

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

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 = collect_query() ▷ Obtain Query (§4.1)
5:   T = sample(Q, A, M', E) ▷ Sample Trajectory (§4.2)
6:   S = get_signal(Q, T, E) ▷ Learning Signal (§5.1)
7:   M' = train(Q, T, S, M') ▷ Model Training (§5.2)
8: end for
9: return M'
```

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.

The primary highlight of this work is its *process-centric organization*. Instead of a purely taxonomic approach, this work is structured around the steps

of a typical agent training procedure, including key topics such as environment setup (§3), data collection (§4), and learning signal generation (§5). 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. With this practical organization, we hope to provide actionable guidance and clear design choices for better implementing and optimizing the entire agent training process.¹

2 Background

2.1 LLM-Based Agents

An LLM-based agent is an AI system designed to interact with the environment to achieve target goals. These agents leverage the inherent understanding and reasoning capabilities of LLMs and interact with the environment in an iterative way. At each time step t , the agent receives a representation of the environment’s current state $s_t \in \mathcal{S}$ and determines an action $a_t \in \mathcal{A}$, where \mathcal{S} and \mathcal{A} denote the state and action spaces, respectively. The action decision is made by the underlying LLM, which processes the environmental observations and the interaction history as input, and generates the action decision as output. The action will be executed, leading the environment to a new state s_{t+1} , where the agent will make further decisions. The interaction will continue until a termination state s_T , resulting in a sequence of states and actions, formally known as a trajectory: $\tau = (s_0, a_0, s_1, a_1, \dots, s_T)$.

2.2 Agent Training

Although LLMs have remarkable instruction-following capabilities, allowing them to be directly prompted for action decisions without specific parameter updates, the effectiveness of this approach is constrained, especially when the LLM is unfamiliar with the state and action representations of the target environment. To enable more stable and improved results, dedicated agent training is usually essential. While the overall agent training process shares similarities with the ordinary LLM training, the multi-step environmental interactions required by agent tasks introduce specific considerations.

¹The detailed paper collection procedure is described in Appendix A. In this work, we primarily focus on approaches that involve parameter updates to the main agent policy, while other related aspects are discussed in Appendix B.

3 Environment

The environment is a key component that distinguishes agent tasks from standard non-interactive text generation tasks. An agent operates in a multi-turn manner by iteratively receiving observations from the environment and deciding actions that further influence the environment’s state. Consequently, the representations of the observations and actions directly define the input and output spaces of the underlying decision model. A key motivation for agent training is precisely the LLM’s unfamiliarity with these specific environmental settings.

The nature of the environment is inherently tied to the target task, with each task type putting distinct 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.

3.1 Observation Representation

Observation Processing. The current state’s observation 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

164 augmenting inputs with element attributes and
165 textual choices. OmniParser (Lu et al., 2024b)
166 specifically fine-tunes models to parse screen-
167 shots and provide function semantic annotations.

- 168 • **Simplification** methods are essential because
169 raw observations are often complex and redund-
170 ant, a concern amplified by the context length
171 constraints of LLMs. Synapse (Zheng et al.,
172 2024b) extracts task-relevant information from
173 raw states to form cleaner observations. Simi-
174 larly, AgentOccam (Yang et al., 2025a) demon-
175 strates that providing less redundant yet infor-
176 mative observations can improve agent perfor-
177 mance. With similar objectives, LCoW (Lee
178 et al., 2025) trains a contextualization module
179 to enhance this process, FocusAgent (Kerboua
180 et al., 2025) uses a lightweight retriever to extract
181 the most relevant lines from the observation, and
182 Prune4Web (Zhang et al., 2025e) generates exe-
183 cutable Python scoring programs to dynamically
184 filter elements in observations.
185

186 **History Processing.** In addition to the current
187 observation, the record of past interaction steps is
188 a critical input source for the next action decision.
189 A primary challenge here is that as the number of
190 interaction steps increases, the cumulative history
191 representation can quickly grow, making it ineffi-
192 cient and difficult for the LLM to process the entire
193 context effectively. To manage the growing history:

- 194 • **Truncation** provides a straightforward solution,
195 by adopting a local window that includes only
196 the details of the most recent steps and omitting
197 previous steps (Tang et al., 2025b).
- 198 • **Summarization** provides a more informative ap-
199 proach by converting the previous trajectory into
200 a more concise representation. This approach can
201 be applied periodically (Yen et al., 2025), trig-
202 gered by heuristic rules such as approaching the
203 token budget (Wu et al., 2025e; Lu et al., 2025b),
204 or dynamically decided based by the agent itself
205 (Sun et al., 2025c; Ye et al., 2025).

206 The history management directly relates to the
207 broader concept of agent memory (Zhang et al.,
208 2025j) and context engineering (Mei et al., 2025),
209 as the previous trajectory can be regarded as the
210 agent’s short-term memory, necessitating efficient
211 encoding and retrieval to optimize the sequential
212 decision process.

3.2 Action Representation 213

214 The output generated by the LLM directly decides
215 the next action. For LLM-based agents, the output
216 space is *inherently diverse and structured*. This
217 is one of the key differences to classical agents,
218 which often operate with a few fixed action options.
219 The output decision is typically formalized in the
220 style of **function calling** (Patil et al., 2024), where
221 the LLM is required to select and parameterize a
222 function from either a pre-defined action list or a
223 tool-use library. To facilitate reliable parsing and
224 execution, these output representations often uti-
225 lize structured text formats such as JSON or other
226 special formats. Moreover, **code-based represen-**
227 **tation** is a natural and elegant way to represent
228 the output (Wang et al., 2024d), and it can enable
229 flexible action compositions by generating and ex-
230 ecuting programs (Nguyen et al., 2025b). Beyond
231 the format, the **design of the action space** is a criti-
232 cal area of investigation, for which there have been
233 some recent discussions on building more agent-
234 friendly APIs and tools (Song et al., 2025c; Lù
235 et al., 2025; Zhang et al., 2025a).

3.3 Environment Modeling 236

237 While the primary goal of agent training is to learn
238 a good policy $\pi(a_t|s_t)$ that decides the next action,
239 an agent system may also benefit from training
240 a model that predicts the transition of the envi-
241 ronment $p(s_{t+1}|s_t, a_t)$. This component, often re-
242 ferred to as a world model (Ha and Schmidhuber,
243 2018), predicts the next state s_{t+1} given the current
244 state s_t and the action taken a_t . Such **environ-**
245 **ment modeling** can be directly used in the agent’s
246 planning process, enabling decisions to be made
247 with a better internal simulation and awareness of
248 potential future outcomes (Chae et al., 2025a; Gu
249 et al., 2025). Furthermore, it can also enhance the
250 training process through better data synthesis (Fang
251 et al., 2025b) and environment synthesis (Guo et al.,
252 2025b; Liu et al., 2025a), or by serving as the tar-
253 get for specialized intermediate pre-training stages
254 (Copet et al., 2025).

Takeaway for §3: Environment design funda-
mentally shapes agent training by defining the
observation, history, and action representations,
often requiring specialized processing and, in
some cases, explicit environment modeling.

4 Data

Data quality is one of the most important factors for agent training. While **manual annotation** could offer the best way to ensure quality (Lu et al., 2024a; Rawles et al., 2023; Deng et al., 2023; Li et al., 2024), it is inherently tedious and cost-prohibitive. Moreover, agent tasks require not only the input queries, but also the multi-step interaction trajectories, posing great challenges to the data curation process. Consequently, recent research has focused on automated and scalable LLM-based **data synthesis** approaches (Tan et al., 2024). The following subsections² will illustrate how these techniques are specifically applied to agent tasks.

4.1 Queries

Similarly to standard LLM tasks, agent training first requires the synthesis of queries that define the goal that the agent needs to complete. We categorize related methods based on the complexity and the grounding level of the synthesizing process.

Direct Synthesis. This approach uses the inherent instruction-following capabilities of LLMs to synthesize new task queries in a direct way. The simplest method involves **direct prompting**, where an LLM is asked to generate target tasks, often inspired by existing examples. This is a technique widely adopted for agent task synthesis, drawing inspiration from Self-Instruct (Wang et al., 2023; Patel et al., 2024; He et al., 2025). Considering that agent tasks are intrinsically grounded in a specific environment, a natural and necessary extension is **context-aware querying**. Here, the synthesizing LLM is provided with environmental context, such as website names and screenshots (Zhou et al., 2025), crawled web pages (Wu et al., 2025b), pre-training documents (Cen et al., 2025), or multi-hop paths in knowledge graphs (Lu et al., 2025c). However, queries constructed via these direct methods may still suffer from limitations, potentially being too simple or lacking intricate connections to the target environment’s full complexity.

Iterative Extension. To address the simplicity limitation of previous methods and improve the overall difficulty of the synthesized queries, an **iterative extension** approach can be adopted. This technique usually starts with a seed item and incre-

²It 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.

mentally transforms it into a more complex goal through repeated refinement. The core goal is to generate *challenging but realistic* query-answer pairs, a process heavily investigated in recent work of deep search agents (Shi et al., 2025a; Gao et al., 2025b; Li et al., 2025b; Tao et al., 2025). The extension process itself is typically mediated by an LLM, often *equipped with external tools*. For example, WebExplorer (Liu et al., 2025b) operates by prompting LLMs equipped with search and browsing tools to guide the construction. Moreover, this approach can be used to build more complex agent tasks; AgentSynth (Xie et al., 2025), for example, iteratively forms a sequence of subtasks that are then summarized into a composite task for computer-use agents. Nevertheless, this approach may produce unnatural queries, and building complex yet natural tasks remains a challenge.

Exploratory Construction. To obtain queries that are more connected to the target environments, the construction itself can be viewed as an agent task, and a dedicated data construction agent can be designed to explore the environment. This typically follows a **reverse construction** process: firstly interacting with the environment with an exploration-based policy and then reversely synthesizing the query with the exploration trajectory. This approach has been adopted for a variety of target environments and tasks, including search (Wang et al., 2025c), web (Murty et al., 2024b; Pahuja et al., 2025; Trabucco et al., 2025), GUI (Su et al., 2025a; Sun et al., 2025b; Ramrakhya et al., 2025), and tool-use (Zhai et al., 2025). Although this method can create queries that are closely related to the target environment, there still remain challenges such as constructing non-trivial and challenging tasks and controlling the extra exploration cost.

4.2 Trajectories

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

Trajectory Conversion. A simple way to collect trajectories is through **trajectory conversion**, by repurposing and transforming existing trajectories originally collected for related tasks (Gandhi et al., 2024). For agent applications, this may be challenging due to the format differences across various target 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 such approaches effectively utilize 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.

Takeaway for §4: Agent training requires scalable synthesis of both task queries and interaction trajectories, with recent methods using LLM-driven generation, environment-grounded exploration, and sampling-based rollouts to construct diverse and high-quality training data.

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

446 methods like reinforcement learning with verifi- 496
447 able rewards (RLVR) approaches (Lambert et al., 497
448 2024; Guo et al., 2025a). If the target task is to 498
449 provide a single final answer, signals can be readily 499
450 obtained by matching the predicted answer with 500
451 the gold reference answer (Jin et al., 2025; Song 501
452 et al., 2025a; Chen et al., 2025; Qian et al., 2025; 502
453 Wei et al., 2025). In addition to final answers, there 503
454 can also be external signals that provide **partial** 504
455 **evaluation** of the predicted trajectory, such as enti- 505
456 tity matching ratios (Zhao et al., 2025) or sub-goal 506
457 scaffolding (Luo et al., 2025b) in complex search 507
458 tasks. An interesting connection is that external 508
459 signals can often be obtained **by construction** dur- 509
460 ing the data synthesis phase (§4), where the queries 510
461 and the corresponding target answers or evaluation 511
462 functions are constructed simultaneously. 512

463 **Model-Based Feedback.** The learning signals 513
464 can also be generated by models. It is common to 514
465 use stronger models in the way of LLM-as-a-judge 515
466 (Gu et al., 2024), such as providing judgments as 516
467 RL rewards (Lee et al., 2024) or filtering highest- 517
468 quality data (Zeng et al., 2024; Trabucco et al., 518
469 2025). Moreover, the agent or its underlying LLM 519
470 to be trained can also serve as an internal source, 520
471 generating its own learning signals. One of the 521
472 key metrics is **confidence**, which can potentially 522
473 indicate the correctness of generated trajectories. 523
474 This is typically calculated by checking the gener- 524
475 ation probability of the model outputs (Wu et al., 525
476 2025d) or by verifying consistency across multiple 526
477 generation paths (self-consistency; Huang et al., 527
478 2023; Kang et al., 2025). Leveraging the LLM’s 528
479 general-purpose ability to reason and evaluate, the 529
480 model can be treated as its own judge in a **self-** 530
481 **rewarding** setting (Yuan et al., 2024), generating 531
482 reward signals based on its own assessment of the 532
483 generated output. Furthermore, LLMs have demon- 533
484 strated strong abilities of **self-reflection** (Shinn 534
485 et al., 2023; Madaan et al., 2023), and this inherent 535
486 capacity can be utilized to provide internal critique 536
487 (Zhang et al., 2025f; Li et al., 2025c). 537

488 5.1.2 Signal Granularity 538

489 Agent tasks are unique in that they consist of multi- 539
490 turn trajectories containing interleaved action pre- 540
491 dictions and environmental observations, rather 541
492 than just a linear sequence of generated tokens 542
493 as in plain LLM tasks. Consequently, learning 543
494 signals can be provided at various levels, includ- 544
495 ing trajectory-level, step-level, or even token-level. 545
546

496 Considering the multi-step characteristic of agent 497
498 tasks, the most interesting discussions are on the 498
499 comparison between *trajectory-level versus step-* 499
500 *level signals*. This directly mirrors the general 500
501 LLM alignment debate over outcome reward mod- 501
502 els (ORM) versus process reward models (PRM). 502
503 While ORM signals are typically easier to obtain, 503
504 PRM signals offer denser supervision (Uesato et al., 504
505 2022; Lightman et al., 2024; Zheng et al., 2025a). 505

506 In agent tasks, ORM provides **trajectory-level** 506
507 **signals**, typically as one aggregated evaluation 507
508 score for the full agent trajectory. Many of the 508
509 verifiable rewards discussed in the previous sec- 509
510 tion (§5.1.1) are examples of this granularity. The 510
511 inherent challenge of this approach is that such sig- 511
512 nals are sparse, often provided only after the full 512
513 trajectory is completed, which may lack sufficient 513
514 supervision to effectively guide the learning for the 514
515 intermediate decision steps. 515

516 Considering **step-level signals** is a natural ex- 516
517 tension, as each step in the agent trajectory corre- 517
518 sponds directly to one action decision. Neverthe- 518
519 less, these finer-grained signals are much more dif- 519
520 ficult to obtain, for which manual annotation (Chae 520
521 et al., 2025b) could be tedious and expensive. We 521
522 categorize the typical automatic approaches used 522
523 to acquire step-level signals in the following: 523

- 524 • **Critic-based** methods use LLM-based reward 524
525 models, where an LLM critic can directly ex- 525
526 amine each step’s details and assign rewards 526
527 (Tan et al., 2025b; Yu et al., 2024; Zhang et al., 527
528 2025i). In addition, specific critic models can be 528
529 explicitly trained to provide these signals. For 529
530 instance, Liu et al. (2025d) optimize an implicit 530
531 PRM to provide step-level rewards as the log 531
532 ratios between the action probabilities with the 532
533 learned PRM and the old policy, while Xi et al. 533
534 (2025c) adopt a temporal difference-based es- 534
535 timation method for PRM training. Similarly, 535
536 Zhou et al. (2024c) learns a high-level value func- 536
537 tion with off-policy RL to provide signals for 537
538 low-level action generation. 538
- 539 • **Search-based** techniques use complex search 539
540 algorithms to yield more fine-grained rewards. 540
541 A common strategy involves generating a search 541
542 tree with the current policy and assigning rewards 542
543 to the intermediate steps according to the rewards 543
544 of its descendants (Hou et al., 2025). The search 544
545 tree can be constructed in various ways, includ- 545
546 ing Monte Carlo sampling (Xiong et al., 2024), 546
entropy-guided rollout (Shen et al., 2025), or it-

erative 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 should 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 utilized 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

647 assumes training with previously collected data
648 and no additional online interactions, intermediate
649 methods are often considered, adopting an iterative
650 and **coarse-grained** approach for balancing data
651 freshness with computational savings (Patel et al.,
652 2024; Aksitov et al., 2024; He et al., 2025). Further-
653 more, to fully utilize available computational re-
654 sources and parallelize operations, **asynchronous**
655 **training** methods have been increasingly employed
656 (Tan et al., 2025a; Gao et al., 2025b; Jiang et al.,
657 2025; Lu et al., 2025a).

658 **Training Stages.** Similar to plain LLM post-
659 training, the agent training procedure can be
660 split into **multiple stages** with different training
661 schemes. A widely adopted scheme involves first
662 performing an offline phase of SFT on expert
663 demonstrations, followed by an online phase of RL
664 to refine the policy through interaction (Vattikonda
665 et al., 2025; Wu et al., 2025b). Furthermore, to
666 enhance the model’s foundational understanding
667 and better support agentic behaviors, some work in-
668 corporates a **continual pre-training** stage that pre-
669 ceedes the main post-training (SFT/RL) fine-tuning
670 (Wang et al., 2025b; Su et al., 2025b; Copet et al.,
671 2025). This additional stage focuses on preparing
672 the model with relevant knowledge and capabilities
673 for effective downstream agent operation.

674 **Training Curriculum.** The data and environ-
675 mental setup can also be adjusted across the train-
676 ing procedure, naturally adopting the idea of **cur-
677 riculum learning** (Bengio et al., 2009). This ap-
678 proach structures the learning process by starting
679 with simpler tasks or data and gradually increas-
680 ing the difficulty. Examples of curriculum learn-
681 ing in agent training include starting with easy
682 samples and then progressively introducing more
683 complex ones (Lai et al., 2024), generating new
684 tasks from previous unsuccessful attempts to cre-
685 ate a self-evolving curriculum (Qi et al., 2025),
686 and incrementally degrading the quality of simu-
687 lated environment observations using a curriculum-
688 based rollout strategy to improve generalization
689 (Sun et al., 2025a).

Takeaway for §5: Agent training is driven by learning signals from experts, environments, or models, whose granularity and form determine the training losses, while training schemes organize data collection, optimization, and curriculum for efficient and effective learning.

6 Conclusions and Future Directions 691

692 This work presents a structured survey of training
693 methodologies for LLM-based agents, organizing
694 existing approaches along the key stages of the
695 training pipeline: environment design, data synthe-
696 sis, and training mechanisms. By examining these
697 components from a process-centric perspective, we
698 aim to clarify the design space of agent training
699 and provide guidance for future research.

700 **Scalable Supervision.** One main difficulty in train-
701 ing LLM-based agents is the lack of sufficient high-
702 quality supervision signals. Agent tasks often re-
703 quire long interaction trajectories in complex en-
704 vironments, making data collection costly while
705 reward signals remain sparse or noisy. One prom-
706 ising direction is reducing interaction cost through
707 simulation, where learned environment models al-
708 low agents to practice and generate trajectories in-
709 ternally (Ding et al., 2025; Liu et al., 2025a; Team
710 et al., 2025). Meanwhile, improving the automatic
711 construction of reliable reward models, especially
712 for tasks without verifiable outcomes, remains a
713 fundamental problem (Leike et al., 2018).

714 **Multi-modal Agent Training.** Many real-world
715 agent applications involve multi-modal environ-
716 ments such as graphical interfaces, visual observa-
717 tions, or audio inputs (Xie et al., 2024). While the
718 overall training procedure remains similar, these
719 settings introduce additional challenges in ground-
720 ing and perception, requiring agents to correctly
721 associate language and symbolic reasoning with
722 sensory inputs (Zheng et al., 2024a). Despite re-
723 cent advances in multi-modal LLMs (Yin et al.,
724 2024b), integrating perception, reasoning, and ac-
725 tion into a unified and reliably trained agent system
726 remains an open challenge.

727 **Life-Long Learning.** One important develop-
728 mental goal for agents is to enable fast and ro-
729 bust adaptation to the dynamic environments en-
730 countered in real-world applications (Zheng et al.,
731 2025b). To achieve this, it will be crucial to in-
732 tegrate agent training with established techniques
733 from life-long or continual learning (Chen and Liu,
734 2018; De Lange et al., 2021; Wang et al., 2024b;
735 Shi et al., 2024). This inherently involves a holi-
736 stic integration of the different components within
737 the agent system, making it crucial to accurately
738 identify which components should remain invariant
739 across tasks and which should be adapted quickly
740 in response to external environmental changes.

741 Limitations

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

765 References

766 Renat Aksitov, Sobhan Miryoosefi, Zonglin Li, Daliang
767 Li, Sheila Babayan, Kavya Kopparapu, Zachary
768 Fisher, Ruiqi Guo, Sushant Prakash, Pranesh Sriniv-
769 asan, Manzil Zaheer, Felix Yu, and Sanjiv Kumar.
770 2024. [ReST meets react: Self-improvement for multi-
771 step reasoning LLM agent](#). In *ICLR 2024 Workshop
772 on Large Language Model (LLM) Agents*.

773 Stefano V Albrecht, Filippos Christianos, and Lukas
774 Schäfer. 2024. *Multi-agent reinforcement learning:
775 Foundations and modern approaches*. MIT Press.

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

780 Greg Brockman, Vicki Cheung, Ludwig Pettersson,
781 Jonas Schneider, John Schulman, Jie Tang, and Woj-
782 ciech Zaremba. 2016. Openai gym. *arXiv preprint
783 arXiv:1606.01540*.

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

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*. 798 799 800 801 802 803 804

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*. 805 806 807 808 809 810 811 812 813 814

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*. 815 816 817 818 819

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*. 820 821 822 823

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*. 824 825 826 827 828 829

Zhiyuan Chen and Bing Liu. 2018. *Lifelong machine
learning*. Morgan & Claypool Publishers. 830 831

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*. 832 833 834 835 836 837

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. 838 839 840 841 842 843

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*. 844 845 846 847 848

849	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.	904
850		905
851		906
852		907
853		908
854	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.	909
855		910
856		911
857		912
858		913
859		914
860	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> .	915
861		916
862		917
863		918
864	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> .	919
865		920
866		921
867		922
868		923
869		924
870	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.	925
871		926
872		927
873		928
874		929
875		930
876		931
877		932
878	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> .	933
879		934
880		935
881		936
882	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.	937
883		938
884		939
885		940
886		941
887		942
888		943
889	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> .	944
890		945
891		946
892		947
893		948
894	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> .	949
895		950
896		951
897		952
898		953
899	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> .	954
900		955
901		956
902		957
903		958
		959
	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> .	904
		905
		906
		907
		908
		909
	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> .	910
		911
		912
		913
		914
		915
	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.	916
		917
		918
		919
		920
	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.	921
		922
		923
		924
		925
		926
		927
		928
		929
	David Ha and Jürgen Schmidhuber. 2018. World models. <i>arXiv preprint arXiv:1803.10122</i> .	930
		931
	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.	932
		933
		934
		935
		936
		937
		938
		939
		940
	Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. <i>arXiv preprint arXiv:1503.02531</i> .	941
		942
		943
	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> .	944
		945
		946
		947
		948
	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> .	949
		950
		951
		952
		953
		954
		955
	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>	956
		957
		958
		959

960					
961					
962					
963	Minda Hu, Tianqing Fang, Jianshu Zhang, Jun-Yu Ma,				
964	Zhisong Zhang, Jingyan Zhou, Hongming Zhang,				
965	Haitao Mi, Dong Yu, and Irwin King. 2025a. Web-				
966	CoT: Enhancing web agent reasoning by reconstruct-				
967	ing chain-of-thought in reflection, branching, and				
968	rollback . In <i>Findings of the Association for Com-</i>				
969	<i>putational Linguistics: EMNLP 2025</i> , pages 5155–				
970	5173, Suzhou, China. Association for Computational				
971	Linguistics.				
972	Xueyu Hu, Tao Xiong, Biao Yi, Zishu Wei, Ruix-				
973	uan Xiao, Yurun Chen, Jiasheng Ye, Meiling Tao,				
974	Xiangxin Zhou, Ziyu Zhao, and 1 others. 2025b.				
975	Os agents: A survey on mllm-based agents for				
976	general computing devices use. <i>arXiv preprint</i>				
977	<i>arXiv:2508.04482</i> .				
978	Jiaxin Huang, Shixiang Gu, Le Hou, Yuexin Wu, Xuezhi				
979	Wang, Hongkun Yu, and Jiawei Han. 2023. Large				
980	language models can self-improve . In <i>Proceedings</i>				
981	<i>of the 2023 Conference on Empirical Methods in Natu-</i>				
982	<i>ral Language Processing</i> , pages 1051–1068, Singa-				
983	apore. Association for Computational Linguistics.				
984	Xu Huang, Weiwen Liu, Xiaolong Chen, Xingmei				
985	Wang, Hao Wang, Defu Lian, Yasheng Wang, Ruim-				
986	ing Tang, and Enhong Chen. 2024. Understanding				
987	the planning of llm agents: A survey. <i>arXiv preprint</i>				
988	<i>arXiv:2402.02716</i> .				
989	Dongfu Jiang, Yi Lu, Zhuofeng Li, Zhiheng Lyu, Ping				
990	Nie, Haozhe Wang, Alex Su, Hui Chen, Kai Zou,				
991	Chao Du, and 1 others. 2025. Verltool: Towards				
992	holistic agentic reinforcement learning with tool use.				
993	<i>arXiv preprint arXiv:2509.01055</i> .				
994	Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Ser-				
995	can O Arik, Dong Wang, Hamed Zamani, and Jiawei				
996	Han. 2025. Search-r1: Training LLMs to reason and				
997	leverage search engines with reinforcement learning .				
998	In <i>Second Conference on Language Modeling</i> .				
999	Minki Kang, Jongwon Jeong, Seanie Lee, Jaewoong				
1000	Cho, and Sung Ju Hwang. 2025. Distilling llm agent				
1001	into small models with retrieval and code tools. <i>arXiv</i>				
1002	<i>preprint arXiv:2505.17612</i> .				
1003	Imene Kerboua, Sahar Omid Shayegan, Megh Thakkar,				
1004	Xing Han Lù, Léo Boisvert, Massimo Caccia, Jérémy				
1005	Espinas, Alexandre Aussem, Véronique Eglin, and				
1006	Alexandre Lacoste. 2025. Focusagent: Simple yet				
1007	effective ways of trimming the large context of web				
1008	agents. <i>arXiv preprint arXiv:2510.03204</i> .				
1009	Yoon Kim and Alexander M. Rush. 2016. Sequence-				
1010	level knowledge distillation . In <i>Proceedings of the</i>				
1011	<i>2016 Conference on Empirical Methods in Natu-</i>				
1012	<i>ral Language Processing</i> , pages 1317–1327, Austin,				
1013	Texas. Association for Computational Linguistics.				
1014	Jing Yu Koh, Stephen Marcus McAleer, Daniel Fried,				
1015	and Ruslan Salakhutdinov. 2025. Tree search for				
	language model agents . <i>Transactions on Machine</i>				1016
	<i>Learning Research</i> .				1017
	Hanyu Lai, Xiao Liu, Iat Long Iong, Shuntian Yao, Yux-				1018
	uan Chen, Pengbo Shen, Hao Yu, Hanchen Zhang,				1019
	Xiaohan Zhang, Yuxiao Dong, and 1 others. 2024.				1020
	Autowebglm: A large language model-based web				1021
	navigating agent. In <i>Proceedings of the 30th ACM</i>				1022
	<i>SIGKDD Conference on Knowledge Discovery and</i>				1023
	<i>Data Mining</i> , pages 5295–5306.				1024
	Nathan Lambert, Jacob Morrison, Valentina Pyatkin,				1025
	Shengyi Huang, Hamish Ivison, Faeze Brahma,				1026
	Lester James V Miranda, Alisa Liu, Nouha Dziri,				1027
	Shane Lyu, and 1 others. 2024. Tulu 3: Pushing fron-				1028
	tiers in open language model post-training. <i>arXiv</i>				1029
	<i>preprint arXiv:2411.15124</i> .				1030
	Dongjun Lee, Juyong Lee, Kyuyoung Kim, Jihoon Tack,				1031
	Jinwoo Shin, Yee Whye Teh, and Kimin Lee. 2025.				1032
	Learning to contextualize web pages for enhanced				1033
	decision making by LLM agents . In <i>The Thirteenth</i>				1034
	<i>International Conference on Learning Representa-</i>				1035
	<i>tions</i> .				1036
	Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas				1037
	Mesnard, Johan Ferret, Kellie Ren Lu, Colton Bishop,				1038
	Ethan Hall, Victor Carbune, Abhinav Rastogi, and				1039
	Sushant Prakash. 2024. RLAIF vs. RLHF: Scaling				1040
	reinforcement learning from human feedback with AI				1041
	feedback . In <i>Proceedings of the 41st International</i>				1042
	<i>Conference on Machine Learning</i> , volume 235 of				1043
	<i>Proceedings of Machine Learning Research</i> , pages				1044
	26874–26901. PMLR.				1045
	Jan Leike, David Krueger, Tom Everitt, Miljan Martić,				1046
	Vishal Maini, and Shane Legg. 2018. Scalable agent				1047
	alignment via reward modeling: a research direction.				1048
	<i>arXiv preprint arXiv:1811.07871</i> .				1049
	Sergey Levine, Aviral Kumar, George Tucker, and Justin				1050
	Fu. 2020. Offline reinforcement learning: Tutorial,				1051
	review, and perspectives on open problems. <i>arXiv</i>				1052
	<i>preprint arXiv:2005.01643</i> .				1053
	Chenliang Li, Adel Elmahdy, Alex Boyd, Zhongruo				1054
	Wang, Alfredo Garcia, Parminder Bhatia, Taha Kass-				1055
	Hout, Cao Xiao, and Mingyi Hong. 2025a. St-				1056
	ppo: Stabilized off-policy proximal policy optimiza-				1057
	tion for multi-turn agents training . <i>arXiv preprint</i>				1058
	<i>arXiv:2511.20718</i> .				1059
	Kuan Li, Zhongwang Zhang, Huifeng Yin, Liwen				1060
	Zhang, Litu Ou, Jialong Wu, Wenbiao Yin, Baixuan				1061
	Li, Zhengwei Tao, Xinyu Wang, and 1 others. 2025b.				1062
	Websailor: Navigating super-human reasoning for				1063
	web agent . <i>arXiv preprint arXiv:2507.02592</i> .				1064
	Shiyu Li, Yang Tang, Yifan Wang, Peiming Li, and				1065
	Xi Chen. 2025c. Reseek: A self-correcting frame-				1066
	work for search agents with instructive rewards .				1067
	<i>arXiv preprint arXiv:2510.00568</i> .				1068
	Wei Li, William E Bishop, Alice Li, Christopher Rawles,				1069
	Folawiyo Campbell-Ajala, Divya Tyamagundlu, and				1070
	Oriana Riva. 2024. On the effects of data scale on				1071

1072	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> .	1126
1073			1127
1074	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> .		1128
1075			1129
1076			1130
1077		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> .	1131
1078			1132
1079	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> .		1133
1080			1134
1081			1135
1082		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> .	1136
1083			1137
1084	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.		1138
1085		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.	1139
1086			1140
1087			1141
1088			1142
1089			1143
1090			1144
1091	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> .	1145
1092			1146
1093			1147
1094			
1095			
1096	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> .	1148
1097			1149
1098			1150
1099			1151
1100			1152
1101	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> .	1154
1102			1155
1103			1156
1104			1157
1105			
1106	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> .	1158
1107			1159
1108			1160
1109			1161
1110		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.	1162
1111			1163
1112	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.		1164
1113			1165
1114			1166
1115			1167
1116			
1117	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> .	1168
1118			1169
1119			1170
1120			1171
1121	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.	1173
1122			1174
1123			1175
1124			1176
1125			1177
			1178
			1179

1294	Ram Ramrakhya, Andrew Szot, Omar Attia, Yuhao Yang, Anh Nguyen, Bogdan Mazoure, Zhe Gan, Harsh Agrawal, and Alexander Toshev. 2025. Scaling synthetic task generation for agents via exploration. <i>arXiv preprint arXiv:2509.25047</i> .	1350
1295		1351
1296		1352
1297		1353
1298		1354
1299	Christopher Rawles, Alice Li, Daniel Rodriguez, Oriana Riva, and Timothy Lillicrap. 2023. Androidinthewild: A large-scale dataset for android device control. <i>Advances in Neural Information Processing Systems</i> , 36:59708–59728.	1355
1300		
1301		
1302		
1303		
1304	Jérémy Scheurer, Jon Ander Campos, Jun Shern Chan, Angelica Chen, Kyunghyun Cho, and Ethan Perez. 2022. Training language models with language feedback. <i>arXiv preprint arXiv:2204.14146</i> .	
1305		
1306		
1307		
1308	Leyang Shen, Yang Zhang, Chun Kai Ling, Xiaoyan Zhao, and Tat-Seng Chua. 2025. Carl: Critical action focused reinforcement learning for multi-step agent. <i>arXiv preprint arXiv:2512.04949</i> .	
1309		
1310		
1311		
1312	Dingfeng Shi, Jingyi Cao, Qianben Chen, Weichen Sun, Weizhen Li, Hongxuan Lu, Fangchen Dong, Tianrui Qin, King Zhu, Minghao Liu, and 1 others. 2025a. Taskcraft: Automated generation of agentic tasks. <i>arXiv preprint arXiv:2506.10055</i> .	
1313		
1314		
1315		
1316		
1317	Haizhou Shi, Zihao Xu, Hengyi Wang, Weiyi Qin, Wenyuan Wang, Yibin Wang, Zifeng Wang, Sayna Ebrahimi, and Hao Wang. 2024. Continual learning of large language models: A comprehensive survey. <i>ACM Computing Surveys</i> .	
1318		
1319		
1320		
1321		
1322	Yaorui Shi, Yuxin Chen, Siyuan Wang, Sihang Li, Hengxing Cai, Qi Gu, Xiang Wang, and An Zhang. 2025b. Look back to reason forward: Revisitable memory for long-context llm agents. <i>arXiv preprint arXiv:2509.23040</i> .	
1323		
1324		
1325		
1326		
1327	Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language agents with verbal reinforcement learning. <i>Advances in Neural Information Processing Systems</i> , 36:8634–8652.	
1328		
1329		
1330		
1331		
1332	Mohit Shridhar, Kingdi Yuan, Marc-Alexandre Cote, Yonatan Bisk, Adam Trischler, and Matthew Hausknecht. 2021. {ALFW}orld: Aligning text and embodied environments for interactive learning. In <i>International Conference on Learning Representations</i> .	
1333		
1334		
1335		
1336		
1337		
1338	Huatong Song, Jinhao Jiang, Yingqian Min, Jie Chen, Zhipeng Chen, Wayne Xin Zhao, Lei Fang, and Jirong Wen. 2025a. R1-searcher: Incentivizing the search capability in llms via reinforcement learning. <i>arXiv preprint arXiv:2503.05592</i> .	
1339		
1340		
1341		
1342		
1343	Yifan Song, Da Yin, Xiang Yue, Jie Huang, Sujian Li, and Bill Yuchen Lin. 2024. Trial and error: Exploration-based trajectory optimization of LLM agents. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 7584–7600, Bangkok, Thailand. Association for Computational Linguistics.	
1344		
1345		
1346		
1347		
1348		
1349		
	Yueqi Song, Ketan Ramaneti, Zaid Sheikh, Ziru Chen, Boyu Gou, Tianbao Xie, Yiheng Xu, Danyang Zhang, Apurva Gandhi, Fan Yang, and 1 others. 2025b. Agent data protocol: Unifying datasets for diverse, effective fine-tuning of llm agents. <i>arXiv preprint arXiv:2510.24702</i> .	1355
	Yueqi Song, Frank F. Xu, Shuyan Zhou, and Graham Neubig. 2025c. Beyond browsing: API-based web agents. In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 11066–11085, Vienna, Austria. Association for Computational Linguistics.	1356
		1357
		1358
		1359
		1360
		1361
	Hongjin Su, Ruoxi Sun, Jinsung Yoon, Pengcheng Yin, Tao Yu, and Sercan O Arik. 2025a. Learn-by-interact: A data-centric framework for self-adaptive agents in realistic environments. In <i>The Thirteenth International Conference on Learning Representations</i> .	1362
		1363
		1364
		1365
		1366
	Liangcai Su, Zhen Zhang, Guangyu Li, Zhuo Chen, Chenxi Wang, Maojia Song, Xinyu Wang, Kuan Li, Jialong Wu, Xuanzhong Chen, and 1 others. 2025b. Scaling agents via continual pre-training. <i>arXiv preprint arXiv:2509.13310</i> .	1367
		1368
		1369
		1370
		1371
	Hao Sun, Zile Qiao, Jiayan Guo, Xuanbo Fan, Yingyan Hou, Yong Jiang, Pengjun Xie, Yan Zhang, Fei Huang, and Jingren Zhou. 2025a. Zerosearch: Incentivize the search capability of llms without searching. <i>arXiv preprint arXiv:2505.04588</i> .	1372
		1373
		1374
		1375
		1376
	Qiushi Sun, Kanzhi Cheng, Zichen Ding, Chuanyang Jin, Yian Wang, Fangzhi Xu, Zhenyu Wu, Chengyou Jia, Liheng Chen, Zhoumianze Liu, Ben Kao, Guohao Li, Junxian He, Yu Qiao, and Zhiyong Wu. 2025b. OS-genesis: Automating GUI agent trajectory construction via reverse task synthesis. In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5555–5579, Vienna, Austria. Association for Computational Linguistics.	1377
		1378
		1379
		1380
		1381
		1382
		1383
		1384
		1385
		1386
	Weiwei Sun, Miao Lu, Zhan Ling, Kang Liu, Xuesong Yao, Yiming Yang, and Jiecao Chen. 2025c. Scaling long-horizon llm agent via context-folding. <i>arXiv preprint arXiv:2510.11967</i> .	1387
		1388
		1389
		1390
	Richard S. Sutton and Andrew G. Barto. 2018. <i>Reinforcement Learning: An Introduction</i> . A Bradford Book, Cambridge, MA, USA.	1391
		1392
		1393
	Sijun Tan, Michael Luo, Colin Cai, Tarun Venkat, Kyle Montgomery, Aaron Hao, Tianhao Wu, Arnab Balyan, Manan Roongta, Chenguang Wang, Li Erran Li, Raluca Ada Popa, and Ion Stoica. 2025a. rllm: A framework for post-training language agents. Notion Blog.	1394
		1395
		1396
		1397
		1398
		1399
	Weiting Tan, Xinghua Qu, Ming Tu, Meng Ge, Andy T Liu, Philipp Koehn, and Lu Lu. 2025b. Process-supervised reinforcement learning for interactive multimodal tool-use agents. <i>arXiv preprint arXiv:2509.14480</i> .	1400
		1401
		1402
		1403
		1404

1405	Zhen Tan, Dawei Li, Song Wang, Alimohammad	Le Sellier de Chezelles, Nicolas Gontier, Miguel	1461
1406	Beigi, Bohan Jiang, Amrita Bhattacharjee, Man-	Muñoz-Mármol, Sahar Omidi Shayegan, Stefania	1462
1407	sooreh Karami, Jundong Li, Lu Cheng, and Huan Liu.	Raimondo, and 1 others. 2025. How to train your llm	1463
1408	2024. Large language models for data annotation and	web agent: A statistical diagnosis. <i>arXiv preprint</i>	1464
1409	synthesis: A survey . In <i>Proceedings of the 2024 Con-</i>	<i>arXiv:2507.04103</i> .	1465
1410	<i>ference on Empirical Methods in Natural Language</i>		
1411	<i>Processing</i> , pages 930–957, Miami, Florida, USA.	Guoqing Wang, Sunhao Dai, Guangze Ye, Zeyu Gan,	1466
1412	Association for Computational Linguistics.	Wei Yao, Yong Deng, Xiaofeng Wu, and Zhenzhe	1467
1413	Fei Tang, Haolei Xu, Hang Zhang, Siqi Chen, Xingyu	Ying. 2025a. Information gain-based policy optimiza-	1468
1414	Wu, Yongliang Shen, Wenqi Zhang, Guiyang Hou,	tion: A simple and effective approach for multi-turn	1469
1415	Zeqi Tan, Yuchen Yan, and 1 others. 2025a. A sur-	llm agents. <i>arXiv preprint arXiv:2510.14967</i> .	1470
1416	vey on (m) llm-based gui agents. <i>arXiv preprint</i>		
1417	<i>arXiv:2504.13865</i> .	Haoming Wang, Haoyang Zou, Huatong Song, Jiazhan	1471
1418	Qiaoyu Tang, Hao Xiang, Le Yu, Bowen Yu, Yaojie Lu,	Feng, Junjie Fang, Junting Lu, Longxiang Liu, Qinyu	1472
1419	Xianpei Han, Le Sun, WenJuan Zhang, Pengbo Wang,	Luo, Shihao Liang, Shijue Huang, and 1 others.	1473
1420	Shixuan Liu, and 1 others. 2025b. Beyond turn limits:	2025b. Ui-tars-2 technical report: Advancing gui	1474
1421	Training deep search agents with dynamic context	agent with multi-turn reinforcement learning. <i>arXiv</i>	1475
1422	window. <i>arXiv preprint arXiv:2510.08276</i> .	<i>preprint arXiv:2509.02544</i> .	1476
1423	Xiangru Tang, Tianrui Qin, Tianhao Peng, Ziyang Zhou,	Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao	1477
1424	Daniel Shao, Tingting Du, Xinming Wei, Peng Xia,	Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang,	1478
1425	Fang Wu, He Zhu, and 1 others. 2025c. Agent kb:	Xu Chen, Yankai Lin, and 1 others. 2024a. A survey	1479
1426	Leveraging cross-domain experience for agentic prob-	on large language model based autonomous agents.	1480
1427	lem solving. <i>arXiv preprint arXiv:2507.06229</i> .	<i>Frontiers of Computer Science</i> , 18(6):186345.	1481
1428	Zhengwei Tao, Jialong Wu, Wenbiao Yin, Junkai	Liyuan Wang, Xingxing Zhang, Hang Su, and Jun Zhu.	1482
1429	Zhang, Baixuan Li, Haiyang Shen, Kuan Li, Li-	2024b. A comprehensive survey of continual learn-	1483
1430	wen Zhang, Xinyu Wang, Yong Jiang, and 1 others.	ing: Theory, method and application. <i>IEEE transac-</i>	1484
1431	2025. Webshaper: Agentic data synthesizing via	<i>tions on pattern analysis and machine intelligence</i> ,	1485
1432	information-seeking formalization. <i>arXiv preprint</i>	46(8):5362–5383.	1486
1433	<i>arXiv:2507.15061</i> .	Rui Wang, Ce Zhang, Jun-Yu Ma, Jianshu Zhang, Hon-	1487
1434	Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen,	gru Wang, Yi Chen, Boyang Xue, Tianqing Fang,	1488
1435	Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru	Zhisong Zhang, Hongming Zhang, and 1 others.	1489
1436	Chen, Yuankun Chen, Yutian Chen, and 1 others.	2025c. Explore to evolve: Scaling evolved aggrega-	1490
1437	2025. Kimi k2: Open agentic intelligence. <i>arXiv</i>	tion logic via proactive online exploration for deep	1491
1438	<i>preprint arXiv:2507.20534</i> .	research agents. <i>arXiv preprint arXiv:2510.14438</i> .	1492
1439	Guiyao Tie, Zeli Zhao, Dingjie Song, Fuyang Wei,	Ruiyi Wang and Prithviraj Ammanabrolu. 2025. A prac-	1493
1440	Rong Zhou, Yurou Dai, Wen Yin, Zhejiang Yang,	titioner’s guide to multi-turn agentic reinforcement	1494
1441	Jiangyue Yan, Yao Su, and 1 others. 2025. Large	learning. <i>arXiv preprint arXiv:2510.01132</i> .	1495
1442	language models post-training: Surveying techni-	Shuhe Wang, Shengyu Zhang, Jie Zhang, Runyi Hu,	1496
1443	ques from alignment to reasoning. <i>arXiv preprint</i>	Xiaoya Li, Tianwei Zhang, Jiwei Li, Fei Wu, Guoyin	1497
1444	<i>arXiv:2503.06072</i> .	Wang, and Eduard Hovy. 2024c. Reinforcement	1498
1445	Brandon Trabucco, Gunnar Sigurdsson, Robinson Pi-	learning enhanced llms: A survey. <i>arXiv preprint</i>	1499
1446	ramuthu, and Ruslan Salakhutdinov. 2025. Insta:	<i>arXiv:2412.10400</i> .	1500
1447	Towards internet-scale training for agents. <i>arXiv</i>	Xingyao Wang, Yangyi Chen, Lifan Yuan, Yizhe Zhang,	1501
1448	<i>preprint arXiv:2502.06776</i> .	Yunzhu Li, Hao Peng, and Heng Ji. 2024d. Exe-	1502
1449	Khanh-Tung Tran, Dung Dao, Minh-Duong Nguyen,	cutable code actions elicit better LLM agents . In	1503
1450	Quoc-Viet Pham, Barry O’Sullivan, and Hoang D	<i>Forty-first International Conference on Machine</i>	1504
1451	Nguyen. 2025. Multi-agent collaboration mech-	<i>Learning</i> .	1505
1452	anisms: A survey of llms. <i>arXiv preprint</i>	Yanlin Wang, Wanjun Zhong, Yanxian Huang, Ensheng	1506
1453	<i>arXiv:2501.06322</i> .	Shi, Min Yang, Jiachi Chen, Hui Li, Yuchi Ma, Qianx-	1507
1454	Jonathan Uesato, Nate Kushman, Ramana Kumar, Fran-	iang Wang, and Zibin Zheng. 2025d. Agents in soft-	1508
1455	cis Song, Noah Siegel, Lisa Wang, Antonia Creswell,	ware engineering: Survey, landscape, and vision. <i>Au-</i>	1509
1456	Geoffrey Irving, and Irina Higgins. 2022. Solv-	<i>tomated Software Engineering</i> , 32(2):1–36.	1510
1457	ing math word problems with process-and outcome-	Yiming Wang, Da Yin, Yuedong Cui, Ruichen Zheng,	1511
1458	based feedback. <i>arXiv preprint arXiv:2211.14275</i> .	Zhiqian Li, Zongyu Lin, Di Wu, Xueqing Wu,	1512
1459	Dheeraj Vattikonda, Santhoshi Ravichandran, Emiliano	Chenchen Ye, Yu Zhou, and 1 others. 2025e. Llms as	1513
1460	Penaloza, Hadi Nekoei, Megh Thakkar, Thibault	scalable, general-purpose simulators for evolving dig-	1514
		ital agent training. <i>arXiv preprint arXiv:2510.14969</i> .	1515

1516	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.	1573
1517		1574
1518		1575
1519		1576
1520		
1521		1577
1522		1578
1523		1579
1524	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> .	1580
1525		1581
1526		1582
1527		
1528		1583
1529		1584
1530	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> .	1585
1531		1586
1532		1587
1533		1588
1534		
1535	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> .	1589
1536		1590
1537		1591
1538		1592
1539		1593
1540		1594
1541		
1542	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> .	1595
1543		1596
1544		1597
1545		1598
1546	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> .	1599
1547		1600
1548		1601
1549		
1550		1602
1551		1603
1552		1604
1553	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> .	1605
1554		1606
1555		1607
1556		1608
1557		
1558	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.	1609
1559		1610
1560		1611
1561		1612
1562		1613
1563		1614
1564		1615
1565	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> .	1616
1566		1617
1567		1618
1568		1619
1569		1620
1570		1621
1571	Xixi Wu, Kuan Li, Yida Zhao, Liwen Zhang, Litu Ou, Huifeng Yin, Zhongwang Zhang, Xinmiao	1622
1572		
	Yu, Dingchu Zhang, Yong Jiang, and 1 others. 2025e. Resum: Unlocking long-horizon search intelligence via context summarization. <i>arXiv preprint arXiv:2509.13313</i> .	1623
		1624
		1625
		1626
		1627
		1628
	Yunjia Xi, Jianghao Lin, Yongzhao Xiao, Zheli Zhou, Rong Shan, Te Gao, Jiachen Zhu, Weiwen Liu, Yong Yu, and Weinan Zhang. 2025a. A survey of llm-based deep search agents: Paradigm, optimization, evaluation, and challenges. <i>arXiv preprint arXiv:2508.05668</i> .	1629
		1630
		1631
		1632
	Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, and 1 others. 2025b. The rise and potential of large language model based agents: A survey. <i>Science China Information Sciences</i> , 68(2):121101.	1633
		1634
		1635
		1636
	Zhiheng Xi, Chenyang Liao, Guanyu Li, Yajie Yang, Wenxiang Chen, Zhihao Zhang, Binghai Wang, Senjie Jin, Yuhao Zhou, Jian Guan, and 1 others. 2025c. Agentprm: Process reward models for llm agents via step-wise promise and progress. <i>arXiv preprint arXiv:2511.08325</i> .	1637
		1638
		1639
	Jingxu Xie, Dylan Xu, Xuandong Zhao, and Dawn Song. 2025. Agentsynth: Scalable task generation for generalist computer-use agents. <i>arXiv preprint arXiv:2506.14205</i> .	1640
		1641
		1642
	Junlin Xie, Zhihong Chen, Ruifei Zhang, Xiang Wan, and Guanbin Li. 2024. Large multimodal agents: A survey. <i>arXiv preprint arXiv:2402.15116</i> .	1643
		1644
		1645
	Weimin Xiong, Yifan Song, Qingxiu Dong, Bingchan Zhao, Feifan Song, XWang, and Sujian Li. 2025. MPO: Boosting LLM agents with meta plan optimization . In <i>Findings of the Association for Computational Linguistics: EMNLP 2025</i> , pages 3914–3935, Suzhou, China. Association for Computational Linguistics.	1646
		1647
		1648
		1649
		1650
	Weimin Xiong, Yifan Song, Xiutian Zhao, Wenhao Wu, Xun Wang, Ke Wang, Cheng Li, Wei Peng, and Sujian Li. 2024. Watch every step! LLM agent learning via iterative step-level process refinement . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 1556–1572, Miami, Florida, USA. Association for Computational Linguistics.	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

1629	Sikuan Yan, Xiufeng Yang, Zuchao Huang, Ercong Nie,	multimodal large language models. <i>National Science</i>	1685
1630	Zifeng Ding, Zonggen Li, Xiaowen Ma, Kristian Ker-	<i>Review</i> , 11(12):nwae403.	1686
1631	sting, Jeff Z Pan, Hinrich Schütze, and 1 others. 2025.		
1632	Memory-r1: Enhancing large language model agents	Yi Yu, Liuyi Yao, Yuexiang Xie, Qingquan Tan, Ji-	1687
1633	to manage and utilize memories via reinforcement	aqi Feng, Yaliang Li, and Libing Wu. 2026. Agen-	1688
1634	learning. <i>arXiv preprint arXiv:2508.19828</i> .	tic memory: Learning unified long-term and short-	1689
		term memory management for large language model	1690
1635	Jianwei Yang, Hao Zhang, Feng Li, Xueyan Zou, Chun-	agents. <i>arXiv preprint arXiv:2601.01885</i> .	1691
1636	yuan Li, and Jianfeng Gao. 2023. Set-of-mark		
1637	prompting unleashes extraordinary visual grounding	Yuanqing Yu, Zhefan Wang, Weizhi Ma, Shuai Wang,	1692
1638	in gpt-4v. <i>arXiv preprint arXiv:2310.11441</i> .	Chuhan Wu, Zhiqiang Guo, and Min Zhang. 2024.	1693
		Steptool: Enhancing multi-step tool usage in llms	1694
1639	Ke Yang, Yao Liu, Sapana Chaudhary, Rasool Fakoor,	through step-grained reinforcement learning. <i>arXiv</i>	1695
1640	Pratik Chaudhari, George Karypis, and Huzefa Rang-	<i>preprint arXiv:2410.07745</i> .	1696
1641	wala. 2025a. Agentoccam: A simple yet strong base-		
1642	line for LLM-based web agents . In <i>The Thirteenth</i>	Zhaochen Yu, Ling Yang, Jiaru Zou, Shuicheng Yan,	1697
1643	<i>International Conference on Learning Representa-</i>	and Mengdi Wang. 2025. Demystifying reinforce-	1698
1644	<i>tions</i> .	ment learning in agentic reasoning. <i>arXiv preprint</i>	1699
		<i>arXiv:2510.11701</i> .	1700
1645	Ruihan Yang, Fanghua Ye, Jian Li, Siyu Yuan, Yikai	Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho,	1701
1646	Zhang, Zhaopeng Tu, Xiaolong Li, and Deqing Yang.	Xian Li, Sainbayar Sukhbaatar, Jing Xu, and Ja-	1702
1647	2025b. The lighthouse of language: Enhancing LLM	son E Weston. 2024. Self-rewarding language mod-	1703
1648	agents via critique-guided improvement . In <i>The</i>	els . In <i>Forty-first International Conference on Ma-</i>	1704
1649	<i>Thirty-ninth Annual Conference on Neural Informa-</i>	<i>chine Learning</i> .	1705
1650	<i>tion Processing Systems</i> .		
1651	Zonghan Yang, Peng Li, Ming Yan, Ji Zhang, Fei Huang,	Aohan Zeng, Mingdao Liu, Rui Lu, Bowen Wang, Xiao	1706
1652	and Yang Liu. 2024. React meets actre: When lan-	Liu, Yuxiao Dong, and Jie Tang. 2024. AgentTun-	1707
1653	guage agents enjoy training data autonomy. <i>arXiv</i>	ing: Enabling generalized agent abilities for LLMs .	1708
1654	<i>preprint arXiv:2403.14589</i> .	In <i>Findings of the Association for Computational</i>	1709
		<i>Linguistics: ACL 2024</i> , pages 3053–3077, Bangkok,	1710
1655	Shunyu Yao, Howard Chen, John Yang, and Karthik	Thailand. Association for Computational Linguistics.	1711
1656	Narasimhan. 2022. Webshop: Towards scalable real-		
1657	world web interaction with grounded language agents.	Yunpeng Zhai, Shuchang Tao, Cheng Chen, Anni Zou,	1712
1658	<i>Advances in Neural Information Processing Systems</i> ,	Ziqian Chen, Qingxu Fu, Shinji Mai, Li Yu, Jiaji	1713
1659	35:20744–20757.	Deng, Zouying Cao, and 1 others. 2025. Agente-	1714
		volver: Towards efficient self-evolving agent system.	1715
1660	Rui Ye, Zhongwang Zhang, Kuan Li, Huifeng Yin,	<i>arXiv preprint arXiv:2511.10395</i> .	1716
1661	Zhengwei Tao, Yida Zhao, Liangcai Su, Liwen		
1662	Zhang, Zile Qiao, Xinyu Wang, and 1 others.	Chaoyun Zhang, Shilin He, Liqun Li, Si Qin, Yu Kang,	1717
1663	2025. Agentfold: Long-horizon web agents with	Qingwei Lin, Saravan Rajmohan, and Dongmei	1718
1664	proactive context management . <i>arXiv preprint</i>	Zhang. 2025a. Api agents vs. gui agents: Divergence	1719
1665	<i>arXiv:2510.24699</i> .	and convergence . <i>arXiv preprint arXiv:2503.11069</i> .	1720
1666	Asaf Yehudai, Lilach Eden, Alan Li, Guy Uziel, Yilun	Chaoyun Zhang, Shilin He, Jiaxu Qian, Bowen Li,	1721
1667	Zhao, Roy Bar-Haim, Arman Cohan, and Michal	Liqun Li, Si Qin, Yu Kang, Minghua Ma, Guyue	1722
1668	Shmueli-Scheuer. 2025. Survey on evaluation of llm-	Liu, Qingwei Lin, and 1 others. 2024a. Large lan-	1723
1669	based agents. <i>arXiv preprint arXiv:2503.16416</i> .	guage model-brained gui agents: A survey. <i>arXiv</i>	1724
		<i>preprint arXiv:2411.18279</i> .	1725
1670	Howard Yen, Ashwin Paranjape, Mengzhou Xia, The-	Chaoyun Zhang, Liqun Li, Shilin He, Xu Zhang,	1726
1671	jas Venkatesh, Jack Hessel, Danqi Chen, and Yuhao	Bo Qiao, Si Qin, Minghua Ma, Yu Kang, Qing-	1727
1672	Zhang. 2025. Lost in the maze: Overcoming con-	wei Lin, Saravan Rajmohan, Dongmei Zhang, and	1728
1673	text limitations in long-horizon agentic search. <i>arXiv</i>	Qi Zhang. 2025b. UFO: A UI-focused agent for win-	1729
1674	<i>preprint arXiv:2510.18939</i> .	dows OS interaction . In <i>Proceedings of the 2025</i>	1730
		<i>Conference of the Nations of the Americas Chap-</i>	1731
1675	Da Yin, Faeze Brahman, Abhilasha Ravichander, Khy-	<i>ter of the Association for Computational Linguistics:</i>	1732
1676	athi Chandu, Kai-Wei Chang, Yejin Choi, and	<i>Human Language Technologies (Volume 1: Long Pa-</i>	1733
1677	Bill Yuchen Lin. 2024a. Agent lumos: Unified and	<i>pers)</i> , pages 597–622, Albuquerque, New Mexico.	1734
1678	modular training for open-source language agents .	Association for Computational Linguistics.	1735
1679	In <i>Proceedings of the 62nd Annual Meeting of the</i>		
1680	<i>Association for Computational Linguistics (Volume 1:</i>	Guibin Zhang, Hejia Geng, Xiaohang Yu, Zhenfei Yin,	1736
1681	<i>Long Papers)</i> , pages 12380–12403, Bangkok, Thai-	Zaibin Zhang, Zelin Tan, Heng Zhou, Zhongzhi Li,	1737
1682	land. Association for Computational Linguistics.	Xiangyuan Xue, Yijiang Li, and 1 others. 2025c. The	1738
		landscape of agentic reinforcement learning for llms:	1739
1683	Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun,	A survey. <i>arXiv preprint arXiv:2509.02547</i> .	1740
1684	Tong Xu, and Enhong Chen. 2024b. A survey on		

1741	Guibin Zhang, Junhao Wang, Junjie Chen, Wangchunshu Zhou, Kun Wang, and Shuicheng Yan. 2025d. Agentracer: Who is inducing failure in the llm agentic systems? <i>arXiv preprint arXiv:2509.03312</i> .	1796
1742		1797
1743		1798
1744		1799
1745	Jianguo Zhang, Tian Lan, Rithesh Murthy, Zhiwei Liu, Weiran Yao, Ming Zhu, Juntao Tan, Thai Hoang, Zuxin Liu, Liangwei Yang, and 1 others. 2024b. Agentohana: Design unified data and training pipeline for effective agent learning. <i>arXiv preprint arXiv:2402.15506</i> .	1800
1746		1801
1747		1802
1748		1803
1749		1804
1750		1805
1751	Jiayuan Zhang, Kaiquan Chen, Zhihao Lu, Enshen Zhou, Qian Yu, and Jing Zhang. 2025e. Prune4web: Dom tree pruning programming for web agent. <i>arXiv preprint arXiv:2511.21398</i> .	1806
1752		1807
1753		1808
1754		1809
1755	Kai Zhang, Xiangchao Chen, Bo Liu, Tianci Xue, Zeyi Liao, Zhihan Liu, Xiyao Wang, Yuting Ning, Zhaorun Chen, Xiaohan Fu, and 1 others. 2025f. Agent learning via early experience. <i>arXiv preprint arXiv:2510.08558</i> .	1810
1756		1811
1757		1812
1758		1813
1759		1814
1760	Shaokun Zhang, Ming Yin, Jieyu Zhang, Jiale Liu, Zhiguang Han, Jingyang Zhang, Beibin Li, Chi Wang, Huazheng Wang, Yiran Chen, and Qingyun Wu. 2025g. Which agent causes task failures and when? on automated failure attribution of LLM multi-agent systems. In <i>Forty-second International Conference on Machine Learning</i> .	1815
1761		1816
1762		1817
1763		1818
1764		1819
1765		1820
1766		1821
1767	Wenlin Zhang, Xiaopeng Li, Yingyi Zhang, Pengyue Jia, Yichao Wang, Huifeng Guo, Yong Liu, and Xiangyu Zhao. 2025h. Deep research: A survey of autonomous research agents. <i>arXiv preprint arXiv:2508.12752</i> .	1822
1768		1823
1769		1824
1770		1825
1771		1826
1772	Yaocheng Zhang, Haohuan Huang, Zijun Song, Yuanheng Zhu, Qichao Zhang, Zijie Zhao, and Dongbin Zhao. 2025i. Criticsearch: Fine-grained credit assignment for search agents via a retrospective critic. <i>arXiv preprint arXiv:2511.12159</i> .	1827
1773		1828
1774		1829
1775		1830
1776		1831
1777	Zeyu Zhang, Quanyu Dai, Xiaohe Bo, Chen Ma, Rui Li, Xu Chen, Jieming Zhu, Zhenhua Dong, and Ji-Rong Wen. 2025j. A survey on the memory mechanism of large language model-based agents. <i>ACM Transactions on Information Systems</i> , 43(6):1–47.	1832
1778		1833
1779		1834
1780		1835
1781		1836
1782	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, and 1 others. 2023. A survey of large language models. <i>arXiv preprint arXiv:2303.18223</i> .	1837
1783		1838
1784		1839
1785		1840
1786		1841
1787	Yida Zhao, Kuan Li, Xixi Wu, Liwen Zhang, Dingchu Zhang, Baixuan Li, Maojia Song, Zhuo Chen, Chenxi Wang, Xinyu Wang, and 1 others. 2025. Repurposing synthetic data for fine-grained search agent supervision. <i>arXiv preprint arXiv:2510.24694</i> .	1842
1788		1843
1789		1844
1790		1845
1791		1846
1792	Boyuan Zheng, Boyu Gou, Jihyung Kil, Huan Sun, and Yu Su. 2024a. Gpt-4v (ision) is a generalist web agent, if grounded. In <i>International Conference on Machine Learning</i> , pages 61349–61385. PMLR.	1847
1793		1848
1794		1849
1795		1850
	Congming Zheng, Jiachen Zhu, Zhuoying Ou, Yuxiang Chen, Kangning Zhang, Rong Shan, Zeyu Zheng, Mengyue Yang, Jianghao Lin, Yong Yu, and 1 others. 2025a. A survey of process reward models: From outcome signals to process supervisions for large language models. <i>arXiv preprint arXiv:2510.08049</i> .	1851
		1852
	Junhao Zheng, Chengming Shi, Xidi Cai, Qiuke Li, Duzhen Zhang, Chenxing Li, Dong Yu, and Qianli Ma. 2025b. Lifelong learning of large language model based agents: A roadmap. <i>arXiv preprint arXiv:2501.07278</i> .	1853
		1854
		1855
		1856
	Longtao Zheng, Rundong Wang, Xinrun Wang, and Bo An. 2024b. Synapse: Trajectory-as-exemplar prompting with memory for computer control. In <i>The Twelfth International Conference on Learning Representations</i> .	1857
		1858
		1859
		1860
		1861
	Xuhui Zheng, Kang An, Ziliang Wang, Yuhang Wang, and Yichao Wu. 2025c. StepSearch: Igniting LLMs search ability via step-wise proximal policy optimization. In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 21816–21841, Suzhou, China. Association for Computational Linguistics.	1862
		1863
		1864
		1865
		1866
		1867
		1868
		1869
		1870
		1871
		1872
		1873
		1874
		1875
		1876
		1877
		1878
		1879
		1880
		1881
		1882
		1883
		1884
		1885
		1886
		1887
		1888
		1889
		1890
		1891
		1892
		1893
		1894
		1895
		1896
		1897
		1898
		1899
		1900

A Detailed Procedure of Paper Collecting

We begin the paper collection process by searching for relevant papers on arXiv. Figure 1 illustrates the quarterly volume of related papers by searching specific keywords on arXiv. The figure shows a rapid increase in publications on LLM-based agents. Together with this overall growth, research on agent training has also progressed rapidly, highlighting the need for a timely review of existing methodologies for training LLM-based agents.

To collect the statistics shown in Figure 1, we adopt a two-stage filtering process for arXiv papers. First, we perform a coarse-grained search using the keyword “agent” in the CS categories. We then apply a fine-grained filtering step to judge whether a paper focuses on LLM-based agents and whether it involves agent training. For this second step, we employ QWEN3-8B with the following prompt.

PROMPT FOR STATISTICS COLLECTING

Role: You are an expert research assistant with a specialization in Large Language Models and AI Agents. Your task is to carefully read a paper’s title and abstract and classify it according to a precise set of criteria.

Task: You will be given the title and abstract of a research paper. Based only on this information, you must:

- Determine if the paper is about LLM Agents.
- Determine if the paper involves Agent Training.
- Extract a list of the paper’s main keywords.
- Respond only with a single, valid JSON object containing your classifications.

Definitions:

1. “is_agent”: (bool) Set this to true if the paper describes an LLM-based system that exhibits autonomous behavior. Mark this as false if no LLM is utilized.
2. “is_training”: (bool) This field is dependent on the first. If “is_agent” is false, this must also be false. If “is_agent” is true, set this to true only if the paper discusses modifying the LLM’s model weights to improve its agentic capabilities.
 - Includes: Fine-tuning (full-parameter), Instruction tuning (on agent-specific data), Reinforcement Learning (e.g., DPO, PPO, GRPO), Parameter-Efficient Fine-Tuning (PEFT, e.g., LoRA, QLoRA), Adapter training, or training a new agent-specific model from scratch.

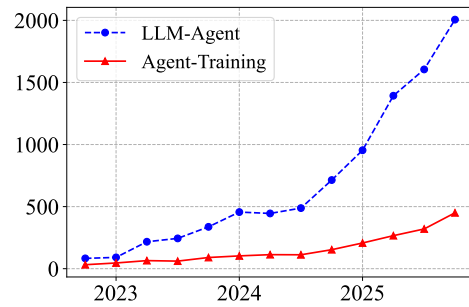


Figure 1: Trends of LLM-based agent research on arXiv. The data is aggregated per season (three months), illustrating the number of publications on general LLM-based agents and on dedicated agent training methods.

- Excludes (NOT training): Prompt Engineering (e.g., ReAct, Chain of Thought), In-Context Learning (Zero-shot, Few-shot), using a static pre-trained LLM as a “brain” without updating its weights, or building a system around a black-box LLM.

3. “keywords”: (list[str]) Extract a list of 3-5 of the most important and specific keywords or phrases from the abstract (e.g., “fine-tuning”, “tool use”, “autonomous agent”, “ReAct framework”).

Input:

Title: {title}

Abstract: {abstract}

Output Format (JSON Only): Your response must be only the JSON object, with no introductory text or explanation.

```
{
  "is_agent": bool,
  "is_training": bool,
  "keywords": ["keyword1", "keyword2", "keyword3"]
}
```

For the paper collection in this survey, we further adopt a recursive reference expansion (snowballing) approach by examining frequently cited papers encountered during reading. Our goal is to provide a timely and comprehensive survey of the related work, with particular attention to recent arxiv papers. We have conducted a thorough audit to replace arXiv citations with peer-reviewed conference or journal versions wherever available. Moreover, we also monitor newly available arxiv papers and include related ones.

B Other Aspects

Training Target. While our focus is on the training of the main policy model, which is the compo-

1882 nent responsible for deciding the next action based
1883 on the current input, an agent system often contains
1884 **other trainable components** that also require ded-
1885 icated training. These auxiliary components can
1886 be trained in similar ways to achieve specialized
1887 functionalities, such as managing memory (Zhang
1888 et al., 2025j; Yan et al., 2025; Shi et al., 2025b;
1889 Yu et al., 2026), performing explicit error correc-
1890 tion (Wang et al., 2025g), or generating improved
1891 plans (Qu et al., 2025b; Xiong et al., 2025). For
1892 systems structured with multiple components, the
1893 training signals for a specific component can of-
1894 ten be obtained by running the full system while
1895 keeping other components frozen. The complexity
1896 may be further amplified in **multi-agent systems**
1897 (Guo et al., 2024; Tran et al., 2025), which have
1898 become widely adopted, and whose effective train-
1899 ing requires careful consideration (Albrecht et al.,
1900 2024; Park et al., 2025; Motwani et al., 2025; Hong
1901 et al., 2025).

1902 **Prompting.** While this work has primarily fo-
1903 cused on traditional model learning approaches
1904 that involve tuning model parameters, the remark-
1905 able advancements in LLMs necessitate discussing
1906 methods that leverage their inherent capabilities.
1907 Specifically, the strong in-context learning and
1908 instruction-following abilities of modern LLMs
1909 enable significant improvements in agent systems
1910 through prompting (Liu et al., 2023). The key ad-
1911 vantage of **prompting-based** approaches is their
1912 flexibility and light-weight nature. They can be
1913 easily adopted to accumulate experience (Tang
1914 et al., 2025c), perform reflection (Shinn et al.,
1915 2023; Madaan et al., 2023) and incorporate external
1916 knowledge (Mialon et al., 2023).

1917 **C About AI Assistants**

1918 AI is only used for linguistic polishing; the concep-
1919 tual framework, literature synthesis, and organiza-
1920 tion are entirely the original work of the authors.