

Article

Not peer-reviewed version

Beyond Scaling: A Survey on Data-Efficient Agentic Learning

Yaqing Wang^{*}, [Zhenlin Luo](#), Peiyao Zhao, Yunfeng Cai, Quanming Yao

Posted Date: 8 May 2026

doi: 10.20944/preprints202605.0477.v1

Keywords: machine learning; few-shot learning; agent-based and multi-agent systems



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC, OpenAlex.

Copyright: This open access article is published under a [Creative Commons CC BY 4.0 license](#), which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Article

Beyond Scaling: A Survey on Data-Efficient Agentic Learning

Yaqing Wang^{1,*}, Zhenlin Luo^{1,2}, Peiyao Zhao^{1,3}, Yunfeng Cai¹ and Quanming Yao⁴

¹ Beijing Institute of Mathematical Sciences and Applications, Beijing, China

² Institute of Statistics and Big Data, Renmin University of China, Beijing, China

³ Yau Mathematical Sciences Center, Tsinghua University, Beijing, China

⁴ Department of Electronic Engineering, Tsinghua University, Beijing, China

* Correspondence: wangyaqing@bimsa.cn

Abstract

LLM-based agents are increasingly deployed across web and GUI automation, embodied decision making, and scientific workflows, yet their progress is often constrained by limited data and interaction. High-quality supervision is costly, and real-environment interactions are expensive, risky, and quickly invalidated by environment drift. This survey studies how to obtain and improve LLM-based agents with fewer samples, fewer labels, and fewer/cheaper interactions. We view agentic learning as a closed-loop decision process where experience arises from both external supervision and online interactions, and data efficiency requires maximizing information yield per unit cost. We then introduce a unified agentic learning framework and organize the literature along three complementary dimensions: experience augmentation, agent structural design, and learning paradigms. This perspective connects design choices to where learning signals come from, how they are utilized, and how adaptation is performed under bounded budgets. We summarize representative benchmarks and synthesize key open challenges, aiming to clarify the emerging landscape and support future progress in data-efficient agentic learning.

Keywords: machine learning; few-shot learning; agent-based and multi-agent systems

1. Introduction

Large language model (LLM)-based agents are rapidly moving beyond prompt-only prototypes into closed-loop systems that can perceive, reason/plan, and act in dynamic environments. This shift marks a transition from “model-as-a-service” to “model-as-an-agent”: success is no longer determined by one-shot generation quality, but by whether the agent can reliably acquire, verify, and refine behaviors while acting under dynamic feedback and long-horizon goals. Recent progress has enabled such agents to operate across a wide range of practical settings—from web and GUI automation to embodied decision making and scientific or medical workflows—where perception, reasoning, and action execution must be coordinated end-to-end [1–4].

In these realistic deployments, the dominant bottleneck is increasingly data efficiency rather than model scaling. Downstream agent tasks are often intrinsically data-scarce: high-quality supervision may be unavailable (e.g., it is unclear how to label every intermediate decision in an interactive trajectory), or feasible but prohibitively expensive (e.g., step-level grounding labels, demonstrations, expert feedback, or verification). Meanwhile, online trial-and-error is not “free data”: it consumes environment steps and tool calls, can require human verification, and may introduce safety or reliability risks, especially in high-stakes scenarios. Compounding the problem, interaction data can become stale under environment or interface drift—a common issue in software automation [1,2]—and similar constraints arise in embodied, scientific, and medical workflows where privacy, expert time, or experimental verification costs dominate [4–6]. As a result, practitioners often have to obtain and

improve agents with whatever limited signal is available, making “how to learn/obtain a capable LLM-based agent in a data-efficient way” a first-order question.

Data efficiency has been a central theme in machine learning, spanning few-shot learning [7] and sample-efficient reinforcement learning (RL) [8]. Yet agentic learning fundamentally broadens what “data” means and where efficiency comes from. Beyond labeled examples, agents learn from human annotations, trajectories, intermediate reasoning traces, tool-use patterns, and verification outcomes [9, 10]. Efficiency is therefore no longer solely a property of a learning algorithm; it emerges from the joint design of (i) experience and how it is generated, transformed, or simulated to reduce reliance on costly supervision, (ii) agent structure—including specialized perceivers and action executors—that reduces wasted interactions and localizes errors, and (iii) learning paradigms that maximize information gain from limited samples, labels, and interactions by governing whether and how model parameters are updated. Recent theoretical analyses further suggest that in-context adaptation can be understood as a form of implicit learning, helping explain strong few-shot generalization behaviors in modern LLMs [11].

Despite the rapidly growing literature, existing reviews of LLM-based agents [12–14] typically emphasize broad architectural scope or focus on individual components such as multi-agent architectures [15,16], feedback mechanisms [17], memory designs [18], or planning patterns [19]. They have not explicitly centered bounded supervision and interaction budgets as the organizing principle that connects techniques across experience acquisition, agent structure, and learning dynamics. At the same time, the growing scale and diversity of recent work make a unified, agent-centric synthesis along this dimension both timely and valuable.

In this survey, we provide a review of data-efficient agentic learning (Figure 1). Concretely, we make the following contributions: We introduce a unified agentic learning framework and a data-efficiency criterion grounded in limited samples, labels, and interactions. We organize the literature into a taxonomy along three complementary dimensions: (i) experience augmentation, (ii) agent structural design, and (iii) learning paradigms, connecting where learning signals originate, how they are utilized, and how adaptation is performed under bounded budgets. Finally, we summarize representative benchmarks across application domains and discuss open challenges that shape future progress.

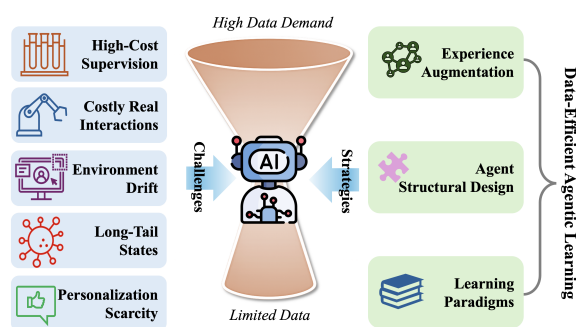


Figure 1. Conceptual overview of data-efficient agentic learning and its three complementary strategies.

2. Overview

We study how to obtain and improve LLM-based agents in a data-efficient way. An LLM-based agent can be viewed as a closed-loop decision-making system that repeatedly perceives the environment, reasons and plans with an LLM, executes actions (often via tools), and receives feedback

through interaction. At time step t , the environment is in state $s_t \in \mathcal{S}$ and the agent follows the interaction loop shown in Figure 2:

$$\begin{aligned} o_t &= P(s_t), \\ (g_t, a_t) &= L_\theta(o_t, m_{t-1}), \\ m_t &= M(m_{t-1}, g_t), \\ a'_t &= E(a_t), \\ s_{t+1} &= T(s_t, a'_t). \end{aligned}$$

Here P denotes a perceiver that maps the environment state to an observation, L_θ is the LLM that produces intermediate reasoning outputs g_t and selects the next action a_t , M is a memory module that maintains and updates the internal state m_t , and E is an action executor that converts the selected action into an executable form. The interaction yields online experience $\mathcal{D}_o = \{(s_t, o_t, a_t, \dots)\}_{t=1}^T$, while the agent may also leverage external experience \mathcal{D}_e collected outside its own interaction loop (e.g., demonstrations, labels, preference feedback, or verified outcomes). We denote the available experience by $\mathcal{D} = \mathcal{D}_o \cup \mathcal{D}_e$.

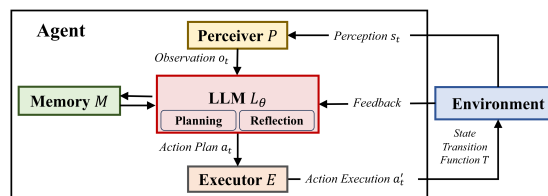


Figure 2. Agentic learning loop with core components.

Then, we define data-efficient agentic learning as follows.

Definition 1 (Data-Efficient Agentic Learning). *Data-efficient agentic learning studies how to obtain and improve LLM-based agents that operate in interactive decision-making settings under limited available experience \mathcal{D} .*

In this survey, we emphasize three coupled aspects. First, agentic interaction refers to a closed-loop process in which an agent repeatedly perceives the environment, reasons and plans with an LLM, executes actions, and incorporates feedback over time. Second, learning signals arise from both online interactions (trajectories, feedback, verification outcomes) and external supervision (demonstrations, labels, preference feedback, curated data). Third, an approach is data-efficient if its performance gains do not rely on collecting large amounts of new supervision or extensive real-environment trial-and-error, but instead improve the information yield per unit data and interaction.

Definition 1 highlights that data-efficiency bottlenecks stem from both costly supervision in \mathcal{D}_e and costly real-environment interaction in \mathcal{D}_o . Accordingly, this survey organizes existing methods along three complementary design levers that act on different parts of the agentic loop. Experience augmentation focuses on expanding the effective experience \mathcal{D} without proportional increases in real interaction (Section 3.1). Agent structural design reorganizes the internal modules and execution protocol (e.g., perceiver, memory, planning, reflection, and action executor) so that interactions become more directed, verifiable, and reusable, reducing redundant trial-and-error (Section 3.2). Learning paradigms characterize how agents are adapted from limited data and interaction (Section 3.3). Section 3 then elaborates this taxonomy and reviews representative methods in each category.

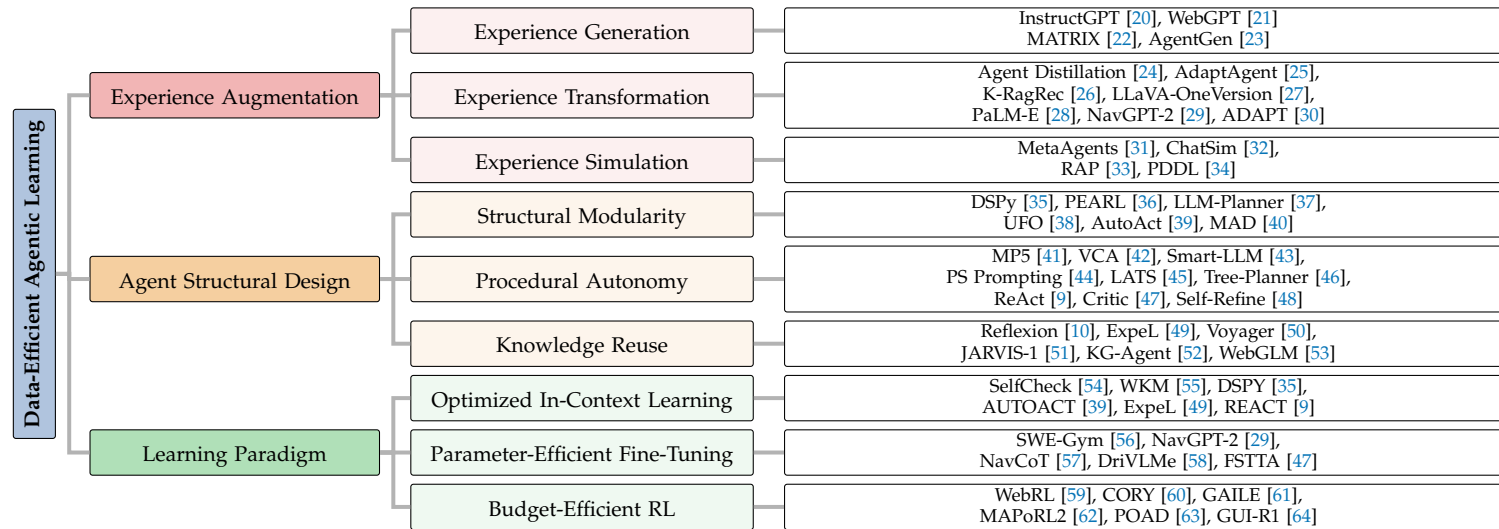


Figure 3. Taxonomy of data-efficient agentic learning.

3. Taxonomy

We now elaborate the taxonomy motivated by Definition 1 and Figure 2. The following three subsections review three complementary aspects of data-efficient agentic learning introduced in Section 2. For each aspect, we summarize its core idea, the type of data scarcity it addresses (samples, labels, or interactions), and representative works, noting that practical systems often combine multiple aspects.

3.1. Experience Augmentation

In data-efficient agentic learning, performance is often constrained by the availability and quality of external experience rather than model capacity. Real-world trajectories rarely cover long-tail states, human supervision is costly, and online interaction incurs substantial time, safety, and resource overhead. Under such constraints, naive trial-and-error yields low information gain and poor generalization.

Experience augmentation addresses this bottleneck by expanding and strengthening the effective experience pool under bounded budgets. Rather than collecting more data, its goal is to increase the density of task-relevant learning signal and reusable behavioral structure per unit of real experience. We organize existing approaches into three categories (Table 1): *experience generation*, *experience transformation*, and *experience simulation*.

Table 1. Experience augmentation strategies for data-efficient agentic learning. S/L/I denote sample/label/interaction.

Category	Core Idea	Data Type	Representative Works
Experience Generation	Create new high-quality data	S, L, I	InstructGPT [20]; WebGPT [21]; MATRIX [22]; AgentGen [23]
Experience Transformation	Increase information density	S, L	Agent Distillation [24]; AdaptAgent [25]; PaLM-E [28]; NavGPT-2 [29]
Experience Simulation	Shift interaction to cheaper surrogates	S, I	MetaAgents [31]; ChatSim [32]; RAP [33]; PDDL [34]

Experience Generation.

This category expands the effective experience pool beyond what is directly collected from the real environment, targeting coverage gaps and recurring failure modes under limited interaction budgets (Figure 4(a)). Instead of unconstrained data synthesis, experience generation focuses on producing high-quality trajectories that expose diverse error cases while maintaining reliable supervision. In LLM-based agents, experience generation often starts from *human-centric* designs, where a small amount of carefully curated demonstrations or preference feedback is used to concentrate supervision on effective behaviors and reduce exploration waste, as exemplified by InstructGPT [20], WebGPT [21]. It can further scale through *model-centric* generation, where agents synthesize additional trajectories at low marginal cost and rely on verification or structured feedback to control error propagation; representative examples include MATRIX [22], which generates interaction data via structured multi-agent simulation, and AgentGen [23], which expands coverage through environment-conditioned trajectory synthesis without additional real-world interaction. These methods improve data efficiency by creating additional training signals to expand coverage under limited supervision.

Experience Transformation.

This category improves data efficiency by enriching and restructuring limited real trajectories with complementary information, allowing each experience item to carry stronger and more reusable learning signals without additional interaction cost (Figure 4(b)). This is achieved by integrating real experience with external supervision or structure, followed by systematic reprocessing such as filtering, rewriting, relabeling, or compression. Existing approaches span multiple mechanisms: *knowledge distillation* transfers reasoning traces or interactive trajectories from stronger teachers (e.g., Agent Distillation [24]); *experience retrieval* reuses relevant demonstrations or structured knowledge as

in-context guidance (e.g., AdaptAgent [25], K-RagRec [26]); *cross-task transfer* enables data-scarce tasks to benefit from skills learned in data-rich domains (e.g., PaLM-E [28], LLaVA-OneVision [27]); and *modality enrichment* aligns multimodal signals to make supervision more explicit and informative (e.g., ADAPT [30], NavGPT-2 [29]). These methods transform existing experience to increase information density and reuse, thereby improving data efficiency.

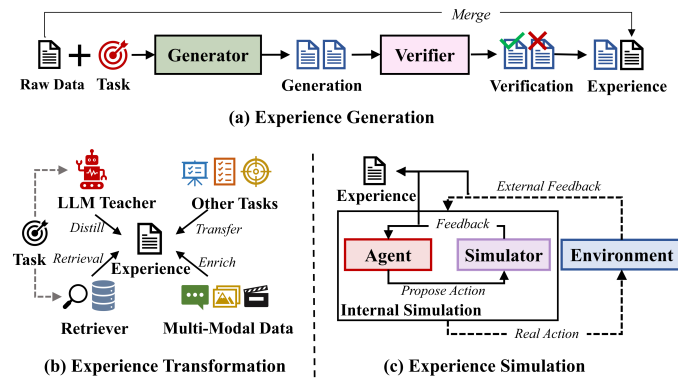


Figure 4. Illustration of experience augmentation strategy for data-efficient agentic learning. (a) Experience generation synthesizes additional training experience to expand coverage under limited interaction budgets. (b) Experience transformation enriches and restructures limited real experience into more reusable training signals. (c) Experience simulation replaces costly real-world interaction with simulated or modeled environments.

Experience Simulation.

This category reduces reliance on expensive real-world interaction by shifting exploration and failure discovery to cheaper surrogate environments (Figure 4(c)). Instead of real trial-and-error, agents explore within explicit simulators or learned world models to obtain diverse trajectories and feedback at lower cost. The objective is not perfect realism, but sufficient diversity and structural fidelity to complement scarce real experience. Systems such as MetaAgents [31] and ChatSim [32] demonstrate the utility of controllable simulated environments for generating targeted and rare interaction scenarios, while world-model-based approaches such as RAP [33] and symbolic planning frameworks like PDDL [34] enable agents to simulate outcomes and validate plans without repeated external execution. In sum, they improve data efficiency by substituting costly real-world interaction with lower-cost interaction sources.

3.2. Agent Structural Design

Agent structural design studies how to reorganize an LLM-based agent’s internal structure and execution protocol while holding data sources fixed, so as to increase the utility of each supervision signal or interaction. Rather than acquiring new experience, it focuses on how the agent perceives, reasons, plans, and executes actions through structured internal modules. The goal is to reduce unnecessary trial-and-error and environment steps by localizing supervision to the most informative stages. This often trades additional inference-time computation for fewer costly external interactions. From a data-efficiency perspective, structural design improves performance by shrinking the decision space, preventing costly error propagation, and enabling reuse of plans, skills, and memories. As a result, performance gains increasingly arise from structured internal self-improvement rather than additional external supervision or interaction. We organize existing designs into three categories (Table 2): *structural modularity*, *procedural autonomy*, and *knowledge reuse*, which respectively emphasize modular composition, controlled decision procedures, and reuse of prior knowledge to avoid redundant exploration.

Table 2. Agent structural design for data-efficient agentic learning.

Category	Core Idea	Core Modules	Representative Works
Structural Modularity	Decompose decision structure	Planner; Action Executor; Critic	DSPy [35]; PEARL [36]; LLM-Planner [37]; UFO [38]; AutoAct [39]; MAD [40]
Procedural Autonomy	Constrain execution procedure	Perceiver; Planner; Controller; Verifier	MP5 [41]; VCA [42]; Plan-and-Solve [44]; LATS [45]; ReAct [9]; Self-Refine [48]
Knowledge Reuse	Reuse prior experience	Memory; Skill library; External KB	Reflexion [10]; ExpeL [49]; Voyager [50]; JARVIS [51]; WebGLM [53]; KG-Agent [52]

Structural Modularity.

This line of work introduces explicit boundaries and interfaces into an agent’s internal workflow, transforming an entangled end-to-end reasoning–action process into coordinated components. From a data-efficiency perspective, modularity reduces global trial-and-error by localizing failures, enabling targeted supervision, and promoting reuse of intermediate artifacts. **Function decoupling**, which factorizes monolithic reasoning into planning, execution, and verification modules, allows errors to be corrected locally without restarting the entire decision loop, as exemplified by DSPy [35] and PEARL [36]; **hierarchical organization**, which separates high-level subgoal planning from low-level execution, compresses long-horizon decision-making via reusable action executors, as in LLM-Planner [37] and UFO [38]; and **role specialization** (Figure 5 (a)), where different agents or components take stable functional roles and exchange structured feedback, internalizes verification and coordination within the system rather than relying on external supervision, as demonstrated by AutoAct [39] and MAD [40]. These design strategies show that modular composition can substantially reduce external interaction cost and improve sample reuse under fixed data budgets.

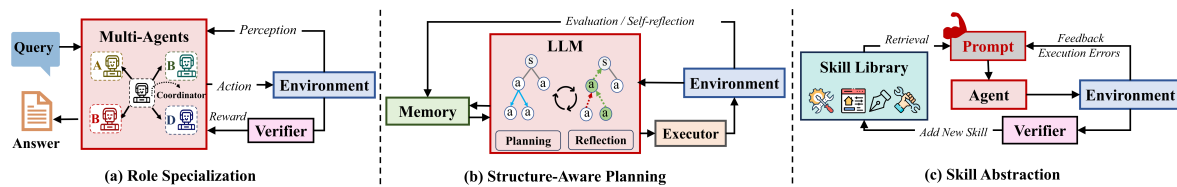


Figure 5. Illustrative examples of agent structural design for data-efficient agentic learning. (a) Role specialization, a representative instantiation of *structural modularity*, where the agent is decomposed into coordinated sub-roles with verifier feedback. (b) Structure-aware planning, a representative instantiation of *procedural autonomy*, where explicit state–action structure with memory, planning, and reflection guides decisions to reduce trial-and-error. (c) Skill abstraction, a representative instantiation of *knowledge reuse*, where the agent retrieves, composes, and verifies reusable skills from a library to avoid repeated low-level interactions.

Procedural Autonomy.

This line of work constrains agent behavior through explicit, reusable decision procedures, replacing unconstrained autoregressive generation with controlled iterative workflows. By deciding what to observe, how to decompose goals, when to act, and when to verify, procedural designs reduce wasted exploration and prevent cascading errors before costly external actions. **Active perception**, which treats perception as a decision policy over what and when to observe, selectively acquires task-relevant information under limited perception budgets, as in MP5 [41] and VCA [42]; **task decomposition and structure-aware planning** (Figure 5 (b)), which break long-horizon goals into verifiable substeps and restrict the search space via explicit plans or trees, reduce blind trial-and-error by enabling backtracking and reuse of partial solutions, as in Plan-and-Solve Prompting [44], Smart-LLM [43], LATS [45], and Tree-Planner [46]; and **execution control and self-verification**, which gate action execution through intermediate checks and critique, prevent error propagation before costly external actions, as in ReAct [9], Self-Refine [48], and Critic [65]; Collectively, these procedural

constraints shift performance gains toward structured internal self-improvement rather than additional external supervision or interaction.

Knowledge Reuse.

This line of work enables agents to avoid re-learning from scratch by converting available priors—past interactions, acquired skills, and external knowledge—into callable and transferable resources for future decisions. From a data-efficiency perspective, it shifts cost from repeated trial-and-error and human correction to reuse of compact representations that generalize across instances. **Memory compression**, which distills long and context-heavy interaction traces into concise, retrievable summaries or rules, helps agents avoid previously encountered failures without repeating costly exploration, as in Reflexion [10] and ExpeL [49]; **skill abstraction** (Figure 5 (c)), which transforms recurring behavior patterns into reusable subroutines or goal-conditioned controllers, enables compositional reuse of action structure across tasks, as in Voyager [50] and JARVIS-1 [51]; and **external knowledge bases**, where agents retrieve factual or structural priors on demand from the Web or knowledge graphs, reduce reliance on parametric memory and additional supervision, as in WebGLM [53] and KG-Agent [52]. Together, these reuse mechanisms substantially reduce redundant exploration and task-specific supervision by amortizing learning signal across time and tasks.

3.3. Learning Paradigm

When experience augmentation and structural design alone are insufficient, learning becomes unavoidable for improving agent performance. However, naive learning in agentic settings often incurs prohibitive supervision and interaction costs, violating the data-efficiency objective. We therefore regard a learning paradigm as data-efficient only when its gains do not rely on large-scale new supervision or extensive real-environment trial-and-error per task, but instead operate under strictly bounded data and interaction budgets.

Under this criterion, we organize existing approaches into three paradigms (Table 3): *optimized in-context learning (ICL)*, which enables inference-time adaptation without parameter updates; *parameter-efficient fine-tuning (PEFT)*, which achieves persistent adaptation via lightweight parameter updates; and *budget-efficient reinforcement learning (RL)*, which improves policies under constrained interaction and credit-assignment budgets. Together, they span a spectrum from transient to persistent adaptation, capturing key trade-offs among learning permanence, data efficiency, and interaction cost.

Table 3. Learning paradigms for data-efficient agentic learning.

Paradigm	Core Idea	Budget Type	Representative Works
Optimized ICL	Adapt without parameter updates	Inference-only	ReAct [9]; SelfCheck [54]; DSPy [35]; ExpeL [49]; WKM [55]; AutoAct [39]
PEFT	Partial parameter adaptation	Limited training data	SWE-Gym [56]; NavGPT-2 [29]; NavCoT [57]; DriVLM [58]; FSTTA [47]
Budget-Efficient RL	Policy learning via a few interaction	Budgeted interaction	WebRL [59]; CORY [60]; GAILE [61]; MAPoRL2 [62]; POAD [63]; GUI-R1 [64]

Optimized In-Context Learning (ICL).

This paradigm enables inference-time adaptation by reusing and restructuring contextual information, eliminating the need for parameter updates (Figure 6(a)). Early agentic prompting frameworks such as ReAct [9] show that interleaving reasoning, action, and observation within context allows agents to incorporate environment feedback without learning. Subsequent work improves data efficiency by optimizing context quality rather than quantity: SelfCheck [54] and DSPy [35] introduce verification and compilation mechanisms to refine prompts based on feedback signals, while knowledge-based approaches (e.g., ExpeL [49], WKM [55]) reuse distilled trial-and-error experience and inject task knowledge as in-context guidance to mitigate blind exploration and hallucinated actions. AutoAct [39] further demonstrates that agents can bootstrap high-quality contextual trajec-

stories through self-instruction and structured reflection under minimal human supervision. Overall, optimized in-context learning enables efficient inference-time adaptation of LLM-based agents without parameter updates; recent theoretical analyses further suggest that such behavior can be understood as implicit learning, exhibiting strong few-shot generalization [11].

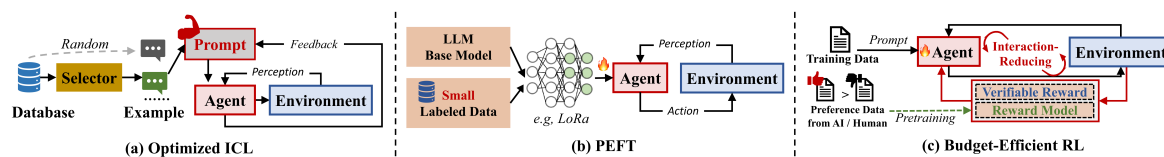


Figure 6. Illustration of learning paradigms for data-efficient agentic learning. (a) Optimized ICL adapts behavior at inference time without parameter updates. (b) PEFT enables persistent adaptation with limited demonstrations via lightweight updates. (c) Budget-efficient RL improves policies under constrained interaction budgets by enhancing learning signals rather than scaling rollouts.

Parameter-Efficient Fine-Tuning (PEFT).

This paradigm enables persistent task adaptation under limited supervision by updating only a small subset of parameters, reducing data demand, optimization cost, and the risk of catastrophic forgetting compared with full fine-tuning (Figure 6(b)). In agentic systems, PEFT allows specialization from a small number of demonstrations or interaction traces while preserving the generality of large pretrained models. Across embodied navigation, autonomous driving, and software agents, works such as NavGPT-2 [29], NavCoT [57], DriVLMe [58], and SWE-Gym [56] demonstrate that adopting lightweight adapters or projections can achieve strong performance with orders of magnitude fewer labeled trajectories. Recent extensions further show that PEFT can be applied online or at test time [47], enabling agents to adapt under strict data budgets without large-scale retraining.

Budget-Efficient Reinforcement Learning (RL).

This paradigm improves agent policies under strictly limited interaction budgets, expensive environment access, and sparse or delayed feedback, where conventional RL that scales rollouts becomes impractical (Figure 6(c)). Rather than increasing interactions, recent approaches enhance data efficiency by restructuring experience, shaping rewards, and improving credit assignment. For example, WebRL [59] and MAPoRL2 [62] demonstrate that self-evolving curriculum reuse, verifier-based rewards, and collaborative training can yield substantial gains with far fewer rollouts. Fine-grained credit assignment methods such as POAD [63] further accelerate learning by yielding finer-grained credit assignment without additional interactions. These strategies have proven particularly effective for GUI agents [64], highlighting that budget-efficient RL is less about scaling interaction and more about extracting maximal learning signal from minimal experience.

4. Applications and Benchmarks

We highlight five representative application domains—Web, GUI, Embodied AI, Medical, and Science—that capture the primary real-world settings in which data-efficient agentic learning is most critical. These domains span common agent interaction modalities, from text-based and vision-language interfaces to long-horizon decision making in physical and scientific workflows. Across all five domains, agents face similar structural constraints: effective interaction trajectories are costly to obtain, fine-grained grounding or expert supervision is expensive, environments and user interfaces evolve over time, and safety, privacy, or experimental considerations restrict large-scale trial-and-error. As a result, agent performance in these settings depends not only on task competence, but on how efficiently learning signals are acquired, reused, and transferred under limited supervision and interaction.

Table 4 lists a small set of widely used benchmarks selected to provide empirical grounding for these challenges. Rather than offering an exhaustive benchmark survey, we focus on benchmarks that (i) adopt relatively stable evaluation protocols, (ii) instantiate explicit perception–decision–action

loops, and (iii) expose dominant sources of data scarcity, including interaction cost, labeling or expert supervision cost, environment or interface drift, and regulatory or safety constraints. These benchmarks therefore serve as representative testbeds for studying how different data-efficient mechanisms—experience augmentation, agent structural design, and learning paradigms—translate into practical gains across application domains.

Table 4. Representative benchmarks for data-efficient agentic learning across application domains. L denotes language-only (text), V-L denotes vision-language, and 3D denotes embodied observations (3D state/environment).

Benchmark	Modality	Task	Link	Year
Web (<i>High interaction cost and rapid environment drift.</i>)				
WebArena [1]	L	Web Navigation	link	2023
Mind2Web [66]	V-L	Web Navigation	link	2023
WebVoyager [67]	V-L	Web Navigation	link	2024
GUI (<i>High annotation cost and UI drift across versions/devices.</i>)				
OSWorld [2]	V-L	Desktop Automation	link	2024
AndroidWorld [68]	V-L	Mobile App Automation	link	2024
ScreenSpot [69]	V-L	GUI Grounding	link	2025
Embodied AI (<i>Costly and safety-constrained real-world interaction.</i>)				
Franka-Kitchen [70]	3D	Manipulation	link	2019
ALFWorld [3]	L	Embodied Task Execution	link	2021
Meta-World [71]	3D	Multi-task Manipulation	link	2025
Medical (<i>Privacy-restricted data and costly expert supervision.</i>)				
ClinicalBench [5]	L	Clinical Prediction	link	2024
MedRAX [72]	V-L	Clinical Diagnosis	link	2025
MedAgentBench [73]	L	Clinical Decision Making	link	2025
Science (<i>Expert-labeled data and expensive experiments.</i>)				
GPQA [6]	L	Scientific Reasoning	link	2024
LAB-Bench [4]	V-L	Biology Research	link	2024
DiscoveryWorld [74]	V-L	Scientific Discovery	link	2024

5. Open Challenges

Data-efficient agentic learning departs from classical notions of sample efficiency because agent data is interactive, sequential, and heterogeneous across tasks, environments, and users. Each step can generate not only observations and rewards, but also reasoning traces, tool-use patterns, and verification signals that affect future decisions. When interaction, supervision, verification, and personalization all carry real cost, several challenges become central to making data-efficient agents practical and reliable.

Long-Horizon Learning.

Many agentic tasks require long-horizon decision making, where errors compound and the burden of verification, credit assignment, and exploration escalates quickly [57]. Although PEFT and budget-efficient RL reduce parameter-update cost, long-horizon settings often still incur substantial interaction and supervision overhead. A key direction is to make horizon a first-class factor in agent design and learning: agents should reason with checkpoints, reuse intermediate artifacts, and allocate verification strategically, rather than relying on naive rollout scaling.

Generalization and Drift.

General-purpose agents must transfer across tasks, tools, environments, and deployment conditions while relying on limited data and interactions. Since exhaustive coverage is infeasible, robust generalization hinges on learning reusable abstractions of reasoning, planning, and interaction rather than fitting individual trajectories [49]. This brings forward practical questions about what should be abstracted (e.g., decomposition patterns and tool-use strategies), how sparse interactions should trigger adaptation, and how to remain stable under distribution shift and environment/UI drift.

Personalization and User-Centric Learning.

Many real-world agents operate in personalized settings, where behavior must adapt to individual users, preferences, and constraints [75]. Personalization is intrinsically data-scarce: each user induces a distinct interaction distribution, and feedback is often implicit, noisy, or delayed. Core challenges include leveraging population-level structure to reduce per-user data needs, maintaining long-term user models under privacy constraints, and ensuring personalization does not erode generalization or safety.

System-Level Efficiency across Deployments.

Most agents are still improved in isolation, causing interaction and supervision costs to scale linearly with the number of deployment instances [60]. A promising direction is to treat data efficiency as a system objective: transform trajectories, skills, and verification outcomes into structured representations that can be selectively shared and reused across instances. Achieving this requires principled selection and aggregation, safeguards against error amplification, and evaluation protocols that measure marginal utility of shared experience rather than isolated task performance.

Self-Evolving Agents.

A long-term goal is open-ended agents that continually improve through interaction with environments, humans, and data [9], blurring the boundary between training-time and test-time learning. In deployment, however, unconstrained self-evolution is unrealistic because interaction, verification, and computation are budgeted. The challenge is to make self-evolution deliberate and economical: agents should decide when to explore, when to verify, and what to retain or reuse so improvement is sustainable and does not regress under drift.

Interaction-Centric Evaluation.

Most benchmarks emphasize final task success and implicitly treat interaction as free, which hides inefficiencies in learning and adaptation. More informative evaluation should report not only success rates, but also interaction steps, tool calls, verification frequency, supervision cost, and performance gain per unit interaction [10]. Such protocols would enable principled comparison of methods that trade computation, verification, and interaction differently.

High-Stakes, Data-Scarce Domains.

Interaction data is costly, sparse, or high-risk rather than simply “limited” in many applications. This is evident in embodied decision making [41], scientific discovery [76], and medical decision support, where unsafe or excessive trial-and-error is unacceptable. These domains call for data-efficient agents that integrate domain priors, rely on verifiable signals, and allocate scarce human supervision where it has the highest leverage.

6. Conclusion

This survey presented an agent-centric view of data-efficient agentic learning, focusing on how to obtain and improve LLM-based agents when supervision and real-world interactions are scarce, expensive, or risky. We framed the design space along three complementary dimensions: experience augmentation, agent structural design, and learning paradigms, which together aim to maximize information yield per unit cost, often trading additional inference-time computation and verification for fewer external interactions and less human supervision. We also summarized representative benchmarks across Web, GUI, embodied, medical, and scientific domains, and discussed open challenges in long-horizon learning under tight budgets, generalization beyond data coverage, user-centric personalization under sparse feedback, and interaction-centric evaluation. We hope this survey helps clarify the emerging landscape and supports the development of robust and deployable data-efficient agents.

References

1. Zhou, S.; Xu, F.F.; Zhu, H.; et al. WebArena: A Realistic Web Environment for Building Autonomous Agents. In Proceedings of the ICLR, 2024.
2. Xie, T.; Zhang, D.; Chen, J.; et al. OSWORLD: Benchmarking multimodal agents for open-ended tasks in real computer environments. In Proceedings of the NeurIPS, 2024.
3. Shridhar, M.; Yuan, X.; Côté, M.A.; et al. ALFWorld: Aligning Text and Embodied Environments for Interactive Learning. In Proceedings of the ICLR, 2021.
4. Laurent, J.M.; Janizek, J.D.; Ruzo, M.; et al. LAB-Bench: Measuring capabilities of language models for biology research. *arXiv preprint arXiv:2407.10362* 2024.
5. Chen, C.; Yu, J.; Chen, S.; et al. ClinicalBench: Can LLMs Beat Traditional ML Models in Clinical Prediction? In Proceedings of the AAAI Workshop, 2024.
6. Rein, D.; Hou, B.L.; Stickland, A.C.; et al. GPQA: A Graduate-Level Google-Proof Q&A Benchmark. In Proceedings of the COLM, 2024.
7. Wang, Y.; Yao, Q.; Kwok, J.T.; et al. Generalizing from a few examples: A survey on few-shot learning. *CSUR* 2020, 53, 1–34.
8. Yu, Y. Towards Sample Efficient Reinforcement Learning. In Proceedings of the IJCAI, 2018, pp. 5739–5743.
9. Yao, S.; Zhao, J.; Yu, D.; et al. ReAct: Synergizing reasoning and acting in language models. In Proceedings of the ICLR, 2022.
10. Shinn, N.; Cassano, F.; Gopinath, A.; et al. Reflexion: Language Agents with Verbal Reinforcement Learning. In Proceedings of the NeurIPS, 2023.
11. Wu, S.; Wang, Y.; Yao, Q. Why In-Context Learning Models are Good Few-Shot Learners? In Proceedings of the ICLR, 2025.
12. Wang, L.; Ma, C.; Feng, X.; et al. A survey on large language model based autonomous agents. *Front. Comput. Sci.* 2024, 18, 186345.
13. Sang, J.; Xiao, J.; Han, J.; et al. Beyond Pipelines: A Survey of the Paradigm Shift toward Model-Native Agentic AI. *J. ACM* 2025.
14. Liu, B.; Li, X.; Zhang, J.; et al. Advances and challenges in foundation agents: From brain-inspired intelligence to evolutionary, collaborative, and safe systems. *arXiv preprint arXiv:2504.01990* 2025.
15. Wu, D.; Wei, X.; Chen, G.; et al. Generative multi-agent collaboration in embodied AI: A systematic review. In Proceedings of the IJCAI, 2025.
16. Guo, T.; Chen, X.; Wang, Y.; et al. Large language model based multi-agents: A survey of progress and challenges. In Proceedings of the IJCAI, 2024.
17. Liu, Z.; Bai, X.; Chen, K.; et al. A survey on the feedback mechanism of LLM-based AI agents. In Proceedings of the IJCAI, 2025.
18. Zhang, Z.; Dai, Q.; Bo, X.; et al. A survey on the memory mechanism of large language model-based agents. *TOIS* 2025, 43, 1–47.
19. Torreno, A.; Onaindia, E.; et al. Cooperative multi-agent planning: A survey. *CSUR* 2017.
20. Ouyang, L.; Wu, J.; Jiang, X.; et al. Training language models to follow instructions with human feedback. In Proceedings of the NeurIPS, 2022.
21. Nakano, R.; Hilton, J.; Balaji, S.; et al. WebGPT: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332* 2021.
22. Tang, S.; Pang, X.; Liu, Z.; et al. Synthesizing post-training data for LLMs through multi-agent simulation. In Proceedings of the ACL, 2025.
23. Hu, M.; Zhao, P.; Xu, C.; et al. AgentGen: Enhancing Planning Abilities for Large Language Model based Agent via Environment and Task Generation. In Proceedings of the KDD, 2025.
24. Kang, M.; Jeong, J.; Lee, S.; et al. Distilling LLM agent into small models with retrieval and code tools. In Proceedings of the NeurIPS, 2025.
25. Verma, G.; Kaur, R.; Srishankar, N.; et al. AdaptAgent: Adapting Multimodal Web Agents with Few-Shot Learning from Human Demonstrations. In Proceedings of the ACL, 2025.
26. Wang, S.; Fan, W.; Feng, Y.; et al. Knowledge graph retrieval-augmented generation for LLM-based recommendation. In Proceedings of the ACL, 2025.
27. Li, B.; Zhang, Y.; Guo, D.; et al. LLaVA-OneVision: Easy Visual Task Transfer. *TMLR* 2025.
28. Driess, D.; Xia, F.; Sajjadi, M.S.; et al. PaLM-E: An Embodied Multimodal Language Model. In Proceedings of the ICML, 2023.

29. Zhou, G.; Hong, Y.; Wang, Z.; et al. NavGPT-2: Unleashing navigational reasoning capability for large vision-language models. In Proceedings of the ECCV, 2024.
30. Lin, B.; Zhu, Y.; Chen, Z.; et al. ADAPT: Vision-Language Navigation With Modality-Aligned Action Prompts. In Proceedings of the CVPR, 2022.
31. Li, Y.; Sun, L.; Zhang, Y. MetaAgents: Large Language Model Based Agents for Decision-Making on Teaming. *HCI* 2025.
32. Wei, Y.; Wang, Z.; Lu, Y.; et al. Editable Scene Simulation for Autonomous Driving via Collaborative LLM-Agents. In Proceedings of the CVPR, 2024.
33. Hao, S.; Gu, Y.; Ma, H.; et al. Reasoning with language model is planning with world model. In Proceedings of the EMNLP, 2023.
34. Guan, L.; Valmeekam, K.; Sreedharan, S.; et al. Leveraging Pre-trained Large Language Models to Construct and Utilize World Models for Model-based Task Planning. In Proceedings of the NeurIPS, 2023.
35. Khattab, O.; Singhvi, A.; Maheshwari, P.; et al. DSPy: Compiling Declarative Language Model Calls into Self-Improving Pipelines. In Proceedings of the NeurIPS, 2023.
36. Sun, S.; Liu, Y.; Wang, S.; et al. PEARL: Prompting large language models to plan and execute actions over long documents. In Proceedings of the EACL, 2024.
37. Song, C.H.; Wu, J.; Washington, C.; et al. LLM-Planner: Few-shot grounded planning for embodied agents with large language models. In Proceedings of the ICCV, 2023.
38. Zhang, C.; Li, L.; He, S.; et al. UFO: A UI-focused agent for Windows OS interaction. In Proceedings of the NAACL, 2025.
39. Qiao, S.; Zhang, N.; Fang, R.; et al. AutoAct: Automatic agent learning from scratch for QA via self-planning. In Proceedings of the ACL, 2024.
40. Liang, T.; He, Z.; Jiao, W.; et al. Encouraging divergent thinking in large language models through multi-agent debate. In Proceedings of the EMNLP, 2024.
41. Qin, Y.; Zhou, E.; Liu, Q.; et al. Mp5: A multi-modal open-ended embodied system in minecraft via active perception. In Proceedings of the CVPR, 2024.
42. Yang, Z.; Chen, D.; Yu, X.; et al. VCA: Video curious agent for long video understanding. In Proceedings of the ICCV, 2025.
43. Kannan, S.S.; Venkatesh, V.L.; Min, B.C. Smart-LLM: Smart multi-agent robot task planning using large language models. In Proceedings of the IROS, 2024.
44. Wang, L.; Xu, W.; Lan, Y.; et al. Plan-and-Solve Prompting: Improving Zero-Shot Chain-of-Thought Reasoning by Large Language Models. In Proceedings of the ACL, 2023.
45. Zhou, A.; Yan, K.; Shlapentokh-Rothman, M.; et al. Language agent tree search unifies reasoning, acting, and planning in language models. In Proceedings of the ICML, 2024.
46. Hu, M.; Mu, Y.; Yu, X.; et al. Tree-Planner: Efficient Close-loop Task Planning with Large Language Models. In Proceedings of the ICLR, 2024.
47. Gao, J.; Yao, X.; Xu, C. Fast-slow test-time adaptation for online vision-and-language navigation. In Proceedings of the ICML, 2024.
48. Madaan, A.; Tandon, N.; Gupta, P.; et al. Self-Refine: Iterative Refinement with Self-Feedback. In Proceedings of the NeurIPS, 2023.
49. Zhao, A.; Huang, D.; Xu, Q.; et al. ExpeL: LLM agents are experiential learners. In Proceedings of the AACL, 2024.
50. Wang, G.; Xie, Y.; Jiang, Y.; et al. Voyager: An Open-Ended Embodied Agent with Large Language Models. *TMLR* 2024.
51. Wang, Z.; Cai, S.; Liu, A.; et al. JARVIS-1: Open-world multi-task agents with memory-augmented multimodal language models. *TPAMI* 2024.
52. Jiang, J.; Zhou, K.; Zhao, W.X.; et al. KG-Agent: An efficient autonomous agent framework for complex reasoning over knowledge graph. In Proceedings of the ACL, 2025.
53. Liu, X.; Lai, H.; Yu, H.; et al. WebGLM: Towards an efficient web-enhanced question answering system with human preferences. In Proceedings of the KDD, 2023.
54. Miao, N.; Teh, Y.W.; Rainforth, T. SelfCheck: Using LLMs to Zero-Shot Check Their Own Step-by-Step Reasoning. In Proceedings of the ICLR, 2024.
55. Qiao, S.; Fang, R.; Zhang, N.; et al. Agent Planning with World Knowledge Model. In Proceedings of the NeurIPS, 2024.

56. Pan, J.; Wang, X.; Neubig, G.; et al. Training Software Engineering Agents and Verifiers with SWE-Gym. In Proceedings of the ICML, 2024.
57. Lin, B.; Nie, Y.; Wei, Z.; et al. NavCoT: Boosting LLM-based vision-and-language navigation via learning disentangled reasoning. *TPAMI* **2025**.
58. Huang, Y.; Sansom, J.; Ma, Z.; et al. DriVLMe: Enhancing LLM-based autonomous driving agents with embodied and social experiences. In Proceedings of the IROS, 2024.
59. Qi, Z.; Liu, X.; Iong, I.L.; et al. WebRL: Training LLM Web Agents via Self-Evolving Online Curriculum Reinforcement Learning. In Proceedings of the ICLR, 2025.
60. Ma, H.; Hu, T.; Pu, Z.; et al. Coevolving with the Other You: Fine-Tuning LLM with Sequential Cooperative Multi-Agent Reinforcement Learning. In Proceedings of the NeurIPS, 2024.
61. Feng, P.; He, Y.; Huang, G.; et al. AGILE: A Novel Reinforcement Learning Framework of LLM Agents. In Proceedings of the NeurIPS, 2024.
62. Park, C.; Han, S.; Guo, X.; et al. MAPoRL: Multi-agent post-co-training for collaborative large language models with reinforcement learning. In Proceedings of the ACL, 2025.
63. Wen, M.; Wan, Z.; Wang, J.; et al. Reinforcing LLM Agents via Policy Optimization with Action Decomposition. In Proceedings of the NeurIPS, 2024.
64. Luo, R.; Wang, L.; He, W.; et al. GUI-R1: A generalist R1-style vision-language action model for GUI agents. *arXiv preprint arXiv:2504.10458* **2025**.
65. Gou, Z.; Shao, Z.; Gong, Y.; et al. CRITIC: Large Language Models Can Self-Correct with Tool-Interactive Critiquing. In Proceedings of the ICLR, 2024.
66. Deng, X.; Gu, Y.; Zheng, B.; et al. Mind2Web: Towards a Generalist Agent for the Web. In Proceedings of the NeurIPS, 2023.
67. He, H.; Yao, W.; Ma, K.; et al. WebVoyager: Building an End-to-End Web Agent with Large Multimodal Models. In Proceedings of the ACL, 2024.
68. Rawles, C.; Clinckemaillie, S.; Chang, Y.; et al. AndroidWorld: A Dynamic Benchmarking Environment for Autonomous Agents. In Proceedings of the ICLR, 2025.
69. Li, K.; Meng, Z.; Lin, H.; Luo, Z.; et al. ScreenSpot-Pro: GUI Grounding for Professional High-Resolution Computer Use. In Proceedings of the MM, 2025.
70. Gupta, A.; Kumar, V.; Lynch, C.; et al. Relay Policy Learning: Solving Long-Horizon Tasks via Imitation and Reinforcement Learning. In Proceedings of the CoRL, 2020.
71. McLean, R.; Chatzaroulas, E.; McCutcheon, L.; et al. Meta-World+: An Improved, Standardized, RL Benchmark. In Proceedings of the NeurIPS, 2025.
72. Fallahpour, A.; Ma, J.; Munim, A.; et al. MedRAX: Medical Reasoning Agent for Chest X-ray. In Proceedings of the ICML, 2025.
73. Jiang, Y.; Black, K.C.; Geng, G.; et al. MedAgentBench: A Virtual EHR Environment to Benchmark Medical LLM Agents. *NEJMAI* **2025**.
74. Jansen, P.; Côté, M.A.; Khot, T.; et al. DiscoveryWorld: A Virtual Environment for Developing and Evaluating Automated Scientific Discovery Agents. In Proceedings of the NeurIPS, 2024.
75. Nie, H.; Wang, Y.; Zhou, M.; et al. Adaptive Preference Arithmetic: Modeling Dynamic Preference Strengths for LLM Agent Personalization. In Proceedings of the NeurIPS, 2025.
76. Swanson, K.; Wu, W.; Bulaong, N.L.; et al. The Virtual Lab of AI agents designs new SARS-CoV-2 nanobodies. *Nature* **2025**, *646*, 716–723.

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.