

A Survey of the Training Process of LLM-Based Agents

Anonymous ACL submission

Abstract

Autonomous agents based on large language models (LLMs) are becoming an increasingly prevalent paradigm for tackling complex and real-world tasks. Despite the remarkable zero-shot capabilities of modern LLMs, specialized agent training is often essential to obtain reliable and improved performance for specific target tasks. In this work, we present a structured survey dedicated exclusively to the training process of LLM-based agents. We establish a clear taxonomy by examining key methodological steps: environment setups, data preparation strategies, the formulation of effective learning signals, as well as the training objectives and schemes. Finally, we conclude with discussions on potential future directions.

1 Introduction

Autonomous agents based on large language models (LLMs) have been applied to a wide variety of real-world applications (Wang et al., 2024a; Xi et al., 2025b; Luo et al., 2025a). These agents are designed to complete target tasks by engaging in an interactive loop with their environment. They utilize the LLMs as the core cognitive engine to process environmental perceptions and decide subsequent actions. This paradigm has shown great promise for the development of AI assistants that can automate workflows across diverse domains.

Earlier agent systems primarily depend on prompt engineering techniques applied to the underlying LLMs (Hong et al., 2023; Qian et al., 2024; Wu et al., 2024). Although the powerful context-understanding capability of LLMs can offer viable initial solutions, this paradigm is inherently insufficient for the growing complexity and specificity of real-world applications. A key point is that agents need to interact with external environments to complete the target task. In practice, the LLMs may be unfamiliar with the target environmental observations or action spaces, which are not adequately

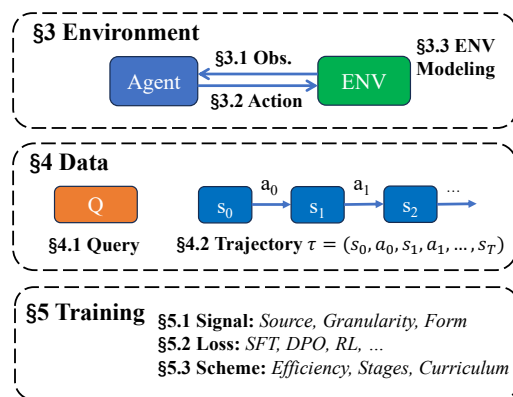


Figure 1: An overview of the survey.

represented in their pre-training or standard post-training stages. This mismatch can substantially degrade task performance and success rates. This motivates the need for dedicated *agent training*, which could fundamentally adapt these LLMs to effectively handle the target agent tasks.

Although there have been many existing surveys for LLM-based agents, few of them place a specific and central focus on the training process of LLM-based agents. General surveys cover agent research *in a broad manner* (Wang et al., 2024a; Xi et al., 2025b; Luo et al., 2025a) or focus on *generalized improvement strategies* (Gao et al., 2025a; Fang et al., 2025a; Du et al., 2025), while others focus on *specific aspects*, such as memory (Zhang et al., 2025j), planning (Huang et al., 2024), tool-use (Qin et al., 2024; Qu et al., 2025a), multi-agent (Guo et al., 2024; Tran et al., 2025) and evaluation (Yehudai et al., 2025), *specific algorithms*, such as reinforcement learning (Zhang et al., 2025c), as well as *specific applications*, such as multi-modal (Xie et al., 2024), web (Ning et al., 2025) and GUI (Zhang et al., 2024a; Tang et al., 2025a; Nguyen et al., 2025a) agents. In this work, we aim to provide a literature review dedicated exclusively to the training process of LLM-based agents.

Figure 1 illustrates the overall structure of this

Algorithm 1 A simplified agent training procedure.

Input: Environment E , Agent System A , Underlying Large Language Model M , Training Iteration I .

Output: Trained Model M' .

```
1:  $M' = M$  ▷ Start with an Initial Model
2: for iter in range( $I$ ) do
3:    $E$ .setup() ▷ Setup Environment (§3)
4:    $Q = \text{collect\_query}()$  ▷ Obtain Query (§4.1)
5:    $T = \text{sample}(Q, A, M', E)$  ▷ Sample Trajectory (§4.2)
6:    $S = \text{get\_signal}(Q, T, E)$  ▷ Learning Signal (§5.1)
7:    $M' = \text{train}(Q, T, S, M')$  ▷ Model Training (§5.2)
8: end for
9: return  $M'$ 
```

069 survey. The primary highlight of this work is its
070 *process-centric organization*. Instead of a purely
071 taxonomic approach, this work is structured around
072 the steps of a typical agent training procedure, in-
073 cluding key topics such as environment setup (§3),
074 data collection (§4), and learning signal generation
075 (§5). With this practical organization, we hope
076 to provide actionable guidance and clear design
077 choices for better implementing and optimizing the
078 entire agent training process.

079 2 Background

080 2.1 LLM-Based Agents

081 An LLM-based agent is an AI system designed
082 to interact with the environment to achieve target
083 goals. These agents leverage the inherent under-
084 standing and reasoning capabilities of LLMs and
085 interact with the environment in an iterative way.
086 At each time step t , the agent receives a represen-
087 tation of the environment’s current state $s_t \in \mathcal{S}$
088 and determines an action $a_t \in \mathcal{A}$, where \mathcal{S} and \mathcal{A}
089 denote the state and action spaces, respectively.
090 The action decision is made by the underlying
091 LLM, which processes the environmental obser-
092 vations and the interaction history as input, and
093 generates the action decision as output. The ac-
094 tion will be executed, leading the environment to
095 a new state s_{t+1} , where the agent will make fur-
096 ther decisions. The interaction will continue until
097 a termination state s_T , resulting in a sequence of
098 states and actions, formally known as a trajectory:
099 $\tau = (s_0, a_0, s_1, a_1, \dots, s_T)$.

100 2.2 Agent Training

101 Although LLMs have remarkable instruction-
102 following capabilities, allowing them to be directly
103 prompted for action decisions without specific pa-
104 rameter updates, the effectiveness of this approach
105 is constrained, especially when the LLM is unfa-

miliar with the state and action representations of
the target environment. To enable more stable and
improved results, dedicated agent training is usu-
ally essential. We outline a typical agent training
process in Algorithm 1, and the remainder of this
survey is structured to examine each of the main
steps in this process. While the overall process
shares similarities with the ordinary LLM training,
the multi-step environmental interactions required
by agent tasks introduce specific considerations.

116 3 Environment

117 The environment is a key component that funda-
118 mentally differentiates agent tasks from standard
119 non-interactive text generation tasks. An agent op-
120 erates in a multi-turn manner by iteratively receiv-
121 ing observations from the environment and decid-
122 ing actions that further influence the environment’s
123 state. Consequently, the representations of these
124 observations and actions directly define the input
125 and output spaces of the underlying decision model.
126 A primary reason why we need agent training is
127 precisely the LLM’s potential unfamiliarity with
128 these specific environmental settings.

The nature of the environment is inherently tied
to the target task, with each task type putting dis-
tinct specifications. For example:

- **Search agents** (Xi et al., 2025a; Zhang et al., 2025h) typically interact with external retrievers or search engines, where actions are formalized as search queries and observations are the corresponding retrieved texts.
- **Web agents** (Ning et al., 2025) interact with a browser. They receive observations of web pages, and their actions include typical browser operations like clicking, scrolling, and typing.
- The scope can be extended to **GUI applications** (Zhang et al., 2024a) for general computer use (Hu et al., 2025b) if visual inputs of screenshots are incorporated.
- In **software engineering** tasks (Liu et al., 2024; Wang et al., 2025d), the environment includes a code repository, enabling agents to apply actions such as code editing and repository manipulation.

149 3.1 Observation Representation

Observation Processing. The current state’s ob-
servation from the environment is a crucial input
signal that directly influences the agent’s decision.
Compared to general natural language inputs, a

unique property of agent observations is their tendency to be *highly structured and complex*. This complexity often requires specialized processing steps to enable the model to accurately understand the current state. To process input observations:

- **Augmentation** methods enrich the original raw observations to ease understanding. For example, techniques like Set-of-Mark (Yang et al., 2023), which adds special marks to annotate visual elements, have been utilized in agent applications (Zhang et al., 2025b). SeeAct (Zheng et al., 2024a) enhances element grounding by augmenting inputs with element attributes and textual choices. OmniParser (Lu et al., 2024b) specifically fine-tunes models to parse screenshots and provide function semantic annotations.
- **Simplification** methods are essential because raw observations are often complex and redundant, a concern amplified by the context length constraints of LLMs. Synapse (Zheng et al., 2024b) extracts task-relevant information from raw states to form cleaner observations. Similarly, AgentOccam (Yang et al., 2025a) demonstrates that providing less redundant yet informative observations can improve agent performance. With similar objectives, LCoW (Lee et al., 2025) trains a contextualization module to enhance this process, FocusAgent (Kerboua et al., 2025) leverages a lightweight retriever to extract the most relevant lines from the observation, and Prune4Web (Zhang et al., 2025e) generates executable Python scoring programs to dynamically filter elements in observations.

History Processing. In addition to the current observation, the record of past interaction steps is a critical input source for the next action decision. A primary challenge here is that as the number of interaction steps increases, the cumulative history representation can quickly grow, making it inefficient and difficult for the LLM to process the entire context effectively. To manage the growing history:

- **Truncation** provides a straightforward solution, by adopting a local window that includes only the details of the most recent steps and omitting previous steps (Tang et al., 2025b).
- **Summarization** provides a more informative approach by converting the previous trajectory into a more concise representation. This approach can be applied periodically (Yen et al., 2025), triggered by heuristic rules such as approaching the

token budget (Wu et al., 2025e; Lu et al., 2025b), or dynamically decided based by the agent itself (Sun et al., 2025c; Ye et al., 2025).

The history management directly relates to the broader concept of agent memory (Zhang et al., 2025j) and context engineering (Mei et al., 2025), as the previous trajectory can be regarded as the agent’s short-term memory, necessitating efficient encoding and retrieval to optimize the sequential decision process.

3.2 Action Representation

The output generated by the LLM directly decides the next action. For LLM-based agents, the output space is *inherently diverse and structured*. This is one of the key differences to classical agents, which often operate with a few fixed action options. The output decision is typically formalized in the style of **function calling** (Patil et al., 2024), where the LLM is required to select and parameterize a function from either a pre-defined action list or a tool-use library. To facilitate reliable parsing and execution, these output representations often utilize structured text formats such as JSON or other special formats. Moreover, **code-based representation** is a natural and elegant way to represent the output (Wang et al., 2024d), and it can enable flexible action compositions by generating and executing programs (Nguyen et al., 2025b). Beyond the format, the **design of the action space** is a critical area of investigation, for which there have been some recent discussions on building more agent-friendly APIs and tools (Song et al., 2025c; Lù et al., 2025; Zhang et al., 2025a).

3.3 Environment Modeling

While the primary goal of agent training is to learn a good policy $\pi(a_t|s_t)$ that decides the next action, an agent system may also benefit from training a model that predicts the transition of the environment $p(s_{t+1}|s_t, a_t)$. This component, often referred to as a world model (Ha and Schmidhuber, 2018), predicts the next state s_{t+1} given the current state s_t and the action taken a_t . Such **environment modeling** can be directly used in the agent’s planning process, enabling decisions to be made with a better internal simulation and awareness of potential future outcomes (Chae et al., 2025a; Gu et al., 2025). Furthermore, it can also enhance the training process through better data synthesis (Fang et al., 2025b) and environment synthesis (Guo et al.,

254 2025b; Liu et al., 2025a), or by serving as the tar- 300
255 get for specialized intermediate pre-training stages 301
256 (Copet et al., 2025). 302

257 4 Data 303

258 Data quality is one of the most important factors for 304
259 agent training. While **manual annotation** could of- 305
260 fer the best way to ensure quality (Lu et al., 2024a; 306
261 Rawles et al., 2023; Deng et al., 2023; Li et al., 307
262 2024), it is inherently tedious and cost-prohibitive. 308
263 Moreover, agent tasks require not only the input 309
264 queries, but also the multi-step interaction trajec- 310
265 tories, posing great challenges to the data curation 311
266 process. Consequently, recent research has focused 312
267 on automated and scalable LLM-based **data syn-** 313
268 **thesis** approaches (Tan et al., 2024). The following 314
269 subsections¹ will illustrate how these techniques 315
270 are specifically applied to agent tasks. 316

271 4.1 Queries 317

272 Similarly to standard LLM tasks, agent training 318
273 first requires the synthesis of queries that define the 319
274 goal that the agent needs to complete. We categor- 320
275 ize related methods based on the complexity and 321
276 the grounding level of the synthesizing process. 322

277 **Direct Synthesis.** This approach leverages the in- 323
278 herent instruction-following capabilities of LLMs 324
279 to synthesize new task queries in a direct way. The 325
280 simplest method involves **direct prompting**, where 326
281 an LLM is asked to generate target tasks, often in- 327
282 spired by existing examples. This is a technique 328
283 widely adopted for agent task synthesis, drawing 329
284 inspiration from Self-Instruct (Wang et al., 2023; 330
285 Patel et al., 2024; He et al., 2025). Considering that 331
286 agent tasks are intrinsically grounded in a specific 332
287 environment, a natural and necessary extension 333
288 is **context-aware querying**. Here, the synthesiz- 334
289 ing LLM is provided with environmental context, 335
290 such as website names and screenshots (Zhou et al., 336
291 2025), crawled web pages (Wu et al., 2025b), pre- 337
292 training documents (Cen et al., 2025), or multi-hop 338
293 paths in knowledge graphs (Lu et al., 2025c). How- 339
294 ever, queries constructed via these direct methods 340
295 may still suffer from limitations, potentially being 341
296 too simple or lacking intricate connections to the 342
297 target environment’s full complexity. 343

298 **Iterative Extension.** To address the simplicity 344
299 limitation of previous methods and improve the 345

¹It is also important to perform *data filtering* to ensure the quality of the training data. This is closely tied to the construction of training signals, which will be covered in §5.1. 346
347
348
349

overall difficulty of the synthesized queries, an **it-** 300
257 **erative extension** approach can be adopted. This 301
258 technique usually starts with a seed item and incre- 302
259 mentally transforms it into a more complex goal 303
260 through repeated refinement. The core goal is to 304
261 generate *challenging but realistic* query-answer 305
262 pairs, a process heavily investigated in recent work 306
263 of deep search agents (Shi et al., 2025a; Gao et al., 307
264 2025b; Li et al., 2025b; Tao et al., 2025). The 308
265 extension process itself is typically mediated by 309
266 an LLM, often *equipped with external tools*. For 310
267 example, WebExplorer (Liu et al., 2025b) oper- 311
268 ates by prompting LLMs equipped with search and 312
269 browsing tools to guide the construction. Moreover, 313
270 this approach can be used to build more complex 314
agent tasks; AgentSynth (Xie et al., 2025), for ex- 315
ample, iteratively forms a sequence of subtasks 316
that are then summarized into a composite task for 317
computer-use agents. Nevertheless, this approach 318
may produce unnatural queries, and building com- 319
plex yet natural tasks remains a challenge. 320

Exploratory Construction. To obtain queries 321
that are more connected to the target environments, 322
the construction itself can be viewed as an agent 323
task, and a dedicated data construction agent can be 324
designed to explore the environment. This typically 325
follows a **reverse construction** process: firstly in- 326
teracting with the environment with an exploration- 327
based policy and then reversely synthesizing the 328
query with the exploration trajectory. This ap- 329
proach has been adopted for a variety of target en- 330
vironments and tasks, including search (Wang et al., 331
2025c), web (Murty et al., 2024b; Pahuja et al., 332
2025; Trabucco et al., 2025), GUI (SU et al., 2025; 333
Sun et al., 2025b; Ramrakhya et al., 2025), and 334
tool-use (Zhai et al., 2025). Although this method 335
can create queries that are closely related to the 336
target environment, there still remain challenges 337
such as constructing non-trivial and challenging 338
tasks and controlling the extra exploration cost. 339

340 4.2 Trajectories 340

341 In addition to the task queries, agent training also 341
342 requires the interaction trajectories to support the 342
343 full multi-turn decision process of the agent. 343

Data Conversion. A simple way to collect tra- 344
jectories is through **data conversion**, by repurpos- 345
ing and transforming existing datasets originally 346
collected for related tasks (Gandhi et al., 2024). 347
For agent applications, this may be challenging 348
due to the format differences across various tar- 349

get tasks and environments. This complexity has driven standardization efforts, with work such as ADP (Song et al., 2025b) and AgentOhana (Zhang et al., 2024b) aiming to define *unified data formats and protocols* to facilitate better data reuse. Furthermore, when the format discrepancy between the original and desired data is large, **LLM-based transformation** can be utilized, where an LLM acts as a robust parser and translator to complete the conversion (Yin et al., 2024a). While data conversion effectively utilizes existing resources, the availability of relevant and high-quality source data is often lacking, necessitating more target-specific trajectory construction approaches.

Sampling. The most fundamental method to collect agent trajectories is **sampling**, which involves using a designated policy to run target tasks directly within the corresponding environment. For the underlying LLM policy used during the sampling process, there are two primary choices: utilizing a stronger teacher model to generate high-quality data that guides the learner (Chen et al., 2023; Zeng et al., 2024), or employing the current version of the agent model itself to enable direct self-improvement (Patel et al., 2024; Fang et al., 2025b). While the simplest implementation involves independently sampling each trajectory in parallel, more complex algorithms like **tree search** (Zhou et al., 2024a; Koh et al., 2025; Lin et al., 2025a; Hou et al., 2025) can be employed to perform deeper exploration and generate better paths.

Special Construction. There have also been a range of **specialized approaches** developed to construct high-quality agent trajectories by leveraging diverse external knowledge or employing special refinements. Some methods focus on **knowledge-based transformation**, such as Synatra (Ou et al., 2024), which uses procedural knowledge resources like WikiHow as a source to be transformed into trajectories, and AgentTrek (Xu et al., 2025), which crawls internet tutorials and converts them into structured representations suitable for guided replay. Other techniques center on refinement and bootstrapping for **quality improvement**: BAGEL (Murty et al., 2024a) applies a round-trip bootstrapping approach that transforms between synthetic instructions and refined trajectories, while UI-Simulator (Wang et al., 2025e) employs a trajectory wrapper that reconstructs raw rollouts into high-quality training instances by inferring user instructions and generating a coherent reasoning

process. Finally, methods like ActRe (Yang et al., 2024) focus on **reasoning injection** by randomly sampling an action and subsequently generating a detailed explanation to compose a full reason-then-act trajectory, and WebCoT (Hu et al., 2025a) curates trajectory data that exemplifies special abilities by reconstructing the agent’s reasoning algorithms into chain-of-thought rationales. These specialized techniques enable the creation of richly annotated agent trajectories.

5 Training

5.1 Signal

The acquisition and utilization of high-quality learning signals is one of the most important factors for successful agent training. In this sub-section, we examine learning signals from three complementary perspectives: **source** (where the signals originate), **granularity** (the level at which the signals are provided) and **form** (the representation of the signals themselves).

5.1.1 Signal Source

Expert Supervision. The most straightforward way to obtain training signals relies on the direct **expert supervision**, including manual annotation or knowledge distillation from superior “teacher” models (Hinton et al., 2015; Kim and Rush, 2016). One typical way is to collect *oracle trajectories* using a teacher model and use them as direct demonstrations to train the agent policy (Chen et al., 2023; Zeng et al., 2024).

Environment-Based Feedback. In many agent scenarios, training supervision is provided by pre-defined external signals from environments. LLM-based agent tasks differ fundamentally from plain LLM tasks in that the agents are interacting with the environment, which can directly provide **environmental learning signals**. The precise form of such signals depends on the target environment: for example, they may be directly determined by changes in the environmental state (Brockman et al., 2016; Shridhar et al., 2021) or calculated by pre-defined evaluation functions (Yao et al., 2022; Zhou et al., 2024b). A typical and valuable characteristic of these external signals is that they could be **rule-based or verifiable**, making them suitable for methods like reinforcement learning with verifiable rewards (RLVR) approaches (Lambert et al., 2024; Guo et al., 2025a). If the target task is to provide a single final answer, signals can be readily

obtained by matching the predicted answer with the gold reference answer (Jin et al., 2025; Song et al., 2025a; Chen et al., 2025; Qian et al., 2025; Wei et al., 2025). In addition to final answers, there can also be external signals that provide **partial evaluation** of the predicted trajectory, such as entity matching ratios (Zhao et al., 2025) or sub-goal scaffolding (Luo et al., 2025b) in complex search tasks. An interesting connection is that external signals can often be obtained **by construction** during the data synthesis phase (§4), where the queries and the corresponding target answers or evaluation functions are constructed simultaneously.

Model-Based Feedback. The learning signals can also be generated by models. It is common to use stronger models in the way of LLM-as-a-judge (Gu et al., 2024), such as providing judgments as RL rewards (Lee et al., 2024) or filtering highest-quality data (Zeng et al., 2024; Trabucco et al., 2025). Moreover, the agent or its underlying LLM to be trained can also serve as an internal source, generating its own learning signals. One of the key metrics is **confidence**, which can potentially indicate the correctness of generated trajectories. This is typically calculated by checking the generation probability of the model outputs (Wu et al., 2025d) or by verifying consistency across multiple generation paths (self-consistency; Huang et al., 2023; Kang et al., 2025). Leveraging the LLM’s general-purpose ability to reason and evaluate, the model can be treated as its own judge in a **self-rewarding** setting (Yuan et al., 2024), generating reward signals based on its own assessment of the generated output. Furthermore, LLMs have demonstrated strong abilities of **self-reflection** (Shinn et al., 2023; Madaan et al., 2023), and this inherent capacity can be leveraged to provide internal critique (Zhang et al., 2025f; Li et al., 2025c).

5.1.2 Signal Granularity

Agent tasks are unique in that they consist of multi-turn trajectories containing interleaved action predictions and environmental observations, rather than just a linear sequence of generated tokens as in plain LLM tasks. Consequently, learning signals can be provided at various levels, including trajectory-level, step-level, or even token-level. Considering the multi-step characteristic of agent tasks, the most interesting discussions are on the comparison between *trajectory-level versus step-level signals*. This directly mirrors the general

LLM alignment debate over outcome reward models (ORM) versus process reward models (PRM). While ORM signals are typically easier to obtain, PRM signals offer denser supervision (Uesato et al., 2022; Lightman et al., 2024; Zheng et al., 2025a).

In agent tasks, ORM provides **trajectory-level signals**, typically as one aggregated evaluation score for the full agent trajectory. Many of the verifiable rewards discussed in the previous section (§5.1.1) are examples of this granularity. The inherent challenge of this approach is that such signals are sparse, often provided only after the full trajectory is completed, which may lack sufficient supervision to effectively guide the learning for the intermediate decision steps.

Considering **step-level signals** is a natural extension, as each step in the agent trajectory corresponds directly to one action decision. Nevertheless, these finer-grained signals are much more difficult to obtain, for which manual annotation (Chae et al., 2025b) could be tedious and expensive. We categorize the typical automatic approaches used to acquire step-level signals in the following:

- **Model-based** methods use LLM-based reward models, where an LLM judge can directly examine each step’s details and assign rewards (Tan et al., 2025b; Yu et al., 2024; Zhang et al., 2025i). In addition, specific PRMs can be explicitly trained to provide these signals. For instance, Liu et al. (2025d) optimize an implicit PRM to provide step-level rewards as the log ratios between the action probabilities with the learned PRM and the old policy, while Xi et al. (2025c) adopt a temporal difference-based estimation method for PRM training. Similarly, Zhou et al. (2024c) learns a high-level value function with off-policy RL to provide signals for low-level action generation.
- **Search-based** techniques leverage complex search algorithms to yield more fine-grained rewards. A common strategy involves generating a search tree with the current policy and assigning rewards to the intermediate steps according to the rewards of its descendants (Hou et al., 2025). The search tree can be constructed in various ways, including Monte Carlo sampling (Xiong et al., 2024), entropy-guided rollout (Shen et al., 2025), or iterative expansion (Lin et al., 2025b; Wu et al., 2025a). Additionally, if step states across different trajectories can be grouped, step-level signals can be calculated using the state rewards within

each group (Feng et al., 2025).

- **Task-specific** methods provide step-level signals uniquely tailored to target tasks. Examples include rewards based on information gain by checking action probabilities depending on the ground-truth answer (Wang et al., 2025a), rewards based on the usefulness and redundancies of retrieved documents for each query action (Zheng et al., 2025c), and specialized rewards generated by detecting suboptimal behaviors for the target search task (Wu et al., 2025c).

Although step-level signals provide denser learning supervision, the performance and stability could be highly dependent on the choice of the learning algorithms (Wang and Ammanabrolu, 2025), and the potential noise inherent in such signals must be carefully considered.

5.1.3 Signal Form

The learning signals can be presented in various forms. The most common form is **numeric rewards**, which provide a scalar judgment on the correctness or goodness of the evaluated trajectory or action. In addition, **natural-language feedback** provides more detailed critique, providing richer and qualitative signals for the training of the LLMs (Scheurer et al., 2022; Chen et al., 2024; Xu et al., 2024) and the agents (Yang et al., 2025b). In scenarios where assigning absolute scores is difficult, it might be easier to provide **preference signals**, indicating which of two actions or trajectories is superior. Such comparative signals are utilized to perform preference learning (Rafailov et al., 2023) for agents (Xiong et al., 2024; Song et al., 2024). Moreover, if expert sources can provide **direct demonstrations** of the correct actions, these signals can be directly leveraged for imitation learning.

A special type of learning signal is the supervision of **error correction**. For example, Lyu et al. (2025) propose a student-centered framework where an expert corrects only the earliest error in a student-generated trajectory, iteratively guiding the erroneous path towards a correct one. However, a central difficulty in obtaining correction signals is in the problem of **error attribution**, which remains an important research question (Cemri et al., 2025; Zhang et al., 2025g,d; Zhu et al., 2025).

5.2 Loss

The loss functions utilized for agent training largely stem from those established for plain LLM training.

The most widely used losses include the standard cross-entropy loss in supervised fine-tuning (SFT), the preference loss used in direct preference optimization (DPO; Rafailov et al., 2023), and various reinforcement learning (RL) losses (Sutton and Barto, 2018). The fundamental details of these functions are omitted here as they are well-known, and unfamiliar readers are referred to related surveys (Zhao et al., 2023; Wang et al., 2024c; Tie et al., 2025; Liu et al., 2025c; Du et al., 2025; Zhang et al., 2025c). Each loss function requires a corresponding type of learning signal (§5.1.3): SFT needs direct demonstrations of the oracle action, DPO requires preference pairs of positive and negative action instances, and RL demands rewards associated with the agent trajectories.

For agent applications, there are several special considerations for the loss. First, since agent trajectories consist of interleaved agent actions and environmental observations, it is common to adopt **observation masking** (or loss masking) during training, as the observation tokens are not directly generated by the model. Jin et al. (2025) show that leveraging loss masking for retrieved tokens can be beneficial. Second, there have been various **modifications** of standard RL losses to adapt to the agent setting and further stabilize training (Wang et al., 2025f; Yu et al., 2025; Deng et al., 2025; Li et al., 2025a). Finally, considering the multi-turn nature of agent trajectories, different loss functions can be applied at different levels, exemplified by ArChER (Zhou et al., 2024c), which adopts a **hierarchical approach** running two distinct RL algorithms in parallel to manage high-level action decision and low-level action token generation.

5.3 Scheme

In Algorithm 1, we only list a simplified version of the agent training scheme; in practice, there are numerous variations with different considerations.

Training Efficiency. In agent training, data collection can be costly, especially in online learning scenarios where trajectory sampling and model updating are interleaved in a fine-grained manner. This poses significant challenges for practical training efficiency. One typical strategy to mitigate this cost is to shift towards **offline learning** approaches (Levine et al., 2020). While pure offline learning assumes training with previously collected data and no additional online interactions, intermediate methods are often considered, adopting an iterative

and **coarse-grained** approach for balancing data freshness with computational savings (Patel et al., 2024; Aksitov et al., 2024; He et al., 2025). Furthermore, to fully utilize available computational resources and parallelize operations, **asynchronous training** methods have been increasingly employed (Tan et al., 2025a; Gao et al., 2025b; Jiang et al., 2025; Lu et al., 2025a).

Training Stages. Similar to plain LLM post-training, the agent training procedure can be split into **multiple stages** with different training schemes. A widely adopted scheme involves first performing an offline phase of SFT on expert demonstrations, followed by an online phase of RL to refine the policy through interaction (Vattikonda et al., 2025; Wu et al., 2025b). Furthermore, to enhance the model’s foundational understanding and better support agentic behaviors, some work incorporates a **continual pre-training** stage that precedes the main post-training (SFT/RL) fine-tuning (Wang et al., 2025b; Su et al., 2025; Copet et al., 2025). This additional stage focuses on preparing the model with relevant knowledge and capabilities for effective downstream agent operation.

Training Curriculum. In addition to the training stages and efficiency considerations, the data and environmental setup can also be adjusted across the training procedure, naturally adopting the idea of **curriculum learning** (Bengio et al., 2009). This approach structures the learning process by starting with simpler tasks or data and gradually increasing the difficulty. Examples of curriculum learning in agent training include starting with easy samples and then progressively introducing more complex ones (Lai et al., 2024), generating new tasks from previous unsuccessful attempts to create a self-evolving curriculum (Qi et al., 2025), and incrementally degrading the quality of simulated environment observations using a curriculum-based rollout strategy to improve generalization (Sun et al., 2025a). This structured progression helps ensure the agent to obtain foundational skills before tackling more challenging scenarios.

6 Future Directions

Scalable Supervision. The main difficulty in training LLM-based agents remains the lack of sufficient high-quality supervision signals that can scale affordably. Agent tasks are inherently challenging because they require tedious interactions

with complex environments, and the reward signals are often sparse and noisy. To reduce the interaction cost, simulations utilizing environment models present a promising direction (Ding et al., 2025; Liu et al., 2025a; Team et al., 2025), allowing the agent to practice and generate data internally. To collect better signals, especially those that are non-verifiable, improving the methods for automatically building accurate reward models remains a foundational problem (Leike et al., 2018).

Multi-modal Agent Training. Although this work does not delve into detailed application scenarios, for many practical agent tasks, the ability to deal with multi-modality (e.g., visual interfaces, audio input) is a requirement (Xie et al., 2024). While the overall training procedure conceptually remains similar, dedicated consideration should be given to grounding issues (Zheng et al., 2024a), where the agent should correctly associate language and symbols with the inputs. The agents need to possess strong core capabilities in understanding the environmental interfaces across any modality and making correspondingly grounded decisions. Although recent developments in multi-modal LLMs have shown some promise (Yin et al., 2024b), integrating all these diverse capabilities into one coherent model or system still remains a great challenge.

Life-Long Learning. An important developmental goal for agents is to enable fast and robust adaptation to the dynamic environments encountered in real-world applications (Zheng et al., 2025b). To achieve this, it will be crucial to integrate agent training with established techniques from life-long or continual learning (Chen and Liu, 2018; De Lange et al., 2021; Wang et al., 2024b; Shi et al., 2024). This inherently involves a holistic integration of the different components within the agent system, making it crucial to accurately identify which components should remain invariant across tasks and which should be adapted quickly in response to external environmental changes.

7 Conclusion

This work provides a structured survey focusing on the training processes of LLM-based agents. We cover key methodological components, including environment setup, data preparation strategies, the design of effective learning signals, training objectives and schemes. We hope that this survey could serve as a practical reference for agent training.

748 Limitations

749 This work has several limitations in its design and
750 scope. Firstly, we specifically focus on the training
751 process of LLM-based agents, omitting other sig-
752 nificant aspects such as architectures, evaluation,
753 and specific applications. This strategic choice is
754 made because many of these areas are already well-
755 covered by existing literature, and we refer read-
756 ers to the other comprehensive agent surveys men-
757 tioned in the introduction (§1). Secondly, the de-
758 scriptions provided in this work are mostly brief to
759 provide a comprehensive coverage within the con-
760 straints of page limits. Our approach is to present
761 related works in meaningful and structured groups
762 rather than describing them in unstructured detail.
763 We hope that this work can serve as a high-level in-
764 dex where further details can be found in the corre-
765 sponding cited papers. Finally, this work is a pure
766 survey without any experiments or empirical re-
767 sults. While performing comparative experiments
768 across various agent training strategies would un-
769 doubtedly provide more meaningful and actionable
770 guidance, we leave this resource-intensive effort
771 for dedicated future work.

772 References

773 Renat Aksitov, Sobhan Miryoosefi, Zonglin Li, Daliang
774 Li, Sheila Babayan, Kavya Kopparapu, Zachary
775 Fisher, Ruiqi Guo, Sushant Prakash, Pranesh Srini-
776 vasan, Manzil Zaheer, Felix Yu, and Sanjiv Kumar.
777 2024. [ReST meets react: Self-improvement for multi-
778 step reasoning LLM agent](#). In *ICLR 2024 Workshop
779 on Large Language Model (LLM) Agents*.

780 Stefano V Albrecht, Filippos Christianos, and Lukas
781 Schäfer. 2024. *Multi-agent reinforcement learning:
782 Foundations and modern approaches*. MIT Press.

783 Yoshua Bengio, Jérôme Louradour, Ronan Collobert,
784 and Jason Weston. 2009. Curriculum learning. In
785 *Proceedings of the 26th annual international confer-
786 ence on machine learning*, pages 41–48.

787 Greg Brockman, Vicki Cheung, Ludwig Pettersson,
788 Jonas Schneider, John Schulman, Jie Tang, and Woj-
789 ciech Zaremba. 2016. Openai gym. *arXiv preprint
790 arXiv:1606.01540*.

791 Mert Cemri, Melissa Z Pan, Shuyi Yang, Lakshya A
792 Agrawal, Bhavya Chopra, Rishabh Tiwari, Kurt
793 Keutzer, Aditya Parameswaran, Dan Klein, Kannan
794 Ramchandran, Matei Zaharia, Joseph E. Gonzalez,
795 and Ion Stoica. 2025. [Why do multi-agent LLM sys-
796 tems fail?](#) In *The Thirty-ninth Annual Conference
797 on Neural Information Processing Systems Datasets
798 and Benchmarks Track*.

Zhepeng Cen, Haolin Chen, Shiyu Wang, Zuxin Liu,
Zhiwei Liu, Ding Zhao, Silvio Savarese, Caim-
ing Xiong, Huan Wang, and Weiran Yao. 2025. [Webscale-rl: Automated data pipeline for scal-
ing rl data to pretraining levels](#). *arXiv preprint
arXiv:2510.06499*. 799 800 801 802 803 804

Hyungjoo Chae, Namyong Kim, Kai Tzu iunn Ong,
Minju Gwak, Gwanwoo Song, Jihoon Kim, Sungh-
wan Kim, Dongha Lee, and Jinyoung Yeo. 2025a. [Web agents with world models: Learning and lever-
aging environment dynamics in web navigation](#). In
*The Thirteenth International Conference on Learning
Representations*. 805 806 807 808 809 810 811

Hyungjoo Chae, Sunghwan Kim, Junhee Cho, Se-
ungone Kim, Seungjun Moon, Gyeom Hwangbo,
Dongha Lim, Minjin Kim, Yeonjun Hwang, Minju
Gwak, Dongwook Choi, Minseok Kang, Gwanhoon
Im, ByeongUng Cho, Hyojun Kim, Jun Hee Han,
Taeyoon Kwon, Minju Kim, Beong woo Kwak, and
2 others. 2025b. [Web-shepherd: Advancing PRMs
for reinforcing web agents](#). In *The Thirty-ninth An-
nual Conference on Neural Information Processing
Systems*. 812 813 814 815 816 817 818 819 820 821

Angelica Chen, Jérémy Scheurer, Jon Ander Campos,
Tomasz Korbak, Jun Shern Chan, Samuel R. Bow-
man, Kyunghyun Cho, and Ethan Perez. 2024. [Learn-
ing from natural language feedback](#). *Transactions on
Machine Learning Research*. 822 823 824 825 826

Baian Chen, Chang Shu, Ehsan Shareghi, Nigel Collier,
Karthik Narasimhan, and Shunyu Yao. 2023. [Fireact:
Toward language agent fine-tuning](#). *arXiv preprint
arXiv:2310.05915*. 827 828 829 830

Mingyang Chen, Linzhuang Sun, Tianpeng Li, Haoze
Sun, Yijie Zhou, Chenzheng Zhu, Haofen Wang,
Jeff Z Pan, Wen Zhang, Huajun Chen, and 1 oth-
ers. 2025. [Learning to reason with search for
llms via reinforcement learning](#). *arXiv preprint
arXiv:2503.19470*. 831 832 833 834 835 836

Zhiyuan Chen and Bing Liu. 2018. *Lifelong machine
learning*. Morgan & Claypool Publishers. 837 838

Jade Copet, Quentin Carbonneaux, Gal Cohen, Jonas
Gehring, Jacob Kahn, Jannik Kossen, Felix Kreuk,
Emily McMilin, Michel Meyer, Yuxiang Wei, and
1 others. 2025. [Cwm: An open-weights llm for re-
search on code generation with world models](#). *arXiv
preprint arXiv:2510.02387*. 839 840 841 842 843 844

Matthias De Lange, Rahaf Aljundi, Marc Masana, Sarah
Parisot, Xu Jia, Aleš Leonardis, Gregory Slabaugh,
and Tinne Tuytelaars. 2021. [A continual learning sur-
vey: Defying forgetting in classification tasks](#). *IEEE
transactions on pattern analysis and machine intelli-
gence*, 44(7):3366–3385. 845 846 847 848 849 850

Wenlong Deng, Yushu Li, Boying Gong, Yi Ren, Chris-
tos Thrampoulidis, and Xiaoxiao Li. 2025. [On
grp collapse in search-rl: The lazy likelihood-
displacement death spiral](#). *arXiv preprint
arXiv:2512.04220*. 851 852 853 854 855

856	Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Sam Stevens, Boshi Wang, Huan Sun, and Yu Su. 2023. Mind2web: Towards a generalist agent for the web. <i>Advances in Neural Information Processing Systems</i> , 36:28091–28114.	911
857		912
858		913
859		914
860		915
861	Jingtao Ding, Yunke Zhang, Yu Shang, Yuheng Zhang, Zefang Zong, Jie Feng, Yuan Yuan, Hongyuan Su, Nian Li, Nicholas Sukiennik, and 1 others. 2025. Understanding world or predicting future? a comprehensive survey of world models. <i>ACM Computing Surveys</i> , 58(3):1–38.	916
862		917
863		918
864		919
865		920
866		921
867	Shangheng Du, Jiabao Zhao, Jinxin Shi, Zhentao Xie, Xin Jiang, Yanhong Bai, and Liang He. 2025. A survey on the optimization of large language model-based agents. <i>arXiv preprint arXiv:2503.12434</i> .	922
868		923
869		924
870		925
871	Jinyuan Fang, Yanwen Peng, Xi Zhang, Yingxu Wang, Xinhao Yi, Guibin Zhang, Yi Xu, Bin Wu, Siwei Liu, Zihao Li, and 1 others. 2025a. A comprehensive survey of self-evolving ai agents: A new paradigm bridging foundation models and lifelong agentic systems. <i>arXiv preprint arXiv:2508.07407</i> .	926
872		927
873		928
874		929
875		930
876		931
877	Tianqing Fang, Hongming Zhang, Zhisong Zhang, Kaixin Ma, Wenhao Yu, Haitao Mi, and Dong Yu. 2025b. WebEvolver: Enhancing web agent self-improvement with co-evolving world model . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 8970–8986, Suzhou, China. Association for Computational Linguistics.	932
878		933
879		934
880		935
881		936
882		937
883		938
884		939
885	Lang Feng, Zhenghai Xue, Tingcong Liu, and Bo An. 2025. Group-in-group policy optimization for LLM agent training . In <i>The Thirty-ninth Annual Conference on Neural Information Processing Systems</i> .	940
886		941
887		942
888		943
889	Saumya Gandhi, Ritu Gala, Vijay Viswanathan, Tongshuang Wu, and Graham Neubig. 2024. Better synthetic data by retrieving and transforming existing datasets . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 6453–6466, Bangkok, Thailand. Association for Computational Linguistics.	944
890		945
891		946
892		947
893		948
894		949
895		950
896	Huan-ang Gao, Jiayi Geng, Wenyue Hua, Mengkang Hu, Xinzhe Juan, Hongzhang Liu, Shilong Liu, Jiahao Qiu, Xuan Qi, Yiran Wu, and 1 others. 2025a. A survey of self-evolving agents: On path to artificial super intelligence. <i>arXiv preprint arXiv:2507.21046</i> .	951
897		952
898		953
899		954
900		955
901	Jiaxuan Gao, Wei Fu, Minyang Xie, Shusheng Xu, Chuyi He, Zhiyu Mei, Banghua Zhu, and Yi Wu. 2025b. Beyond ten turns: Unlocking long-horizon agentic search with large-scale asynchronous rl. <i>arXiv preprint arXiv:2508.07976</i> .	956
902		957
903		958
904		959
905		960
906	Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, and 1 others. 2024. A survey on llm-as-a-judge. <i>arXiv preprint arXiv:2411.15594</i> .	961
907		962
908		963
909		964
910		965
		966
	Yu Gu, Kai Zhang, Yuting Ning, Boyuan Zheng, Boyu Gou, Tianci Xue, Cheng Chang, Sanjari Srivastava, Yanan Xie, Peng Qi, Huan Sun, and Yu Su. 2025. Is your LLM secretly a world model of the internet? model-based planning for web agents . <i>Transactions on Machine Learning Research</i> .	911
		912
		913
		914
		915
		916
	Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025a. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. <i>arXiv preprint arXiv:2501.12948</i> .	917
		918
		919
		920
		921
		922
	Jiacheng Guo, Ling Yang, Peter Chen, Qixin Xiao, Yinjie Wang, Xinzhe Juan, Jiahao Qiu, Ke Shen, and Mengdi Wang. 2025b. Genenv: Difficulty-aligned co-evolution between llm agents and environment simulators . <i>Preprint</i> , arXiv:2512.19682.	923
		924
		925
		926
		927
	Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V. Chawla, Olaf Wiest, and Xiangliang Zhang. 2024. Large language model based multi-agents: A survey of progress and challenges . In <i>Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence, IJCAI-24</i> , pages 8048–8057. International Joint Conferences on Artificial Intelligence Organization. Survey Track.	928
		929
		930
		931
		932
		933
		934
		935
		936
	David Ha and Jürgen Schmidhuber. 2018. World models. <i>arXiv preprint arXiv:1803.10122</i> .	937
		938
	Hongliang He, Wenlin Yao, Kaixin Ma, Wenhao Yu, Hongming Zhang, Tianqing Fang, Zhenzhong Lan, and Dong Yu. 2025. OpenWebVoyager: Building multimodal web agents via iterative real-world exploration, feedback and optimization . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 27545–27564, Vienna, Austria. Association for Computational Linguistics.	939
		940
		941
		942
		943
		944
		945
		946
		947
	Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. <i>arXiv preprint arXiv:1503.02531</i> .	948
		949
		950
	Haoyang Hong, Jiajun Yin, Yuan Wang, Jingnan Liu, Zhe Chen, Ailing Yu, Ji Li, Zhiling Ye, Hansong Xiao, Yefei Chen, and 1 others. 2025. Multi-agent deep research: Training multi-agent systems with m-grpo. <i>arXiv preprint arXiv:2511.13288</i> .	951
		952
		953
		954
		955
	Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, and 1 others. 2023. Metagpt: Meta programming for a multi-agent collaborative framework. In <i>The Twelfth International Conference on Learning Representations</i> .	956
		957
		958
		959
		960
		961
		962
	Zhenyu Hou, Ziniu Hu, Yujiang Li, Rui Lu, Jie Tang, and Yuxiao Dong. 2025. TreeRL: LLM reinforcement learning with on-policy tree search . In <i>Proceedings of the 63rd Annual Meeting of the Association</i>	963
		964
		965
		966

967					
968					
969					
970	Minda Hu, Tianqing Fang, Jianshu Zhang, Jun-Yu Ma,				
971	Zhisong Zhang, Jingyan Zhou, Hongming Zhang,				
972	Haitao Mi, Dong Yu, and Irwin King. 2025a. Web-				
973	CoT: Enhancing web agent reasoning by reconstruct-				
974	ing chain-of-thought in reflection, branching, and				
975	rollback . In <i>Findings of the Association for Com-</i>				
976	<i>putational Linguistics: EMNLP 2025</i> , pages 5155–				
977	5173, Suzhou, China. Association for Computational				
978	Linguistics.				
979	Xueyu Hu, Tao Xiong, Biao Yi, Zishu Wei, Ruix-				
980	uan Xiao, Yurun Chen, Jiasheng Ye, Meiling Tao,				
981	Xiangxin Zhou, Ziyu Zhao, and 1 others. 2025b.				
982	Os agents: A survey on mllm-based agents for				
983	general computing devices use. <i>arXiv preprint</i>				
984	<i>arXiv:2508.04482</i> .				
985	Jiaxin Huang, Shixiang Gu, Le Hou, Yuexin Wu, Xuezhi				
986	Wang, Hongkun Yu, and Jiawei Han. 2023. Large				
987	language models can self-improve . In <i>Proceedings</i>				
988	<i>of the 2023 Conference on Empirical Methods in Natu-</i>				
989	<i>ral Language Processing</i> , pages 1051–1068, Singa-				
990	apore. Association for Computational Linguistics.				
991	Xu Huang, Weiwen Liu, Xiaolong Chen, Xingmei				
992	Wang, Hao Wang, Defu Lian, Yasheng Wang, Ruim-				
993	ing Tang, and Enhong Chen. 2024. Understanding				
994	the planning of llm agents: A survey. <i>arXiv preprint</i>				
995	<i>arXiv:2402.02716</i> .				
996	Dongfu Jiang, Yi Lu, Zhuofeng Li, Zhiheng Lyu, Ping				
997	Nie, Haozhe Wang, Alex Su, Hui Chen, Kai Zou,				
998	Chao Du, and 1 others. 2025. Verltool: Towards				
999	holistic agentic reinforcement learning with tool use.				
1000	<i>arXiv preprint arXiv:2509.01055</i> .				
1001	Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Ser-				
1002	can O Arik, Dong Wang, Hamed Zamani, and Jiawei				
1003	Han. 2025. Search-r1: Training LLMs to reason and				
1004	leverage search engines with reinforcement learning .				
1005	In <i>Second Conference on Language Modeling</i> .				
1006	Minki Kang, Jongwon Jeong, Seanie Lee, Jaewoong				
1007	Cho, and Sung Ju Hwang. 2025. Distilling llm agent				
1008	into small models with retrieval and code tools. <i>arXiv</i>				
1009	<i>preprint arXiv:2505.17612</i> .				
1010	Imene Kerboua, Sahar Omid Shayegan, Megh Thakkar,				
1011	Xing Han Lù, Léo Boisvert, Massimo Caccia, Jérémy				
1012	Espinas, Alexandre Aussem, Véronique Eglin, and				
1013	Alexandre Lacoste. 2025. Focusagent: Simple yet				
1014	effective ways of trimming the large context of web				
1015	agents. <i>arXiv preprint arXiv:2510.03204</i> .				
1016	Yoon Kim and Alexander M. Rush. 2016. Sequence-				
1017	level knowledge distillation . In <i>Proceedings of the</i>				
1018	<i>2016 Conference on Empirical Methods in Natu-</i>				
1019	<i>ral Language Processing</i> , pages 1317–1327, Austin,				
1020	Texas. Association for Computational Linguistics.				
1021	Jing Yu Koh, Stephen Marcus McAleer, Daniel Fried,				
1022	and Ruslan Salakhutdinov. 2025. Tree search for				
	language model agents . <i>Transactions on Machine</i>				1023
	<i>Learning Research</i> .				1024
	Hanyu Lai, Xiao Liu, Iat Long Iong, Shuntian Yao, Yux-				1025
	uan Chen, Pengbo Shen, Hao Yu, Hanchen Zhang,				1026
	Xiaohan Zhang, Yuxiao Dong, and 1 others. 2024.				1027
	Autowebglm: A large language model-based web				1028
	navigating agent. In <i>Proceedings of the 30th ACM</i>				1029
	<i>SIGKDD Conference on Knowledge Discovery and</i>				1030
	<i>Data Mining</i> , pages 5295–5306.				1031
	Nathan Lambert, Jacob Morrison, Valentina Pyatkin,				1032
	Shengyi Huang, Hamish Ivison, Faeze Brahman,				1033
	Lester James V Miranda, Alisa Liu, Nouha Dziri,				1034
	Shane Lyu, and 1 others. 2024. Tulu 3: Pushing fron-				1035
	tiers in open language model post-training. <i>arXiv</i>				1036
	<i>preprint arXiv:2411.15124</i> .				1037
	Dongjun Lee, Juyong Lee, Kyuyoung Kim, Jihoon Tack,				1038
	Jinwoo Shin, Yee Whye Teh, and Kimin Lee. 2025.				1039
	Learning to contextualize web pages for enhanced				1040
	decision making by LLM agents . In <i>The Thirteenth</i>				1041
	<i>International Conference on Learning Representa-</i>				1042
	<i>tions</i> .				1043
	Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas				1044
	Mesnard, Johan Ferret, Kellie Ren Lu, Colton Bishop,				1045
	Ethan Hall, Victor Carbune, Abhinav Rastogi, and				1046
	Sushant Prakash. 2024. RLAIF vs. RLHF: Scaling				1047
	reinforcement learning from human feedback with AI				1048
	feedback . In <i>Proceedings of the 41st International</i>				1049
	<i>Conference on Machine Learning</i> , volume 235 of				1050
	<i>Proceedings of Machine Learning Research</i> , pages				1051
	26874–26901. PMLR.				1052
	Jan Leike, David Krueger, Tom Everitt, Miljan Martic,				1053
	Vishal Maini, and Shane Legg. 2018. Scalable agent				1054
	alignment via reward modeling: a research direction.				1055
	<i>arXiv preprint arXiv:1811.07871</i> .				1056
	Sergey Levine, Aviral Kumar, George Tucker, and Justin				1057
	Fu. 2020. Offline reinforcement learning: Tutorial,				1058
	review, and perspectives on open problems. <i>arXiv</i>				1059
	<i>preprint arXiv:2005.01643</i> .				1060
	Chenliang Li, Adel Elmahdy, Alex Boyd, Zhongruo				1061
	Wang, Alfredo Garcia, Parminder Bhatia, Taha Kass-				1062
	Hout, Cao Xiao, and Mingyi Hong. 2025a. Stp-				1063
	ppo: Stabilized off-policy proximal policy optimiza-				1064
	tion for multi-turn agents training . <i>arXiv preprint</i>				1065
	<i>arXiv:2511.20718</i> .				1066
	Kuan Li, Zhongwang Zhang, Huifeng Yin, Liwen				1067
	Zhang, Litu Ou, Jialong Wu, Wenbiao Yin, Baixuan				1068
	Li, Zhengwei Tao, Xinyu Wang, and 1 others. 2025b.				1069
	Websailor: Navigating super-human reasoning for				1070
	web agent . <i>arXiv preprint arXiv:2507.02592</i> .				1071
	Shiyu Li, Yang Tang, Yifan Wang, Peiming Li, and				1072
	Xi Chen. 2025c. Reseek: A self-correcting frame-				1073
	work for search agents with instructive rewards .				1074
	<i>arXiv preprint arXiv:2510.00568</i> .				1075
	Wei Li, William E Bishop, Alice Li, Christopher Rawles,				1076
	Folawiyo Campbell-Ajala, Divya Tyamagundlu, and				1077
	Oriana Riva. 2024. On the effects of data scale on				1078

1079	ui control agents. <i>Advances in Neural Information Processing Systems</i> , 37:92130–92154.	Miao Lu, Weiwei Sun, Weihua Du, Zhan Ling, Xuesong Yao, Kang Liu, and Jiecao Chen. 2025b. Scaling llm multi-turn rl with end-to-end summarization-based context management. <i>arXiv preprint arXiv:2510.06727</i> .	1133
1080			1134
1081	Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2024. Let’s verify step by step . In <i>The Twelfth International Conference on Learning Representations</i> .	Rui Lu, Zhenyu Hou, Zihan Wang, Hanchen Zhang, Xiao Liu, Yujiang Li, Shi Feng, Jie Tang, and Yuxiao Dong. 2025c. Deepdive: Advancing deep search agents with knowledge graphs and multi-turn rl. <i>arXiv preprint arXiv:2509.10446</i> .	1135
1082			1136
1083			1137
1084			1138
1085			1139
1086	Zongyu Lin, Yao Tang, Xingcheng Yao, Da Yin, Ziniu Hu, Yizhou Sun, and Kai-Wei Chang. 2025a. QLASS: Boosting language agent inference via q-guided stepwise search . In <i>Forty-second International Conference on Machine Learning</i> .	Xing Han Lù, Gaurav Kamath, Marius Mosbach, and Siva Reddy. 2025. Build the web for agents, not agents for the web. <i>arXiv preprint arXiv:2506.10953</i> .	1140
1087			1141
1088			1142
1089			1143
1090			1144
1091	Zongyu Lin, Yao Tang, Xingcheng Yao, Da Yin, Ziniu Hu, Yizhou Sun, and Kai-Wei Chang. 2025b. QLASS: Boosting language agent inference via q-guided stepwise search . In <i>Proceedings of the 42nd International Conference on Machine Learning</i> , volume 267 of <i>Proceedings of Machine Learning Research</i> , pages 37942–37958. PMLR.	Xing Han Lu, Zdeněk Kasner, and Siva Reddy. 2024a. WebLINX: Real-world website navigation with multi-turn dialogue . In <i>Proceedings of the 41st International Conference on Machine Learning</i> , volume 235 of <i>Proceedings of Machine Learning Research</i> , pages 33007–33056. PMLR.	1145
1092			1146
1093			1147
1094			1148
1095			1149
1096			1150
1097			1151
1098	Aixin Liu, Aoxue Mei, Bangcai Lin, Bing Xue, Bingxuan Wang, Bingzheng Xu, Bochao Wu, Bowei Zhang, Chaofan Lin, Chen Dong, and 1 others. 2025a. Deepseek-v3.2: Pushing the frontier of open large language models. <i>arXiv preprint arXiv:2512.02556</i> .	Yadong Lu, Jianwei Yang, Yelong Shen, and Ahmed Awadallah. 2024b. Omniparser for pure vision based gui agent. <i>arXiv preprint arXiv:2408.00203</i> .	1152
1099			1153
1100			1154
1101			1155
1102			1156
1103	Junteng Liu, Yunji Li, Chi Zhang, Jingyang Li, Aili Chen, Ke Ji, Weiyu Cheng, Zijia Wu, Chengyu Du, Qidi Xu, and 1 others. 2025b. Webexplorer: Explore and evolve for training long-horizon web agents. <i>arXiv preprint arXiv:2509.06501</i> .	Junyu Luo, Weizhi Zhang, Ye Yuan, Yusheng Zhao, Junwei Yang, Yiyang Gu, Bohan Wu, Binqi Chen, Ziyue Qiao, Qingqing Long, and 1 others. 2025a. Large language model agent: A survey on methodology, applications and challenges. <i>arXiv preprint arXiv:2503.21460</i> .	1157
1104			1158
1105			1159
1106			1160
1107			1161
1108	Junwei Liu, Kaixin Wang, Yixuan Chen, Xin Peng, Zhenpeng Chen, Lingming Zhang, and Yiling Lou. 2024. Large language model-based agents for software engineering: A survey. <i>arXiv preprint arXiv:2409.02977</i> .	Kun Luo, Hongjin Qian, Zheng Liu, Ziyi Xia, Shitao Xiao, Siqi Bao, Jun Zhao, and Kang Liu. 2025b. Inflow: Reinforcing search agent via reward density optimization. <i>arXiv preprint arXiv:2510.26575</i> .	1162
1109			1163
1110			1164
1111			1165
1112			1166
1113	Keliang Liu, Dingkan Yang, Ziyun Qian, Weijie Yin, Yuchi Wang, Hongsheng Li, Jun Liu, Peng Zhai, Yang Liu, and Lihua Zhang. 2025c. Reinforcement learning meets large language models: A survey of advancements and applications across the llm lifecycle. <i>arXiv preprint arXiv:2509.16679</i> .	Yuanjie Lyu, Chengyu Wang, Jun Huang, and Tong Xu. 2025. From correction to mastery: Reinforced distillation of large language model agents. <i>arXiv preprint arXiv:2509.14257</i> .	1167
1114			1168
1115			1169
1116			1170
1117			1171
1118			1172
1119	Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. <i>ACM computing surveys</i> , 55(9):1–35.	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, and 1 others. 2023. Self-refine: Iterative refinement with self-feedback. <i>Advances in Neural Information Processing Systems</i> , 36:46534–46594.	1173
1120			1174
1121			1175
1122			1176
1123			1177
1124	Xiaoqian Liu, Ke Wang, Yuchuan Wu, Fei Huang, Yongbin Li, Junge Zhang, and Jianbin Jiao. 2025d. Agentic reinforcement learning with implicit step rewards. <i>arXiv preprint arXiv:2509.19199</i> .	Lingrui Mei, Jiayu Yao, Yuyao Ge, Yiwei Wang, Baolong Bi, Yujun Cai, Jiazhi Liu, Mingyu Li, Zhong-Zhi Li, Duzhen Zhang, and 1 others. 2025. A survey of context engineering for large language models. <i>arXiv preprint arXiv:2507.13334</i> .	1178
1125			1179
1126			1180
1127			1181
1128	Han Lu, Zichen Liu, Shaopan Xiong, Yancheng He, Wei Gao, Yanan Wu, Weixun Wang, Jiashun Liu, Yang Li, Haizhou Zhao, and 1 others. 2025a. Part ii: Roll flash—accelerating rlvr and agentic training with asynchrony. <i>arXiv preprint arXiv:2510.11345</i> .	Grégoire Mialon, Roberto Dessi, Maria Lomeli, Christoforos Nalmpantis, Ramakanth Pasunuru, Roberta Raileanu, Baptiste Roziere, Timo Schick, Jane Dwivedi-Yu, Asli Celikyilmaz, Edouard Grave, Yann LeCun, and Thomas Scialom. 2023. Augmented language models: a survey . <i>Transactions on Machine Learning Research</i> . Survey Certification.	1182
1129			1183
1130			1184
1131			1185
1132			1186

1412	Zhen Tan, Dawei Li, Song Wang, Alimohammad	Le Sellier de Chezelles, Nicolas Gontier, Miguel	1468
1413	Beigi, Bohan Jiang, Amrita Bhattacharjee, Man-	Muñoz-Mármol, Sahar Omidi Shayegan, Stefania	1469
1414	sooreh Karami, Jundong Li, Lu Cheng, and Huan Liu.	Raimondo, and 1 others. 2025. How to train your llm	1470
1415	2024. Large language models for data annotation and	web agent: A statistical diagnosis. <i>arXiv preprint</i>	1471
1416	synthesis: A survey . In <i>Proceedings of the 2024 Con-</i>	<i>arXiv:2507.04103</i> .	1472
1417	<i>ference on Empirical Methods in Natural Language</i>		
1418	<i>Processing</i> , pages 930–957, Miami, Florida, USA.	Guoqing Wang, Sunhao Dai, Guangze Ye, Zeyu Gan,	1473
1419	Association for Computational Linguistics.	Wei Yao, Yong Deng, Xiaofeng Wu, and Zhenzhe	1474
		Ying. 2025a. Information gain-based policy optimiza-	1475
1420	Fei Tang, Haolei Xu, Hang Zhang, Siqi Chen, Xingyu	tion: A simple and effective approach for multi-turn	1476
1421	Wu, Yongliang Shen, Wenqi Zhang, Guiyang Hou,	llm agents. <i>arXiv preprint arXiv:2510.14967</i> .	1477
1422	Zeqi Tan, Yuchen Yan, and 1 others. 2025a. A sur-		
1423	vey on (m) llm-based gui agents. <i>arXiv preprint</i>	Haoming Wang, Haoyang Zou, Huatong Song, Jiazhan	1478
1424	<i>arXiv:2504.13865</i> .	Feng, Junjie Fang, Junting Lu, Longxiang Liu, Qinyu	1479
		Luo, Shihao Liang, Shijue Huang, and 1 others.	1480
1425	Qiaoyu Tang, Hao Xiang, Le Yu, Bowen Yu, Yaojie Lu,	2025b. Ui-tars-2 technical report: Advancing gui	1481
1426	Xianpei Han, Le Sun, WenJuan Zhang, Pengbo Wang,	agent with multi-turn reinforcement learning. <i>arXiv</i>	1482
1427	Shixuan Liu, and 1 others. 2025b. Beyond turn limits:	<i>preprint arXiv:2509.02544</i> .	1483
1428	Training deep search agents with dynamic context		
1429	window. <i>arXiv preprint arXiv:2510.08276</i> .	Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao	1484
		Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang,	1485
1430	Xiangru Tang, Tianrui Qin, Tianhao Peng, Ziyang Zhou,	Xu Chen, Yankai Lin, and 1 others. 2024a. A survey	1486
1431	Daniel Shao, Tingting Du, Xinming Wei, Peng Xia,	on large language model based autonomous agents.	1487
1432	Fang Wu, He Zhu, and 1 others. 2025c. Agent kb:	<i>Frontiers of Computer Science</i> , 18(6):186345.	1488
1433	Leveraging cross-domain experience for agentic prob-		
1434	lem solving. <i>arXiv preprint arXiv:2507.06229</i> .	Liyuan Wang, Xingxing Zhang, Hang Su, and Jun Zhu.	1489
		2024b. A comprehensive survey of continual learn-	1490
1435	Zhengwei Tao, Jialong Wu, Wenbiao Yin, Junkai	ing: Theory, method and application. <i>IEEE transac-</i>	1491
1436	Zhang, Baixuan Li, Haiyang Shen, Kuan Li, Li-	<i>tions on pattern analysis and machine intelligence</i> ,	1492
1437	wen Zhang, Xinyu Wang, Yong Jiang, and 1 others.	46(8):5362–5383.	1493
1438	2025. Webshaper: Agentic data synthesizing via		
1439	information-seeking formalization. <i>arXiv preprint</i>	Rui Wang, Ce Zhang, Jun-Yu Ma, Jianshu Zhang, Hon-	1494
1440	<i>arXiv:2507.15061</i> .	gru Wang, Yi Chen, Boyang Xue, Tianqing Fang,	1495
		Zhisong Zhang, Hongming Zhang, and 1 others.	1496
1441	Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen,	2025c. Explore to evolve: Scaling evolved aggrega-	1497
1442	Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru	tion logic via proactive online exploration for deep	1498
1443	Chen, Yuankun Chen, Yutian Chen, and 1 others.	research agents. <i>arXiv preprint arXiv:2510.14438</i> .	1499
1444	2025. Kimi k2: Open agentic intelligence. <i>arXiv</i>		
1445	<i>preprint arXiv:2507.20534</i> .	Ruiyi Wang and Prithviraj Ammanabrolu. 2025. A prac-	1500
		titioner’s guide to multi-turn agentic reinforcement	1501
1446	Guiyao Tie, Zeli Zhao, Dingjie Song, Fuyang Wei,	learning. <i>arXiv preprint arXiv:2510.01132</i> .	1502
1447	Rong Zhou, Yurou Dai, Wen Yin, Zhejiang Yang,		
1448	Jiangyue Yan, Yao Su, and 1 others. 2025. Large	Shuhe Wang, Shengyu Zhang, Jie Zhang, Runyi Hu,	1503
1449	language models post-training: Surveying techni-	Xiaoya Li, Tianwei Zhang, Jiwei Li, Fei Wu, Guoyin	1504
1450	ques from alignment to reasoning. <i>arXiv preprint</i>	Wang, and Eduard Hovy. 2024c. Reinforcement	1505
1451	<i>arXiv:2503.06072</i> .	learning enhanced llms: A survey. <i>arXiv preprint</i>	1506
		<i>arXiv:2412.10400</i> .	1507
1452	Brandon Trabucco, Gunnar Sigurdsson, Robinson Pi-	Xingyao Wang, Yangyi Chen, Lifan Yuan, Yizhe Zhang,	1508
1453	ramuthu, and Ruslan Salakhutdinov. 2025. Insta:	Yunzhu Li, Hao Peng, and Heng Ji. 2024d. Exe-	1509
1454	Towards internet-scale training for agents. <i>arXiv</i>	cutable code actions elicit better LLM agents . In	1510
1455	<i>preprint arXiv:2502.06776</i> .	<i>Forty-first International Conference on Machine</i>	1511
		<i>Learning</i> .	1512
1456	Khanh-Tung Tran, Dung Dao, Minh-Duong Nguyen,		
1457	Quoc-Viet Pham, Barry O’Sullivan, and Hoang D	Yanlin Wang, Wanjun Zhong, Yanxian Huang, Ensheng	1513
1458	Nguyen. 2025. Multi-agent collaboration mech-	Shi, Min Yang, Jiachi Chen, Hui Li, Yuchi Ma, Qianx-	1514
1459	anisms: A survey of llms. <i>arXiv preprint</i>	iang Wang, and Zibin Zheng. 2025d. Agents in soft-	1515
1460	<i>arXiv:2501.06322</i> .	ware engineering: Survey, landscape, and vision. <i>Au-</i>	1516
		<i>tomated Software Engineering</i> , 32(2):1–36.	1517
1461	Jonathan Uesato, Nate Kushman, Ramana Kumar, Fran-	Yiming Wang, Da Yin, Yuedong Cui, Ruichen Zheng,	1518
1462	cis Song, Noah Siegel, Lisa Wang, Antonia Creswell,	Zhiqian Li, Zongyu Lin, Di Wu, Xueqing Wu,	1519
1463	Geoffrey Irving, and Irina Higgins. 2022. Solv-	Chenchen Ye, Yu Zhou, and 1 others. 2025e. Llms as	1520
1464	ing math word problems with process-and outcome-	scalable, general-purpose simulators for evolving dig-	1521
1465	based feedback. <i>arXiv preprint arXiv:2211.14275</i> .	ital agent training. <i>arXiv preprint arXiv:2510.14969</i> .	1522
1466	Dheeraj Vattikonda, Santhoshi Ravichandran, Emiliano		
1467	Penaloza, Hadi Nekoei, Megh Thakkar, Thibault		

1523	Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. Self-instruct: Aligning language models with self-generated instructions . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 13484–13508, Toronto, Canada. Association for Computational Linguistics.	1580
1524		1581
1525		1582
1526		1583
1527		
1528		1584
1529		1585
1530		1586
		1587
1531	Zihan Wang, Kangrui Wang, Qineng Wang, Pingyue Zhang, Linjie Li, Zhengyuan Yang, Xing Jin, Kefan Yu, Minh Nhat Nguyen, Licheng Liu, and 1 others. 2025f. Ragen: Understanding self-evolution in llm agents via multi-turn reinforcement learning. <i>arXiv preprint arXiv:2504.20073</i> .	1588
1532		1589
1533		
1534		1590
1535		1591
1536		1592
		1593
1537	Ziliang Wang, Kang An, Xuhui Zheng, Faqiang Qian, Weikun Zhang, Cijun Ouyang, Jialu Cai, Yuhang Wang, and Yichao Wu. 2025g. Erase to improve: Erasable reinforcement learning for search-augmented llms. <i>arXiv preprint arXiv:2510.00861</i> .	1594
1538		1595
1539		
1540		1596
1541		1597
		1598
1542	Yuxiang Wei, Olivier Duchenne, Jade Copet, Quentin Carbonneaux, LINGMING ZHANG, Daniel Fried, Gabriel Synnaeve, Rishabh Singh, and Sida Wang. 2025. SWE-RL: Advancing LLM reasoning via reinforcement learning on open software evolution . In <i>The Thirty-ninth Annual Conference on Neural Information Processing Systems</i> .	1599
1543		1600
1544		1601
1545		
1546		1602
1547		1603
1548		1604
		1605
1549	Feijie Wu, Weiwu Zhu, Yuxiang Zhang, Soumya Chatterjee, Jiarong Zhu, Fan Mo, Rodin Luo, and Jing Gao. 2025a. Portool: Tool-use llm training with rewarded tree. <i>arXiv preprint arXiv:2510.26020</i> .	1606
1550		1607
1551		1608
1552		
1553	Jialong Wu, Baixuan Li, Runnan Fang, Wenbiao Yin, Liwen Zhang, Zhenglin Wang, Zhengwei Tao, Dingchu Zhang, Zekun Xi, Xiangru Tang, Yong Jiang, Pengjun Xie, Fei Huang, and Jingren Zhou. 2025b. Webdancer: Towards autonomous information seeking agency . In <i>The Thirty-ninth Annual Conference on Neural Information Processing Systems</i> .	1609
1554		1610
1555		1611
1556		1612
1557		1613
1558		1614
1559		1615
		1616
1560	Peilin Wu, Mian Zhang, Kun Wan, Wentian Zhao, Kaiyu He, Xinya Du, and Zhiyu Chen. 2025c. Hiprag: Hierarchical process rewards for efficient agentic retrieval augmented generation. <i>arXiv preprint arXiv:2510.07794</i> .	1617
1561		1618
1562		1619
1563		1620
1564		1621
		1622
1565	Peilin Wu, Mian Zhang, Xinlu Zhang, Xinya Du, and Zhiyu Chen. 2025d. Search wisely: Mitigating sub-optimal agentic searches by reducing uncertainty . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 19734–19745, Suzhou, China. Association for Computational Linguistics.	1623
1566		
1567		1624
1568		1625
1569		1626
1570		1627
1571		1628
		1629
1572	Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, and 1 others. 2024. Autogen: Enabling next-gen llm applications via multi-agent conversations. In <i>First Conference on Language Modeling</i> .	1630
1573		1631
1574		1632
1575		1633
1576		1634
1577		1635
1578	Xixi Wu, Kuan Li, Yida Zhao, Liwen Zhang, Litu Ou, Huifeng Yin, Zhongwang Zhang, Xinmiao	
1579		
		1630
		1631
		1632
		1633
		1634
		1635
		1636
		1637
		1638
		1639
		1640
		1641
		1642
		1643
		1644
		1645
		1646
		1647
		1648
		1649
		1650
		1651
		1652
		1653
		1654
		1655
		1656
		1657
		1658
		1659
		1660
		1661
		1662
		1663
		1664
		1665
		1666
		1667
		1668
		1669
		1670
		1671
		1672
		1673
		1674
		1675
		1676
		1677
		1678
		1679
		1680
		1681
		1682
		1683
		1684
		1685
		1686
		1687
		1688
		1689
		1690
		1691
		1692
		1693
		1694
		1695
		1696
		1697
		1698
		1699
		1700
		1701
		1702
		1703
		1704
		1705
		1706
		1707
		1708
		1709
		1710
		1711
		1712
		1713
		1714
		1715
		1716
		1717
		1718
		1719
		1720
		1721
		1722
		1723
		1724
		1725
		1726
		1727
		1728
		1729
		1730
		1731
		1732
		1733
		1734
		1735
		1736
		1737
		1738
		1739
		1740
		1741
		1742
		1743
		1744
		1745
		1746
		1747
		1748
		1749
		1750
		1751
		1752
		1753
		1754
		1755
		1756
		1757
		1758
		1759
		1760
		1761
		1762
		1763
		1764
		1765
		1766
		1767
		1768
		1769
		1770
		1771
		1772
		1773
		1774
		1775
		1776
		1777
		1778
		1779
		1780
		1781
		1782
		1783
		1784
		1785
		1786
		1787
		1788
		1789
		1790
		1791
		1792
		1793
		1794
		1795
		1796
		1797
		1798
		1799
		1800

1636	Sikuan Yan, Xiufeng Yang, Zuchao Huang, Ercong Nie,	multimodal large language models. <i>National Science</i>	1692
1637	Zifeng Ding, Zonggen Li, Xiaowen Ma, Kristian Ker-	<i>Review</i> , 11(12):nwae403.	1693
1638	sting, Jeff Z Pan, Hinrich Schütze, and 1 others. 2025.		
1639	Memory-r1: Enhancing large language model agents	Yuanqing Yu, Zhefan Wang, Weizhi Ma, Shuai Wang,	1694
1640	to manage and utilize memories via reinforcement	Chuhan Wu, Zhiqiang Guo, and Min Zhang. 2024.	1695
1641	learning. <i>arXiv preprint arXiv:2508.19828</i> .	Steptool: Enhancing multi-step tool usage in llms	1696
		through step-grained reinforcement learning. <i>arXiv</i>	1697
1642	Jianwei Yang, Hao Zhang, Feng Li, Xueyan Zou, Chun-	<i>preprint arXiv:2410.07745</i> .	1698
1643	yuan Li, and Jianfeng Gao. 2023. Set-of-mark		
1644	prompting unleashes extraordinary visual grounding	Zhaochen Yu, Ling Yang, Jiaru Zou, Shuicheng Yan,	1699
1645	in gpt-4v. <i>arXiv preprint arXiv:2310.11441</i> .	and Mengdi Wang. 2025. Demystifying reinforce-	1700
		ment learning in agentic reasoning. <i>arXiv preprint</i>	1701
1646	Ke Yang, Yao Liu, Sapana Chaudhary, Rasool Fakoor,	<i>arXiv:2510.11701</i> .	1702
1647	Pratik Chaudhari, George Karypis, and Huzefa Rang-		
1648	wala. 2025a. Agentoccam: A simple yet strong base-	Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho,	1703
1649	line for LLM-based web agents . In <i>The Thirteenth</i>	Xian Li, Sainbayar Sukhbaatar, Jing Xu, and Jaso-	1704
1650	<i>International Conference on Learning Representa-</i>	n E Weston. 2024. Self-rewarding language mod-	1705
1651	<i>tions</i> .	els . In <i>Forty-first International Conference on Ma-</i>	1706
		<i>chine Learning</i> .	1707
1652	Ruihan Yang, Fanghua Ye, Jian Li, Siyu Yuan, Yikai	Aohan Zeng, Mingdao Liu, Rui Lu, Bowen Wang, Xiao	1708
1653	Zhang, Zhaopeng Tu, Xiaolong Li, and Deqing Yang.	Liu, Yuxiao Dong, and Jie Tang. 2024. AgentTun-	1709
1654	2025b. The lighthouse of language: Enhancing LLM	ing: Enabling generalized agent abilities for LLMs .	1710
1655	agents via critique-guided improvement . In <i>The</i>	In <i>Findings of the Association for Computational</i>	1711
1656	<i>Thirty-ninth Annual Conference on Neural Informa-</i>	<i>Linguistics: ACL 2024</i> , pages 3053–3077, Bangkok,	1712
1657	<i>tion Processing Systems</i> .	Thailand. Association for Computational Linguistics.	1713
1658	Zonghan Yang, Peng Li, Ming Yan, Ji Zhang, Fei Huang,	Yunpeng Zhai, Shuchang Tao, Cheng Chen, Anni Zou,	1714
1659	and Yang Liu. 2024. React meets actre: When lan-	Ziqian Chen, Qingxu Fu, Shinji Mai, Li Yu, Jiaji	1715
1660	guage agents enjoy training data autonomy. <i>arXiv</i>	Deng, Zouying Cao, and 1 others. 2025. Agente-	1716
1661	<i>preprint arXiv:2403.14589</i> .	volver: Towards efficient self-evolving agent system .	1717
		<i>arXiv preprint arXiv:2511.10395</i> .	1718
1662	Shunyu Yao, Howard Chen, John Yang, and Karthik	Chaoyun Zhang, Shilin He, Liqun Li, Si Qin, Yu Kang,	1719
1663	Narasimhan. 2022. Webshop: Towards scalable real-	Qingwei Lin, Saravan Rajmohan, and Dongmei	1720
1664	world web interaction with grounded language agents.	Zhang. 2025a. Api agents vs. gui agents: Divergence	1721
1665	<i>Advances in Neural Information Processing Systems</i> ,	and convergence . <i>arXiv preprint arXiv:2503.11069</i> .	1722
1666	35:20744–20757.		
1667	Rui Ye, Zhongwang Zhang, Kuan Li, Huifeng Yin,	Chaoyun Zhang, Shilin He, Jiaxu Qian, Bowen Li,	1723
1668	Zhengwei Tao, Yida Zhao, Liangcai Su, Liwen	Liqun Li, Si Qin, Yu Kang, Minghua Ma, Guyue	1724
1669	Zhang, Zile Qiao, Xinyu Wang, and 1 others.	Liu, Qingwei Lin, and 1 others. 2024a. Large lan-	1725
1670	2025. Agentfold: Long-horizon web agents with	guage model-brained gui agents: A survey . <i>arXiv</i>	1726
1671	proactive context management . <i>arXiv preprint</i>	<i>preprint arXiv:2411.18279</i> .	1727
1672	<i>arXiv:2510.24699</i> .		
1673	Asaf Yehudai, Lilach Eden, Alan Li, Guy Uziel, Yilun	Chaoyun Zhang, Liqun Li, Shilin He, Xu Zhang,	1728
1674	Zhao, Roy Bar-Haim, Arman Cohan, and Michal	Bo Qiao, Si Qin, Minghua Ma, Yu Kang, Qing-	1729
1675	Shmueli-Scheuer. 2025. Survey on evaluation of llm-	wei Lin, Saravan Rajmohan, Dongmei Zhang, and	1730
1676	based agents. <i>arXiv preprint arXiv:2503.16416</i> .	Qi Zhang. 2025b. UFO: A UI-focused agent for win-	1731
		dows OS interaction . In <i>Proceedings of the 2025</i>	1732
1677	Howard Yen, Ashwin Paranjape, Mengzhou Xia, The-	<i>Conference of the Nations of the Americas Chap-</i>	1733
1678	jas Venkatesh, Jack Hessel, Danqi Chen, and Yuhao	<i>ter of the Association for Computational Linguistics:</i>	1734
1679	Zhang. 2025. Lost in the maze: Overcoming con-	<i>Human Language Technologies (Volume 1: Long Pa-</i>	1735
1680	text limitations in long-horizon agentic search . <i>arXiv</i>	<i>pers)</i> , pages 597–622, Albuquerque, New Mexico.	1736
1681	<i>preprint arXiv:2510.18939</i> .	Association for Computational Linguistics.	1737
1682	Da Yin, Faeze Brahman, Abhilasha Ravichander, Khy-	Guibin Zhang, Hejia Geng, Xiaohang Yu, Zhenfei Yin,	1738
1683	athi Chandu, Kai-Wei Chang, Yejin Choi, and	Zaibin Zhang, Zelin Tan, Heng Zhou, Zhongzhi Li,	1739
1684	Bill Yuchen Lin. 2024a. Agent lumos: Unified and	Xiangyuan Xue, Yijiang Li, and 1 others. 2025c. The	1740
1685	modular training for open-source language agents .	landscape of agentic reinforcement learning for llms:	1741
1686	In <i>Proceedings of the 62nd Annual Meeting of the</i>	A survey . <i>arXiv preprint arXiv:2509.02547</i> .	1742
1687	<i>Association for Computational Linguistics (Volume 1:</i>		
1688	<i>Long Papers)</i> , pages 12380–12403, Bangkok, Thai-	Guibin Zhang, Junhao Wang, Junjie Chen, Wangchun-	1743
1689	land. Association for Computational Linguistics.	shu Zhou, Kun Wang, and Shuicheng Yan. 2025d.	1744
1690	Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun,	Agentracer: Who is inducing failure in the llm agentic	1745
1691	Tong Xu, and Enhong Chen. 2024b. A survey on	systems? <i>arXiv preprint arXiv:2509.03312</i> .	1746

A Other Aspects for Agent Training

1850

Training Target. While our focus is on the training of the main policy model, which is the component responsible for deciding the next action based on the current input, an agent system often contains **other trainable components** that also require dedicated training. These auxiliary components can be trained in similar ways to achieve specialized functionalities, such as managing memory (Zhang et al., 2025j; Yan et al., 2025; Shi et al., 2025b), performing explicit error correction (Wang et al., 2025g), or generating improved plans (Qu et al., 2025b; Xiong et al., 2025). For systems structured with multiple components, the training signals for a specific component can often be obtained by running the full system while keeping other components frozen. The complexity may be further amplified in **multi-agent systems** (Guo et al., 2024; Tran et al., 2025), which have become widely adopted, and whose effective training requires careful consideration (Albrecht et al., 2024; Park et al., 2025; Motwani et al., 2025; Hong et al., 2025).

1851

1852

1853

1854

1855

1856

1857

1858

1859

1860

Prompting. While this work has primarily focused on traditional model learning approaches that involve tuning model parameters, the remarkable advancements in LLMs necessitate discussing methods that leverage their inherent capabilities. Specifically, the strong in-context learning and instruction-following abilities of modern LLMs enable significant improvements in agent systems through prompting (Liu et al., 2023). The key advantage of **prompting-based** approaches is their flexibility and light-weight nature. They can be easily adopted to accumulate experience (Tang et al., 2025c), perform reflection (Shinn et al., 2023; Madaan et al., 2023) and incorporate external knowledge (Mialon et al., 2023).

1861

1862

1863

1864

1865

1866

1867