optimizations provide no additional improvement. By analyzing the failed tasks in this workload, we find that most failures remaining are due to user instruction following failures, which cannot be directly addressed by environment optimizations. This is consistent with our findings for other models in Figure 5.

The key advantage of system-for-agent is its lower development effort. Although potentially effective, fine-tuning is a resource-intensive and error-prone process [5, 26, 42]. It involves a complex, multi-stage pipeline of data preparation, model training, and evaluation. The tuned model is also prone to overfitting and catastrophic forgetting [17, 25]. In contrast, environment optimizations are lightweight and interpretable system-level modification that only involves modifying tool APIs, thus enabling rapid iterations. We therefore advocate that when building an agent-based system, developers should first prioritize optimizing the agent's environment. Fine-tuning should be considered only if these system-level improvements are insufficient.

## 7 Related Works

***Failure Analysis of LLM Agents.*** Recent works systematically analyze agent failures through two main approaches. First, taxonomies focusing on agent-centric failures (agent-for-system): Cemri et al. [3] identify multi-agent failure modes such as incorrect inter-agent communication, while

Deshpande et al. [7] categorize reasoning failures such as hallucinations and flawed planning. Our work complements this by focusing on failures due to agent-environment interaction (system-for-agent). Another line of research focuses on automated failure localization methods. Zhang et al. [50] use LLMs to identify which agent step caused task failures, while Arabzadeh et al. [2] introduce a framework for measuring agent misalignment with user expectations. Our work complements these works by introducing a subtask abstraction for finer-grained failure decomposition and analysis. In addition, unlike prior work that primarily identifies failures, we use our analysis to propose, implement, and empirically validate concrete environment optimizations that directly mitigate identified failure modes.

***Specialized Agent Environments.*** One line of work shows that by designing specialized agentic environments, one can achieve superior performance. Zhang et. al. [51] observe that when an agent plays games, environment design can significantly impact agent performance. For example, whether the agent has access to screenshots of the gaming interface. Yang et. al. [45] show that under web browsing tasks, by simply refining the agent's observation and action space, the agent's success rate can be improved. SWE-agent [44] designs custom agent-computer interface tools to improve the performance of software engineering agents. Our work is inspired by these works, so we try to find generalizable techniques for environment optimization.

A growing body of works shows that tailoring an agent's environment can significantly boost its accuracy. These works focus on creating highly specialized, domain-specific interfaces. For example, Zhang et al. [51] observe that an agent's success in games is highly dependent on the environment's design, such as whether the agent receives screenshots or only text. In the context of web browsing, Yang et al. [45] show that refining the agent's observation and action space leads to substantial improvements in task success rates. SWE-agent [44] introduces custom agent-computer interface tools specifically designed to improve the performance of agents on software engineering tasks. Our paper is inspired by these works. However, whereas prior efforts focus on developing specialized solutions for specific domains, our research aims to identify generalizable principles for environment optimization across various workloads.

## 8 Conclusion

In this paper, we present the system-for-agent design paradigm, which treats the agent's environment as a first-class component in building more reliable agentic systems. We introduce a systematic methodology for analyzing agent-environment interactions, leading to a taxonomy of common failure modes. Based on this taxonomy, we develop Aegis, a set of targeted, system-level optimizations to mitigate these

failures. Our evaluations demonstrate that these environment optimizations substantially improve agent task success rates across a variety of workloads and models.

# References

[1] Anthropic. Introducing Claude 4, 2025.

[2] Negar Arabzadeh, Siqing Huo, Nikhil Mehta, Qinqyun Wu, Chi Wang, Ahmed Awadallah, Charles L. A. Clarke, and Julia Kiseleva. Assessing and verifying task utility in llm-powered applications, 2024.

[3] Mert Cemri, Melissa Z. Pan, Shuyi Yang, Lakshya A. Agrawal, Bhavya Chopra, Rishabh Tiwari, Kurt Keutzer, Aditya Parameswaran, Dan Klein, Kannan Ramchandran, Matei Zaharia, Joseph E. Gonzalez, and Ion Stoica. Why do multi-agent llm systems fail?, 2025.

[4] Zichen Chen, Jiaao Chen, Jianda Chen, and Misha Sra. Standard benchmarks fail – auditing llm agents in finance must prioritize risk, 2025.

[5] Databricks. Understanding fine-tuning, 2024.

[6] DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanjia Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025.

[7] Darshan Deshpande, Varun Gangal, Hersh Mehta, Jitin Krishnan, Anand Kannappan, and Rebecca Qian. Trail: Trace reasoning and agentic issue localization, 2025.

[8] Ilche Georgievski and Marco Aiello. An overview of hierarchical task network planning, 2014.

[9] Google. Advanced version of Gemini with deep think officially achieves gold-medal standard at the international mathematical olympiad, 2025.

[10] Google. Gemini 2.5 Pro, 2025.

[11] Yanchu Guan, Dong Wang, Zhixuan Chu, Shiyu Wang, Feiyue Ni, Ruihua Song, Longfei Li, Jinjie Gu, and Chenyi Zhuang. Intelligent virtual assistants with llm-based process automation, 2023.

[12] Kung-Hsiang Huang, Akshara Prabhakar, Sidharth Dhawan, Yixin Mao, Huan Wang, Silvio Savarese, Caiming Xiong, Philippe Laban, and Chien-Sheng Wu. 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)*, 2025.

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

[14] Yixing Jiang, Kameron C. Black, Gloria Geng, Danny Park, James Zou, Andrew Y. Ng, and Jonathan H. Chen. Medagentbench: A realistic virtual ehr environment to benchmark medical llm agents, 2025.

[15] Sehyeong Jo and Jungwon Seo. Proxyllm : Llm-driven framework for customer support through text-style transfer, 2024.

[16] Sidhant Kabra. The hidden cost of ignoring llm failures, 2025.

[17] Damjan Kalajdzievski. Scaling laws for forgetting when fine-tuning large language models, 2024.

[18] Thomas Kwa, Ben West, Joel Becker, Amy Deng, Katharyn Garcia, Max Hasin, Sami Jawhar, Megan Kinniment, Nate Rush, Sydney Von Arx, Ryan Bloom, Thomas Broadley, Haoxing Du, Brian Goodrich, Nikola Jurkovic, Luke Harold Miles, Seraphina Nix, Tao Lin, Neev Parikh, David Rein, Lucas Jun Koba Sato, Hjalmar Wijk, Daniel M. Ziegler, Elizabeth Barnes, and Lawrence Chan. Measuring ai ability to complete long tasks, 2025.

[19] Philippe Laban, Hiroaki Hayashi, Yingbo Zhou, and Jennifer Neville. Llms get lost in multi-turn conversation, 2025.

[20] Fangyu Lei, Jixuan Chen, Yuxiao Ye, Ruisheng Cao, Dongchan Shin, Hongjin Su, Zhaoqing Suo, Hongcheng Gao, Wenjing Hu, Pengcheng Yin, Victor Zhong, Caiming Xiong, Ruoxi Sun, Qian Liu, Sida Wang, and Tao Yu. Spider 2.0: Evaluating language models on real-world enterprise text-to-sql workflows, 2025.

[21] Yuanchun Li, Hao Wen, Weijun Wang, Xiangyu Li, Yizhen Yuan, Guohong Liu, Jiacheng Liu, Wenxing Xu, Xiang Wang, Yi Sun, Rui Kong, Yile Wang, Hanfei Geng, Jian Luan, Xuefeng Jin, Zilong Ye, Guanjing Xiong, Fan Zhang, Xiang Li, Mengwei Xu, Zhijun Li, Peng Li, Yang Liu, Ya-Qin Zhang, and Yunxin Liu. Personal llm agents: Insights and survey about the capability, efficiency and security, 2024.

[22] Weiwen Liu, Xu Huang, Xingshan Zeng, Xinlong Hao, Shuai Yu, Dexun Li, Shuai Wang, Weinan Gan, Zhengying Liu, Yuanqing Yu, Zezhong Wang, Yuxian Wang, Wu Ning, Yutai Hou, Bin Wang, Chuhan Wu, Xinzhi Wang, Yong Liu, Yasheng Wang, Duyu Tang, Dandan Tu, Lifeng Shang, Xin Jiang, Ruiming Tang, Defu Lian, Qun Liu, and Enhong Chen. Toolace: Winning the points of llm function calling, 2025.

[23] Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, Shudan Zhang, Xiang Deng, Aohan Zeng, Zhengxiao Du, Chenhui Zhang, Sheng Shen, Tianjun Zhang, Yu Su, Huan Sun, Minlie Huang, Yuxiao Dong, and Jie Tang. Agentbench: Evaluating llms as agents, 2023.

[24] Jiarui Lu, Thomas Holleis, Yizhe Zhang, Bernhard Aumayer, Feng Nan, Felix Bai, Shuang Ma, Shen Ma, Mengyu Li, Guoli Yin, Zirui Wang, and Ruoming Pang. Toolsandbox: A stateful, conversational, interactive evaluation benchmark for llm tool use capabilities, 2025.

[25] Yun Luo, Zhen Yang, Fandong Meng, Yafu Li, Jie Zhou, and Yue Zhang. An empirical study of catastrophic forgetting in large language models during continual fine-tuning, 2025.

[26] Lindsay MacDonald. Rag vs fine tuning: How to choose the right method, 2024.

[27] Huanzhi Mao. Bfcl v3 multi-turn and multi-step function calling evaluation, 2024.

[28] Meta. Introducing meta llama 3: The most capable openly available llm to date, 2024.

[29] Akshat Naik, Patrick Quinn, Guillermo Bosch, Emma Gouné, Francisco Javier Campos Zabala, Jason Ross Brown, and Edward James Young. Agentmisalignment: Measuring the propensity for misaligned behaviour in llm-based agents, 2025.

[30] Nilufer Onder. Hierarchical task network (htn) planning, 2012.

[31] OpenAI. Learning to reason with LLMs, 2024.

[32] OpenAI. Introducing ChatGPT agent: bridging research and action, 2025.

[33] OpenAI. Introducing GPT-4.1 in the API, 2025.

[34] OpenAI. Introducing openai o3 and o4-mini, 2025.

[35] Shishir G Patil, Huanzhi Mao, Fanjia Yan, Charlie Cheng-Jie Ji, Vishnu Suresh, Ion Stoica, and Joseph E. Gonzalez. The berkeley function calling leaderboard (BFCL): From tool use to agentic evaluation of large language models. In *Forty-second International Conference on Machine Learning*, 2025.

[36] Akshara Prabhakar, Zuxin Liu, Ming Zhu, Jianguo Zhang, Tulika Awalgaonkar, Shiyu Wang, Zhiwei Liu, Haolin Chen, Thai Hoang, Juan Carlos Niebles, Shelby Heinecke, Weiran Yao, Huan Wang, Silvio Savarese, and Caiming Xiong. Apigen-mt: Agentic pipeline for multi-turn data generation via simulated agent-human interplay, 2025.

[37] Yangjun Ruan, Honghua Dong, Andrew Wang, Silviu Pitis, Yongchao Zhou, Jimmy Ba, Yann Dubois, Chris J. Maddison, and Tatsunori Hashimoto. Identifying the risks of lm agents with an lm-emulated sandbox, 2024.

[38] Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed Chi, Nathanael Schärli, and Denny Zhou. Large language models can be easily distracted by irrelevant context, 2023.

[39] Wentao Shi, Mengqi Yuan, Junkang Wu, Qifan Wang, and Fuli Feng. Direct multi-turn preference optimization for language agents, 2025.

[40] Zeru Shi, Kai Mei, Mingyu Jin, Yongye Su, Chaoji Zuo, Wenyue Hua, Wujiang Xu, Yujie Ren, Zirui Liu, Mengnan Du, Dong Deng, and Yongfeng Zhang. From commands to prompts: Llm-based semantic file system for aios, 2025.

[41] Deepak Singla. Salesforce study finds llm agents fail 65

[42] Cohere team. Understanding fine-tuning, 2025.

[43] Mingzhe Xing, Rongkai Zhang, Hui Xue, Qi Chen, Fan Yang, and Zhen Xiao. Understanding the weakness of large language model agents within a complex android environment, 2024.

[44] John Yang, Carlos E. Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan, and Ofir Press. Swe-agent: Agent-computer interfaces enable automated software engineering, 2024.

[45] 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, 2025.

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

[47] Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir Radev. Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-sql task, 2019.

[48] Jianguo Zhang, Thai Hoang, Ming Zhu, Zuxin Liu, Shiyu Wang, Tulika Awalgaonkar, Akshara Prabhakar, Haolin Chen, Weiran Yao, Zhiwei Liu, Juntao Tan, Juan Carlos Niebles, Shelby Heinecke, Huan Wang, Silvio Savarese, and Caiming Xiong. Actionstudio: A lightweight framework for data and training of large action models, 2025.

[49] Jianguo Zhang, Tian Lan, Ming Zhu, Zuxin Liu, Thai Hoang, Shirley Kokane, Weiran Yao, Juntao Tan, Akshara Prabhakar, Haolin Chen, Zhiwei Liu, Yihao Feng, Tulika Awalgaonkar, Rithesh Murthy, Eric Hu, Zeyuan Chen, Ran Xu, Juan Carlos Niebles, Shelby Heinecke, Huan Wang, Silvio Savarese, and Caiming Xiong. xlam: A family of large action models to empower ai agent systems, 2024.

[50] Shaokun Zhang, Ming Yin, Jieyu Zhang, Jiale Liu, Zhiguang Han, Jingyang Zhang, Beibin Li, Chi Wang, Huazheng Wang, Yiran Chen, and Qingyun Wu. Which agent causes task failures and when? on automated failure attribution of llm multi-agent systems, 2025.

[51] Yuxuan Zhang, Haoyang Yu, Lanxiang Hu, Haojian Jin, and Hao Zhang. General modular harness for llm agents in multi-turn gaming environments, 2025.

[52] Zhexin Zhang, Shiyao Cui, Yida Lu, Jingzhuo Zhou, Junxiao Yang, Hongning Wang, and Minlie Huang. Agent-safetybench: Evaluating the safety of llm agents, 2025.

[53] Shuyan Zhou, Frank F. Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, Uri Alon, and Graham Neubig. Webarena: A realistic web environment for building autonomous agents, 2024.

[54] Xiaohu Zhu, Qian Li, Lizhen Cui, and Yongkang Liu. Large language model enhanced text-to-sql generation: A survey, 2024.