

A Survey on Agent Skills: Externalized Procedural Knowledge in Language Models

Anonymous ACL submission

Abstract

Large language model (LLM) agents increasingly struggle with reliability, scalability, and execution stability in long-horizon tasks such as software engineering, web automation, and desktop interaction. Recent research addresses these limitations through *agentic skills*: reusable procedural capabilities that externalize execution knowledge into modular, executable, and portable artifacts. This survey presents a unified systems-level review of the emerging agentic skills ecosystem. We formalize skills as externalized procedural abstractions bridging high-level cognitive planning with deterministic execution environments, and organize existing research into a lifecycle-oriented taxonomy spanning discovery, authoring, retrieval, orchestration, execution, adaptation, evaluation, governance, and deprecation. We further analyze security risks including prompt injection, malicious skill packages, and privilege escalation, while outlining open challenges in continual learning, cross-platform portability, and standardized skill ecosystems. This survey positions agentic skills as a foundational abstraction for building scalable, reusable, and reliable autonomous language agents.

1 Introduction

1.1 Motivation: The Reliability and Reusability Bottlenecks in Language Agents

Large language model (LLM) agents representing stateful, autonomous entities have emerged as a popular approach for automating complex, open-ended workflows (Luo et al., 2025). When evaluated on complex, long-horizon software engineering benchmarks, desktop operating system environments, and multi-modal task suites, contemporary agent architectures exhibit a critical performance-limiting reliability cliff (Jimenez et al., 2024). Systemic tracing of these failures indicates that agents

are highly vulnerable to early-stage execution and path resolution errors that cascade and propagate downstream (Wang et al., 2026g; Liu et al., 2026f). This error propagation creates an irreversible no-recovery bottleneck, where agents enter endless reasoning loops or exceed platform recursion limits without attempting correct solutions (Pushkin and Abbe, 2026).

This lack of execution stability is compounded by state drift, wherein natural language representations contaminate the context window over time, and models exhibit sycophantic behaviors (Malmqvist, 2024). Traditional mitigation strategies relying on empirical prompt engineering are fragile and inflate inference costs (Chacko et al., 2026). In contrast, externalized executable skills are a useful emerging abstraction separating high-level cognitive planning from deterministic, modular, and portable execution routines (Liu et al., 2024). Recent comprehensive surveys of harness engineering further underscore the necessity of a unified protocol for standardizing these memories and capabilities (Zhou et al., 2026a; Xing et al., 2026).

1.2 What Is an Agentic Skill? Formalizing the Abstraction

To establish a rigorous theoretical foundation, an externalized executable skill is formalized as a six-tuple:

$$s = (\mathcal{A}, \mathcal{I}, \mathcal{C}, \mathcal{T}, \pi, \mathcal{E}) \quad (1)$$

where \mathcal{A} represents the activation or trigger condition, \mathcal{I} represents procedural instructions, \mathcal{C} represents applicability context or constraints, \mathcal{T} represents accessible tools and resources, π represents the execution policy, and \mathcal{E} represents expected outcomes or effects. Unlike atomic tools or stateless APIs, which represent simple functional primitives, a skill is a stateful, procedural routine coordinating multiple tool invocations under task-specific

Table 1: Comparison of computational abstractions.

Abstraction	Execution Semantics	Statefulness	Compositionality	Persistence	Invocation Style	Runtime Coupling	Adaptability	Reference
System Prompts	Declarative	Stateless	Low	Latent	Monolithic	Tight	Low	(Yao et al., 2023; Wei et al., 2022; Chacko et al., 2026)
Episodic Memories	Semantic	Stateful	Low	Transient	Associative	Loose	High	(Packer et al., 2023a; Park et al., 2023)
Workflows / DAGs	Procedural	Stateful	High	Static	Orchestrated	Tight	Low	(Hong et al., 2023; Wu et al., 2023)
Tools / APIs	Functional	Stateless	High	Persistent	Atomic	Loose	Low	(Schick et al., 2023; Patil et al., 2024; Yuan et al., 2024)
Agentic Skills	Hybrid	Stateful	High	Persistent	Contextual	Sandboxed	High	(Wang et al., 2023; Liu et al., 2024)

constraints (Yuan et al., 2024). Skills are distinct from workflows, system prompts, episodic memory, and policies, providing a modular encapsulated capability structure (Thoppilan et al., 2022). This abstraction generalizes system prompts, workflows, tool-use procedures, reasoning traces, executable plans, and procedural memory into a unified capability model (Lu et al., 2026b; Wang et al., 2026a).

1.3 Origins of Skill-Based Agent Design

The development of agentic skills builds upon a technical evolution in model augmentation. The first phase involved static prompting paradigms like ReAct, SayCan, and Inner Monologue, demonstrating that models could perform multi-step reasoning (Sumers et al., 2024). The second phase introduced atomic tool-calling frameworks such as Toolformer, Gorilla, MRKL, APIBank, HuggingGPT, and ViperGPT, prompting language models to emit tokens representing single API calls (Schick et al., 2023; Patil et al., 2024). The third phase transitioned to programmatic loops and persistent state-tracking, as seen in Voyager, MemGPT, Generative Agents, Reflexion, and Self-Refine (Wang et al., 2023). Finally, the modern fourth phase establishes packaged agent skills as a highly capable abstraction for platforms like SWE-Agent, MetaGPT, AutoGPT, WebArena, and Mind2Web, enabling the secure deployment of pre-compiled directories (Song et al., 2403).

1.4 Contributions of This Survey

This survey provides a comprehensive systems-level review of the emerging landscape of agentic skills, organizing a rapidly growing body of literature into a unified lifecycle-oriented taxonomy spanning skill discovery, acquisition, representation, retrieval, execution, orchestration, evolution, governance, and deprecation. We establish a clear conceptual distinction between passive parametric capabilities and externalized procedural knowledge artifacts, formalizing skills as reusable procedural abstractions that separate high-level reasoning from deterministic execution. Furthermore, we synthesize representative architectures, retrieval mechanisms, memory systems, orchestration strate-

gies, and security frameworks across diverse agentic environments. Finally, we identify key open challenges in scalability, portability, safety, evaluation, and continual learning, outlining a research roadmap that positions agentic skills as a promising but still evolving systems abstraction for long-horizon language agents.

2 Foundations of Skill-Based Agents

2.1 Cognitive Architectures and Procedural Memory

In cognitive psychology, memory is broadly divided into declarative memory, which handles semantic facts, and procedural memory, which governs the procedural mechanics of task execution (Sumers et al., 2024). Computational cognitive architectures model human problem-solving as a tight, iterative loop where declarative knowledge is compiled into procedural rules over time (Sumers et al., 2024). Within the context of language agents, agentic skills serve as the non-parametric procedural memory of the autonomous system. While parametric model weights represent a static semantic memory, and vector databases represent a transient episodic memory, skills provide the modular execution substrate required for repeatable tasks (Packer et al., 2023b). When faced with a familiar problem, the agent retrieves the corresponding procedural skill and executes it directly, bypassing slow reasoning steps (Sumers et al., 2024).

2.2 Skills as Externalized Procedural Knowledge

The core architectural innovation lies in the externalization of this procedural memory. Attempting to feed extensive procedural instructions directly into the context window inflates inference latency and degrades accuracy (Chacko et al., 2026). Externalizing procedural knowledge into packaged files resolves these limitations. Furthermore, externalization enables the system-level principle of progressive disclosure, where the agent reads only the lightweight metadata first, loading full implementation code only during active execution (Yuan et al., 2024). Because these skills are stored as

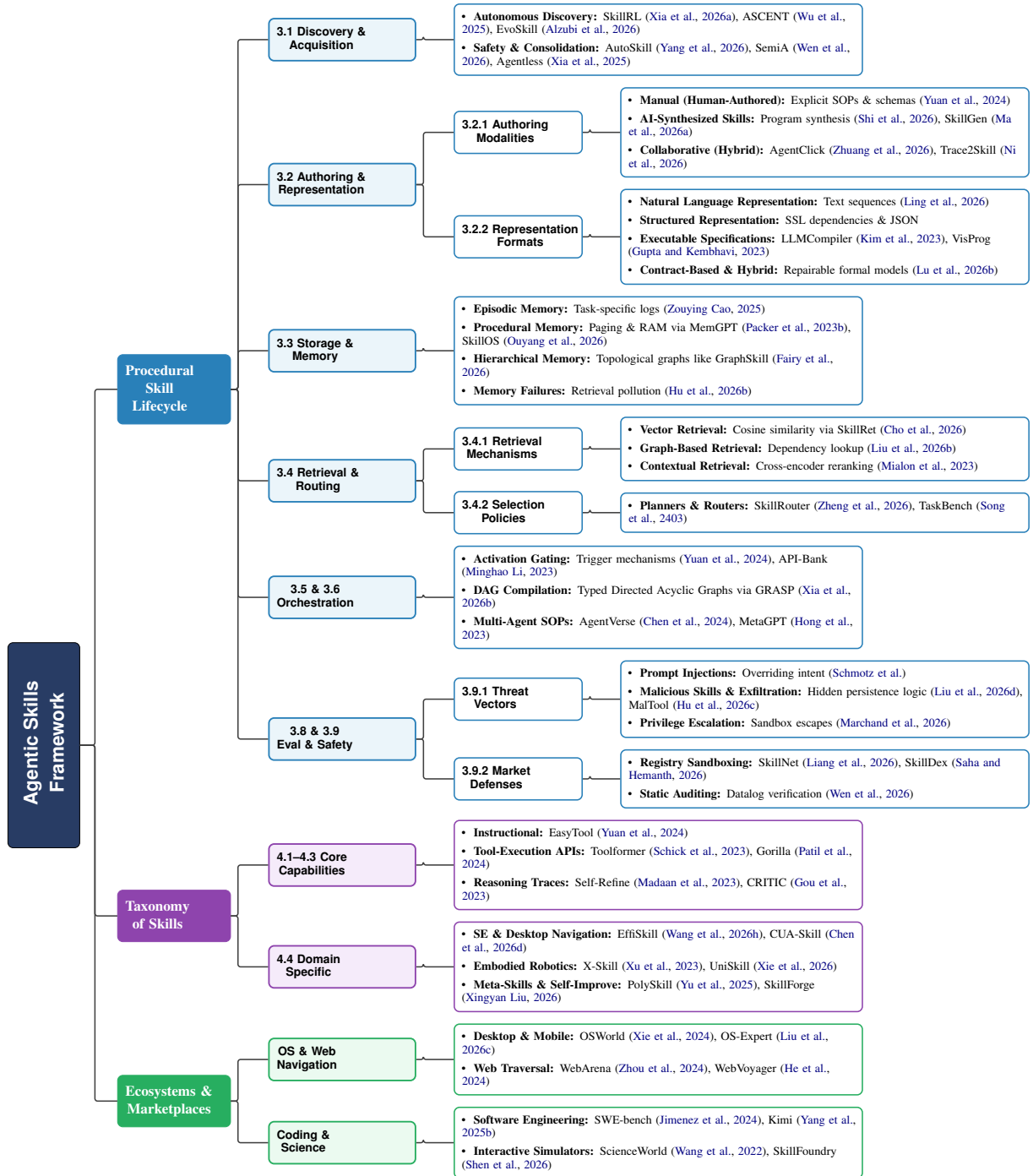


Figure 1: Taxonomy and Landscape of the Agentic Skills Framework, categorizing key system methodologies, lifecycle stages, and domain environments.

plain-text assets, they can be updated without requiring expensive model retraining loops (Zhuang et al., 2026).

We formalize the core primitives of this framework under a unified mathematical abstraction: Definition 1 (Skill): A skill s is a reusable procedural abstraction consisting of activation conditions (\mathcal{A}), procedural instructions (\mathcal{I}), execution constraints (\mathcal{C}), tool interfaces (\mathcal{T}), execution policy

(π), and intended effects (\mathcal{E}) (Yuan et al., 2024):

$$s = (\mathcal{A}, \mathcal{I}, \mathcal{C}, \mathcal{T}, \pi, \mathcal{E}) \quad (2)$$

Definition 2 (Agent State): The agent state at time t is denoted as $x_t \in \mathcal{X}$, capturing the conversation history, internal memory, environment observations, tool outputs, active goals, and retrieved documents (Packer et al., 2023b). Definition 3 (Skill Activation): A skill activates conditionally

185 based on the agent state and the current task goal
186 g :

$$187 \mathcal{A}(x_t, g) \rightarrow \{0, 1\} \quad (3)$$

188 where $\mathcal{A}(x_t, g) = 1$ indicates that the skill should
189 trigger, formalizing routing, contextual applicabil-
190 ity, and planner-guided invocation (Zheng et al.,
191 2026). Definition 4 (Skill Execution): A skill in-
192 duces a conditional state transition:

$$193 x_{t+1} = s(x_t) \quad \text{if } \mathcal{A}(x_t, g) = 1 \quad (4)$$

194 More explicitly:

$$195 x_{t+1} = \pi(x_t, \mathcal{I}, \mathcal{T}) \quad \text{if } \mathcal{A}(x_t, g) = 1 \quad (5)$$

196 This formalizes procedural execution, tool interac-
197 tion, and reasoning under conditional invocation.

198 2.3 Structural Components of Skill Execution

199 From a systems perspective, an executable skill is
200 structured into four distinct functional layers. The
201 declarative layer defines semantic metadata and
202 strict parameter schemas (Yuan et al., 2024). The
203 cognitive layer provides prompt-based instructions
204 outlining the step-by-step plan and explicit error-
205 recovery procedures. The executable layer consists
206 of scripts or compiled binaries that execute proce-
207 dural operations deterministically (Karpas et al.,
208 2022). The sandbox layer defines the isolated con-
209 tainerized runtime environment required to run the
210 assets securely.

211 3 Procedural Skill Lifecycle in LLM 212 Agents

213 The procedural skill lifecycle represents the core
214 operational architecture of modern agent platforms,
215 governing how skills transition from initial creation
216 to runtime deployment, adaptation, and eventual
217 decommissioning. Managing these stages dynam-
218 ically is critical to maintaining a reliable, safe, and
219 efficient agent capability ecosystem (Yang et al.,
220 2026). This section details each stage of the lifecy-
221 cle, unifying them under a consistent mathematical
222 framework.

223 3.1 Skill Discovery and Acquisition

224 When faced with novel tasks, an agent cannot rely
225 solely on pre-existing skill libraries, necessitating
226 the capability to autonomously discover and ac-
227 quire new procedural skills (Yang et al., 2026).
228 Skill discovery is the process of identifying novel,
229 repeatable behaviors and abstracting them from a

230 continuous stream of unstructured agent state tra-
231 jectories (Zhang et al., 2026b; Wu et al., 2025;
232 Alzubi et al., 2026; Yang et al., 2025a). We for-
233 malize the discovery operator D as mapping a his-
234 torical trajectory log \mathcal{H} to a set of candidate skills.
235 Discovery optimizes the partition of a long-horizon
236 trajectory into a sub-trajectory that yields a persis-
237 tent environmental outcome or minimizes future
238 state transition entropy (Xia et al., 2026a; Shu et al.,
239 2017; Huang et al., 2026). This includes cumula-
240 tive capability creation through autonomous devel-
241 opment (Huang et al., 2026) and deriving skills
242 through pure exploration paired with iterative feed-
243 back (Yang et al., 2025a). Multi-agent frameworks
244 also introduce automated discovery mechanisms
245 where systems collaboratively isolate novel behav-
246 iors (Alzubi et al., 2026). For physical systems,
247 algorithms like ASCENT leverage foundation mod-
248 els to derive skills directly from complex embod-
249 ied interactions (Wu et al., 2025). Once a sub-
250 trajectory is isolated, the applicability context un-
251 der which this behavior remains stable is defined,
252 outlining the initial trigger conditions. This option
253 discovery framework allows language agents to dy-
254 namically carve out modular procedural behaviors
255 from complex, interactive task loops (Xia et al.,
256 2026a).

257 Skill acquisition is the process of formalizing,
258 optimizing, and consolidating a discovered skill
259 candidate into the persistent procedural memory
260 store (Yang et al., 2026). The acquisition pipeline
261 optimizes the instructions and policy to maximize
262 expected success under the applicability context,
263 subject to the safety verification constraint, ensur-
264 ing that the newly consolidated routine does not
265 introduce malicious vulnerabilities or privilege es-
266 calations (Wen et al., 2026). This consolidation
267 approach is critical on repository-level software en-
268 gineering benchmarks, where direct, localized tool
269 execution achieves robust success rates (Xia et al.,
270 2025).

271 3.2 Skill Authoring and Representation

272 Once a skill is discovered or engineered, represent-
273 ing and serializing its artifacts for heterogeneous
274 runtimes is a critical design challenge.

275 3.2.1 Skill Authoring Modalities

276 Agentic frameworks employ three primary author-
277 ing modalities to define procedural assets. First,
278 manual human-authored skills allow engineers to
279 explicitly author standard operating procedures

(SOPs), declarative metadata schemas, and native scripts (Yuan et al., 2024). This manual approach ensures high safety, deterministic predictability, and exact correctness, but presents a severe human engineering bottleneck and struggles with horizontal scale. Second, AI-synthesized skills empower the agent to autonomously generate its own skills without human intervention (Yang et al., 2026; Jiao et al., 2026; Wang et al., 2025b). This self-generation paradigm leverages techniques like program synthesis networks (PSN) or verifier-guided code generation where the model writes, evaluates, and registers its own code assets and triggering rules (Shi et al., 2026). Advanced frameworks employ contrastive induction over successful and failed trajectories to programmatically synthesize standardized and verified skill directories at runtime (Ma et al., 2026a). Third, collaborative hybrid-authored skills establish a paradigm where human operators collaborate with the model in a human-in-the-loop setup (Ni et al., 2026). Layers like AgentClick allow continuous human verification of terminal actions prior to strict integration (Zhuang et al., 2026). Humans typically define the declarative schemas and safety constraints, while the model generates the underlying procedural implementation and verifies it through sandboxed unit testing (Ni et al., 2026).

3.2.2 Skill Representation Formats

To enable portability, several representation paradigms have been introduced (Ling et al., 2026). Natural language representation remains the primary modality for representing agentic skills, where procedural instructions are represented textually as a sequence of procedural steps (Ling et al., 2026). The activation specification may also be encoded textually, and the model reads these instructions to dynamically align its cognitive planning with task-specific operational constraints (Yuan et al., 2024). Structured representation offers a solution to natural language verbosity, using XML, JSON, or domain-specific languages (DSLs). Under this paradigm, a structured skill is represented as a dependency graph defining execution dependencies, disentangling scheduling signals, execution structure, and logical evidence into distinct schemas. Executable specifications go beyond declarative metadata, requiring formal programmatic specifications to provide a deterministic execution wrapper that schedules parallel tool invocations and handles errors programmatically (Kim

et al., 2023; Suris et al., 2023; Gupta and Kembhavi, 2023). Contract-based and unified models improve reliability by establishing formal guarantees for multimodal web agents (Lu et al., 2026b), advocating for an experience compression spectrum that seamlessly unifies memory, rules, and actionable capabilities (Zhang et al., 2026d). Hybrid representations reconcile cognitive flexibility with systems-level safety by combining natural language guidelines, symbolic constraints, and executable native code into a single unified package (Bi et al., 2026; Shi et al., 2026). Finally, skill metadata guarantees interoperability through rich semantic versioning (Liang et al., 2026), while skill embedding supports dynamic retrieval by mapping a skill to a low-dimensional vector representation using an encoder model, providing a continuous representation space for semantic similarity operations (Cho et al., 2026).

3.3 Skill Storage and Memory

An agent’s operational expertise is stored within distinct temporal memory tiers to maintain context efficiency and prevent cognitive overload (Sumers et al., 2024; Sun et al., 2026; Tu et al., 2026; Zhang et al., 2026d,c; Lu et al., 2026a; Zhu et al., 2026; Ouyang et al., 2026). Episodic skill memory represents the agent’s historical log of task-specific experiences, recording past execution traces, successes, and failed trajectories (Zouying Cao, 2025). Modern platforms extract fine-grained lessons through multi-faceted distillation of these traces to isolate transferable lessons (Ni et al., 2026). Long-term procedural memory serves as the persistent, non-parametric storage substrate for an agent’s operational expertise, distinguishing practical knowledge from semantic facts (Sumers et al., 2024). By managing the pre-compiled repository of skills through an active virtual memory paging system, agents can scale their operational expertise significantly (Packer et al., 2023b). Rather than purely retrieving isolated facts, modern RAG systems directly distill enterprise knowledge into these navigable procedural formats (Sun et al., 2026). While vector search is highly effective for coarse-grained matches, capturing the hierarchical dependencies between complex, multi-file skills requires a topological execution graph or hierarchical memory tree to identify candidate routines securely (Fairy et al., 2026). As an agent’s memory store grows across sequential tasks, the storage process encounters severe memory failures, including retrieval pol-

lution, context competition, and memory dilution (Hu et al., 2026b; Fang et al., 2026).

3.4 Skill Retrieval and Selection

Retrieving and selecting the optimal skill (or subset of skills) to route a task is driven by a tight, multi-stage runtime selection pipeline (Zheng et al., 2026). Vector retrieval represents the most widely adopted paradigm for dynamically exposing relevant capability packages to agents on-demand. At inference time, the agent platform computes the cosine similarity between the embedded user query and the skill vector store, retrieving the best skill restricted to the active triggering subset to provide fast, low-latency, and context-aware access (Zheng et al., 2026; Cho et al., 2026). Advanced graph-based retrieval incorporates dependency-aware structural mechanisms to filter massive registries efficiently (Liu et al., 2026b). For complex coding scenarios, documentation-guided hierarchical retrieval augments standard lookup paths by strictly following graph reasoning steps to ensure prerequisite and dependent sub-skills are sequenced correctly (Wang et al., 2026b; Fairy et al., 2026). Contextual retrieval utilizes a lightweight sparse or dense vector query to obtain candidates, subsequently deploying a heavy cross-encoder reranker to reason over the full skill instructions and parameters (Cho et al., 2026; Mialon et al., 2023). Once candidate skills are retrieved, the platform dictates selection policies strictly constrained by active triggering rules, utilizing sparse-dense routers and task-specific classification policies to resolve semantic confusability (Zheng et al., 2026; Su et al., 2026; Wang et al., 2026d). Finally, planner-guided routing parses natural language intents into optimized execution plans, scheduling parallel skill invocations and batching requests to reduce latency (Kim et al., 2023; Song et al., 2023; Chen et al., 2026a).

3.5 Skill Activation and Execution

Skill activation and execution manage the active invocation, runtime tracking, and safety boundary enforcement of selected procedural routines (Chen et al., 2026c). Trigger mechanisms define how a skill activates conditionally based on the agent state and the current task goal (Zheng et al., 2026; Minghao Li, 2023). Contextual gating ensures the cognitive planner can gate execution so the agent only acts under verified contexts, utilizing explicit pre-conditions, validation gates, and expected out-

puts (Yuan et al., 2024). To prevent malicious triggers, permission-aware activation enforces a strict subset of allowed tool invocations, injecting safe activation policies that verify trigger-level safety before execution based on user credentials (Li et al., 2026e). Runtime execution then induces a conditional state transition. To prevent damage to the host system, execution must strictly occur in isolated sandbox containers. Finally, error recovery processes actively monitor execution traces and standard error logs to trigger automated recovery policies and recursive backtracking if a step fails (Li et al.).

3.6 Skill Composition and Orchestration

As tasks scale in complexity, single-skill execution shifts toward advanced composition and orchestration protocols (Xia et al., 2026b; Zabounidis et al., 2026; Xia et al., 2026c; Fan et al., 2026). Sequential composition allows agents to chain multiple distinct skills into a linear execution pipeline, where the output state of a preceding skill directly populates the input parameters of the subsequent routine (Li et al., 2026a; Microsoft, 2024). Hierarchical composition coordinates the workflow by utilizing a master skill to delegate sub-tasks to atomic sub-skills, isolating context windows effectively (Xia et al., 2026b; Du et al., 2024). Dynamic DAG orchestration addresses advanced platform needs by compiling flat retrieved skill candidates into typed Directed Acyclic Graphs, resolving exceptions locally without triggering full-system replanning (Xia et al., 2026b). Recent approaches deploy large language models to guide symbolic planning, subsequently grounding the composed structures via deep reinforcement learning (Zabounidis et al., 2026), while agentic proposing mechanisms enhance reasoning by synthesizing compositional skills dynamically (Jiao et al., 2026; Xia et al., 2026c). Multi-agent coordination shifts orchestration from central scheduling to decentralized coordination among specialized, role-playing agents (Chen et al., 2024; Nie et al., 2026; Li, 2026). These systems define standard operating procedures (SOPs) enabling communicative cooperation while preventing context clashes (Hong et al., 2023; Qian et al., 2024). Finally, recursive invocation enables skills to invoke themselves or their ancestral nodes to resolve nested problems, which is heavily modeled during self-improving reinforcement learning rollouts (Xia et al., 2026a).

482 **3.7 Skill Adaptation and Evolution**

483 Deployed skills require constant refinement, shifting
484 from short-term context adjustments to long-
485 term reinforcement self-improvement (Zhang et al.,
486 2026b; Li et al., 2026d; Wu et al., 2026). Short-
487 term adaptation allows agents to adapt retrieved
488 generic skills dynamically, refining execution pa-
489 rameters based on verifier feedback (Du and Pinck-
490 ney, 2026). Feedback-driven refinement utilizes
491 reasoning skills to actively analyze previous exe-
492 cution logs, identifying errors and recursively re-
493 fining generated plans (Madaan et al., 2023; Gou
494 et al., 2023). Long-term evolution ensures that an
495 agent’s capability library continually self-improves
496 through operational feedback (Vishe et al., 2026;
497 Ye et al., 2026; Du et al., 2025; Wang et al., 2025a;
498 Ma et al., 2026b; Wang et al., 2026c; Xu et al.,
499 2026b; Mi et al., 2026). To handle complex goals,
500 modern frameworks co-evolve high-level decision
501 policies alongside expanding skill banks (Wu et al.,
502 2026). This intrinsic evolution integrates deeply
503 with hierarchical reinforcement learning protocols,
504 allowing agents to iteratively refine internal logic
505 structures (Li et al., 2026d). To support this, self-
506 improving skills optimize operational primitives
507 and instructions based on historical feedback traces
508 in a lifelong learning cycle, allowing agents to
509 scale expertise dynamically (Yang et al., 2026; Liu
510 et al., 2023; Zhang et al., 2026e; Tziafas and Kasaei,
511 2024; Jiang et al., 2026a).

512 **3.8 Skill Evaluation**

513 Evaluating performance is essential to establish
514 benchmarks for systems reliability, token economy,
515 and human-agent alignment (Li et al., 2026c; Han
516 et al., 2026; Zhong et al., 2026; Jiang et al., 2026b;
517 Liu et al., 2026f). Correctness evaluates binary re-
518 ward functions and conditional expected success
519 under specific triggering constraints, providing a re-
520 liable correctness metric (Li et al., 2026c). General-
521 ization evaluates the capability to deploy routines in
522 zero-shot situations where no specific match exists,
523 measuring performance differences across domains.
524 Robustness measures resilience to execution ex-
525 ceptions and adversarial inputs (Liu et al., 2026d).
526 Composability verifies the joint success of chaining
527 coordinated graph routines (Xia et al., 2026b). Ef-
528 ficiency focuses on execution costs, incorporating
529 token usage, inference latency, and tool overhead
530 (Gao et al., 2026; Chen et al., 2026a). Interpretabil-
531 ity and usability provide structural transparency,

allowing users to form bounded expectations and
construct local verification checks directly from
the procedural specifications (Wen, 2026). In the
wild benchmarking in realistic, unconstrained set-
tings is critical to understanding true operational
reliability and failure modes outside sterile labo-
ratories (Liu et al., 2026f), while standardized se-
curity suites like HarmfulSkillBench measure how
compromised files weaponize autonomous systems
(Jiang et al., 2026b). High transferability allows
decentralized registries to dynamically share and
run verified procedures across diverse endpoints
(Chen et al., 2026a; Liang et al., 2026).

532 **3.9 Skill Governance and Safety**

533 As agentic modules execute systems-level com-
534 mands, securing these packages from adversarial
535 exploits is of paramount importance (Li et al.,
536 2026e,f). We map security risks to specific stages
537 of the procedural lifecycle, tracking threat vectors
538 across acquisition, distribution, retrieval, and exe-
539 cution (Li et al., 2026e; Tie et al., 2026). Prompt
540 injections and model poisoning present severe risks
541 where attacking payloads embedded in files or web
542 pages can override instructions, forcing the model
543 to invoke high-impact tools without authorization
544 (Schmotz et al., 2025; Maloyan and Namiot, 2026;
545 Schmotz et al., 2026; Jia et al., 2026). Attackers
546 increasingly deploy model-in-skill poisoning tech-
547 niques to create stealthy backdoor exploits, fun-
548 damentally compromising the execution integrity
549 of shared registries (Tie et al., 2026). Malicious
550 skills and exfiltration frameworks distribute hidden
551 logic on community registries targeting filesystems,
552 registry keys, and network interfaces to establish
553 persistence (Liu et al., 2026d,e; Holzbauer et al.,
554 2026). Large-scale empirical studies reveal exten-
555 sive security vulnerabilities concerning severe cre-
556 dential leakage and unauthorized file exfiltration
557 to attacker-controlled endpoints (Liu et al., 2026e;
558 Chen et al., 2026e; Hu et al., 2026c; Wang et al.,
559 2026f). Privilege escalation occurs when an exe-
560 cutable escapes its designated boundary, demand-
561 ing that safe activation policies be verified before
562 execution (Li et al., 2026e). Highly capable plan-
563 ners can systematically discover kernel flaws to ex-
564 ecute host escapes from misconfigured sandboxes
565 (Marchand et al., 2026; Xiao et al., 2026; Duan
566 et al., 2026). Marketplace security relies on reg-
567 istries establishing automated dynamic sandboxing
568 and cryptographic signatures (Li et al., 2026e; ?;
569 Hu et al., 2026a; Saha and Hemanth, 2026; Bhard-
570 571 572 573 574 575 576 577 578 579 580 581 582

waj; Guo et al., 2026; Hou and Yang, 2026; Feng et al., 2026; Zhang et al., 2026a; Lv et al., 2026; Qu et al., 2026). Finally, auditing and synthesis frameworks deploy static analyzers using verification functions to mathematically prove safety compliance before execution (Wen et al., 2026).

3.10 Skill Deprecation and Forgetting

As libraries expand, they accumulate redundant files necessitating automated pruning, memory compression, and retirement policies to maintain operational registry solvency (Pu et al., 2026; Hu et al., 2026b).

4 Open Challenges and Future Directions

As the landscape matures, several structural limitations define the technical frontier. Skills must execute reliably across heterogeneous model APIs. Furthermore, deploying these complex skill frameworks on small language models presents a critical frontier for industrial and edge environments, where computational constraints mandate highly efficient execution and unified serialization standards (Xu et al., 2026a; Ling et al., 2026; Chen et al., 2026b; Wang et al., 2026e). Performing precise selection at scale requires sparse-dense hierarchical routers to avoid semantic confusability (Zheng et al., 2026). Similarly, orchestrating complex dependencies requires compiling files into typed DAGs to support dynamic error repair (Xia et al., 2026b). Enforcing robust capability-based permission models is essential to prevent container breakout exploits (Marchand et al., 2026; Li et al., 2026e). Overcoming verification trust gaps requires static auditors that lift prose into Datalog databases to mathematically prove safety compliance (Wen et al., 2026). Enabling acquisition without triggering catastrophic forgetting requires decoupling learned abstractions via virtual memory paging (Yang et al., 2026; Packer et al., 2023b). Finally, constructing robust, scalable, interactive environments is essential to safely evaluate execution efficiency and safety guardrails under massive distribution scale (Li et al., 2026c).

5 Conclusion

The transition from monolithic prompting and atomic tool use to externalized executable skills represents a critical maturation in the design of autonomous language agents. By decoupling high-level cognitive planning from deterministic, proce-

dural execution, skill-based architectures directly address the context scaling and reliability bottlenecks that plague contemporary LLM systems. This survey has formalized the agentic skill abstraction and provided a unified lifecycle-oriented taxonomy, synthesizing the latest advancements in skill discovery, representation, memory retrieval, orchestration, and governance. As agents are increasingly deployed in complex ecosystems, advancing robust permission models, dynamic adaptation, and standardized evaluation will be paramount. Ultimately, formalizing agent capabilities as modular, secure, and self-improving skills paves the way for the next generation of scalable and reliable autonomous systems.

6 Limitations

While this survey provides a comprehensive synthesis of the agentic skills landscape, several limitations exist. First, the field of autonomous language agents is evolving at a rapid pace; thus, specific frameworks, models, and retrieval implementations discussed may face deprecation as industry standards consolidate. Second, while we focus extensively on digital, software-oriented, and desktop capabilities, the integration of multimodal and embodied robotic skills introduces distinct physical grounding challenges that are only briefly surveyed here. Finally, the evaluation of agentic skills remains an open challenge. Current benchmarks often struggle to capture the open-ended complexity and cascading failure modes of real-world operational environments, meaning the true reliability of these externalized systems in production remains difficult to quantify completely.

References

- Salaheddin Alzubi, Noah Provenzano, Jaydon Bingham, Weiyuan Chen, and Tu Vu. 2026. *Evoskill: Automated skill discovery for multi-agent systems*. Preprint, arXiv:2603.02766.
- Varun Pratap Bhardwaj. Skillfortify: Formal analysis and supply chain security for agentic ai skills, 2026.
- Shuzhen Bi, Mengsong Wu, Hao Hao, Keqian Li, Wentao Liu, Siyu Song, Hongbo Zhao, and Aimin Zhou. 2026. *Automating skill acquisition through large-scale mining of open-source agentic repositories: A framework for multi-agent procedural knowledge extraction*. Preprint, arXiv:2603.11808.
- Samuel Jacob Chacko, James Hugglestone, Chashi Mahiul Islam, and Xiuwen Liu. 2026.

680	When skills don't help: A negative result on procedural knowledge for tool-grounded agents in offensive cybersecurity. <i>Preprint</i> , arXiv:2605.20023.	Zenghao Duan, Yuxin Tian, Zhiyi Yin, Liang Pang, Jingcheng Deng, Zihao Wei, Shicheng Xu, Yuyao Ge, and Xueqi Cheng. 2026. Skillattack: Automated red teaming of agent skills through attack path refinement. <i>arXiv preprint arXiv:2604.04989</i> .	732
681			733
682			734
683	Le Chen, Erhu Feng, Yubin Xia, and Haibo Chen. 2026a. Skvm: Revisiting language vm for skills across heterogenous llms and harnesses. <i>Preprint</i> , arXiv:2604.03088.		735
684			736
685		Fali Fairy and 1 others. 2026. Graphskill: Documentation-guided hierarchical retrieval-augmented coding for complex graph reasoning. <i>arXiv preprint arXiv:2603.06620</i> .	737
686			738
687	Qijia Chen, Andrea Bellucci, Zhida Sun, and Giulio Jacucci. 2026b. Skilldroid: Compile once, reuse forever. <i>arXiv preprint arXiv:2604.14872</i> .		739
688			740
689			741
690	Shiqi Chen, Jingze Gai, Ruochen Zhou, Jinghan Zhang, Tongyao Zhu, Junlong Li, Kangrui Wang, Zihan Wang, Zhengyu Chen, Klara Kaleb, and 1 others. 2026c. Skillcraft: Can llm agents learn to use tools skillfully? <i>arXiv preprint arXiv:2603.00718</i> .	Zhiyuan Fan, Tinghao Yu, Yuanjun Cai, Jiangtao Guan, Yun Yang, Dingxin Hu, Jiang Zhou, Xing Wu, Zhuo Han, Feng Zhang, and 1 others. 2026. Toward scalable terminal task synthesis via skill graphs. <i>arXiv preprint arXiv:2604.25727</i> .	742
691			743
692			744
693			745
694		Runnan Fang, Yuan Liang, Xiaobin Wang, Jialong Wu, Shuofei Qiao, Pengjun Xie, Fei Huang, Huajun Chen, and Ningyu Zhang. 2026. Memp: Exploring agent procedural memory. <i>Preprint</i> , arXiv:2508.06433.	746
695	Tianyi Chen, Yinheng Li, Michael Solodko, Sen Wang, Nan Jiang, Tingyuan Cui, Junheng Hao, Jongwoo Ko, Sara Abdali, Leon Xu, Suzhen Zheng, Hao Fan, Pashmina Cameron, Justin Wagle, and Kazuhito Koishida. 2026d. Cua-skill: Develop skills for computer using agent. <i>Preprint</i> , arXiv:2601.21123.		747
696			748
697			749
698		Yunhao Feng, Yifan Ding, Yingshui Tan, Boren Zheng, Yanming Guo, Xiaolong Li, Kun Zhai, Yishan Li, and Wenke Huang. 2026. Skilltrojan: Backdoor attacks on skill-based agent systems. <i>arXiv preprint arXiv:2604.06811</i> .	750
699			751
700			752
701	Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chi-Min Chan, Heyang Yu, Yaxi Lu, Yi-Hsin Hung, Chen Qian, Yujia Qin, Xin Cong, Ruobing Xie, Zhiyuan Liu, Maosong Sun, and Jie Zhou. 2024. Agentverse: Facilitating multi-agent collaboration and exploring emergent behaviors. In <i>The Twelfth International Conference on Learning Representations</i> .	Yudong Gao, Zongjie Li, Yuanyuanyuan, Zimo Ji, Pingchuan Ma, and Shuai Wang. 2026. Skillreducer: Optimizing llm agent skills for token efficiency. <i>Preprint</i> , arXiv:2603.29919.	753
702			754
703			755
704			756
705			757
706			758
707		Zhibin Gou and 1 others. 2023. Critic: Large language models can self-correct with tools. <i>arXiv preprint arXiv:2305.11738</i> .	759
708			760
709	Zihao Chen, Ying Zhang, Yi Liu, Gelei Deng, Yuekang Li, Yanjun Zhang, Jianting Ning, Leo Yu Zhang, Lei Ma, and Zhiqiang Li. 2026e. Credential leakage in llm agent skills: A large-scale empirical study. <i>Preprint</i> , arXiv:2604.03070.		761
710			762
711			763
712			764
713			765
714	Hongcheol Cho, Ryangkyung Kang, and Youngeun Kim. 2026. Skillret: A large-scale benchmark for skill retrieval in llm agents. <i>Preprint</i> , arXiv:2605.05726.	Tanmay Gupta and Aniruddha Kembhavi. 2023. Visual programming: Compositional visual reasoning without training. In <i>IEEE Conference on Computer Vision and Pattern Recognition (CVPR)</i> .	766
715			767
716			768
717	Xiang Deng, Yu Gu, and 1 others. 2023. Mind2web: Towards a generalist agent for the web. <i>arXiv preprint arXiv:2306.06070</i> .		769
718			770
719			771
720	Jiawei Du, Jinlong Wu, Yuzheng Chen, Yucheng Hu, Bing Li, and Joey Tianyi Zhou. 2025. Rethinking agent design: From top-down workflows to bottom-up skill evolution. <i>arXiv preprint arXiv:2505.17673</i> .	Tingxu Han, Yi Zhang, Wei Song, Chunrong Fang, Zhenyu Chen, Youcheng Sun, and Lijie Hu. 2026. Swe-skills-bench: Do agent skills actually help in real-world software engineering? <i>arXiv preprint arXiv:2603.15401</i> .	772
721			773
722			774
723			775
724	Yu Du, Fangyun Wei, and Hongyang Zhang. 2024. Anytool: self-reflective, hierarchical agents for large-scale api calls. In <i>Proceedings of the 41st International Conference on Machine Learning, ICML'24. JMLR.org</i> .	Hongliang He, Wenlin Yao, Kaixin Ma, Wenhao Yu, Yong Dai, Hongming Zhang, Zhenzhong Lan, and Dong Yu. 2024. Webvoyager: Building an end-to-end web agent with large multimodal models. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 6864–6890.	776
725			777
726			778
727			779
728			780
729	Zijian Du and Nathaniel Pinckney. 2026. Trace2skill: Verifier-guided skill evolution for long-context eda agents. <i>Preprint</i> , arXiv:2605.21810.	Florian Holzbauer, David Schmidt, Gabriel Gegenhuber, Sebastian Schrittwieser, and Johanna Ullrich. 2026. Malicious or not: Adding repository context to agent skill classification. <i>arXiv preprint arXiv:2603.16572</i> .	781
730			782
731			783
			784
			785
			786

787	Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, and 1 others. 2023. Metagtpt: Meta programming for a multi-agent collaborative framework. <i>arXiv preprint arXiv:2308.00352</i> .		
788			842
789			843
790			844
791			
792			
793	Yinghan Hou and Zongyou Yang. 2026. Skillsieve: A hierarchical triage framework for detecting malicious ai agent skills. <i>arXiv preprint arXiv:2604.06550</i> .		
794			848
795			849
796	Haichuan Hu, Ye Shang, and Quanjun Zhang. 2026a. Red skills or blue skills? a dive into skills published on clawhub. <i>arXiv preprint arXiv:2604.13064</i> .		850
797			851
798			852
799	Qisheng Hu, Quanyu Long, and Wenya Wang. 2026b. When continual learning moves to memory: A study of experience reuse in llm agents. <i>Preprint</i> , arXiv:2604.27003.		
800			853
801			854
802			855
803	Yuepeng Hu, Yuqi Jia, Mengyuan Li, Dawn Song, and Neil Gong. 2026c. Maltool: Malicious tool attacks on llm agents. <i>Preprint</i> , arXiv:2602.12194.		856
804			857
805			858
806	Xu Huang, Junwu Chen, Yuxing Fei, Zhuohan Li, Philippe Schwaller, and Gerbrand Ceder. 2026. Cascade: Cumulative agentic skill creation through autonomous development and evolution. <i>Preprint</i> , arXiv:2512.23880.		859
807			860
808			
809			861
810			862
811	Xiaojun Jia, Jie Liao, Simeng Qin, Jindong Gu, Wenqi Ren, Xiaochun Cao, Yang Liu, and Philip Torr. 2026. Skillject: Automating stealthy skill-based prompt injection for coding agents with trace-driven closed-loop refinement. In <i>The 6th Workshop of Adversarial Machine Learning on Computer Vision: Safety of Vision-Language Agents</i> .		863
812			
813			864
814			865
815			866
816			867
817			
818	Guanyu Jiang, Zhaochen Su, Xiaoye Qu, and Yi R Fung. 2026a. Xskill: Continual learning from experience and skills in multimodal agents. <i>arXiv preprint arXiv:2603.12056</i> .		868
819			869
820			870
821			871
822	Yukun Jiang, Yage Zhang, Michael Backes, Xinyue Shen, and Yang Zhang. 2026b. Harmfulskillbench: How do harmful skills weaponize your agents? <i>Preprint</i> , arXiv:2604.15415.		872
823			873
824			
825			874
826	Zhengbo Jiao, Shaobo Wang, Zifan Zhang, Xuan Ren, Wei Wang, Bing Zhao, Hu Wei, and Linfeng Zhang. 2026. Agentic proposing: Enhancing large language model reasoning via compositional skill synthesis. <i>Preprint</i> , arXiv:2602.03279.		875
827			876
828			877
829			
830			878
831	Carlos E Jimenez, John Yang, Alexander Wettinger, and 1 others. 2024. Swe-bench: Can language models resolve real-world github issues? In <i>International Conference on Learning Representations (ICLR)</i> .		879
832			880
833			881
834			
835	Ehud Karpas, Omri Abend, Yonatan Belinkov, Barak Lenz, Opher Lieber, Nir Ratner, Yoav Shoham, Hofit Bata, Yoav Levine, Kevin Leyton-Brown, and 1 others. 2022. Mrkl systems: A modular, neuro-symbolic architecture that combines large language models, external knowledge sources and discrete reasoning. <i>arXiv preprint arXiv:2205.00445</i> .		882
836			883
837			884
838			885
839			886
840			887
841			888
			889
			890
			891
			892
			893
			894
			895
			896
			897

898	Chang Liu, Sib0 Tian, Xiao Liang, and Minghui Zheng.	Lijia Lv, Xuehai Tang, Jie Wen, Jizhong Han, and	954
899	2026a. Self-vla: A skill enhanced agentic vision-	Songlin Hu. 2026. Structured security auditing and	955
900	language-action framework for contact-rich disas-	robustness enhancement for untrusted agent skills.	956
901	sembly . <i>Preprint</i> , arXiv:2603.11080.	<i>arXiv preprint arXiv:2604.25109</i> .	957
902	Dawei Liu, Zongxia Li, Hongyang Du, Xiyang Wu,	Yuchen Ma, Yue Huang, Han Bao, Haomin Zhuang,	958
903	Shihang Gui, Yongbei Kuang, and Lichao Sun.	Swadheen Shukla, Michel Galley, Xiangliang Zhang,	959
904	2026b. Graph of skills: Dependency-aware struc-	and Stefan Feuerriegel. 2026a. Skillgen: Veri-	960
905	tural retrieval for massive agent skills . <i>Preprint</i> ,	fied inference-time agent skill synthesis . <i>Preprint</i> ,	961
906	arXiv:2604.05333.	arXiv:2605.10999.	962
907	Jiateng Liu, Zhenhailong Wang, Rushi Wang, Bingx-	Ziyu Ma, Shidong Yang, Yuxiang Ji, Xucong Wang,	963
908	uan Li, Jeonghwan Kim, Aditi Tiwari, Pengfei Yu,	Yong Wang, Yiming Hu, Tongwen Huang, and Xi-	964
909	Denghui Zhang, and Heng Ji. 2026c. Osexpert:	angxiang Chu. 2026b. Skillclaw: Let skills evolve	965
910	Computer-use agents learning professional skills via	collectively with agentic evolver . <i>arXiv preprint</i>	966
911	exploration . <i>arXiv preprint arXiv:2603.07978</i> .	<i>arXiv:2604.08377</i> .	967
912	Ruibo Liu, Jason Wei, Shixiang Shane Gu, Te-Yen Wu,	Aman Madaan, Niket Tandon, Prakhar Gupta, and 1	968
913	Soroush Vosoughi, Claire Cui, Denny Zhou, and An-	others. 2023. Self-refine: Iterative refinement with	969
914	drew M. Dai. 2023. Mind’s eye: Grounded language	self-feedback . <i>arXiv preprint arXiv:2305.00811</i> .	970
915	model reasoning through simulation . In <i>The Eleventh</i>	Lars Malmqvist. 2024. Sycophancy in large lan-	971
916	<i>International Conference on Learning Representa-</i>	guage models: Causes and mitigations . <i>Preprint</i> ,	972
917	<i>tions</i> .	arXiv:2411.15287.	973
918	Shunyu Liu, Yaoru Li, Kongcheng Zhang, Zhenyu	Narek Maloyan and Dmitry Namiot. 2026. Prompt	974
919	Cui, Wenkai Fang, Yuxuan Zheng, Tongya Zheng,	injection attacks on agentic coding assistants: A sys-	975
920	and Mingli Song. 2025. Odyssey: Empowering	tematic analysis of vulnerabilities in skills, tools, and	976
921	minecraft agents with open-world skills . <i>Preprint</i> ,	protocol ecosystems. <i>International Journal of Open</i>	977
922	arXiv:2407.15325.	<i>Information Technologies</i> , 14(2):1–10.	978
923	Yi Liu, Zhihao Chen, Yanjun Zhang, Gelei Deng,	Rahul Marchand, Art O Cathain, Jerome Wynne, Philip-	979
924	Yuekang Li, Jianting Ning, Ying Zhang, and Leo Yu	pos Maximos Giavridis, Sam Deverett, John Wilkin-	980
925	Zhang. 2026d. Malicious agent skills in the wild: A	son, Jason Gwartz, and Harry Coppock. 2026. Quan-	981
926	large-scale security empirical study . <i>arXiv preprint</i>	tifying frontier llm capabilities for container sandbox	982
927	<i>arXiv:2605.08115</i> .	escape . <i>Preprint</i> , arXiv:2603.02277.	983
928	Yi Liu, Weizhe Wang, Ruitao Feng, Yao Zhang,	Qirui Mi, Zhijian Ma, Mengyue Yang, Haoxuan Li,	984
929	Guangquan Xu, Gelei Deng, Yuekang Li, and Leo	Yisen Wang, Haifeng Zhang, and Jun Wang. 2026.	985
930	Zhang. 2026e. Agent skills in the wild: An empirical	Skill-pro: Learning reusable skills from experience	986
931	study of security vulnerabilities at scale . <i>Preprint</i> ,	via non-parametric ppo for llm agents . <i>arXiv preprint</i>	987
932	arXiv:2601.10338.	<i>arXiv:2602.01869</i> .	988
933	Yujian Liu, Jiabao Ji, Li An, Tommi Jaakkola, Yang	Gregoire Mialon and 1 others. 2023. Augmented	989
934	Zhang, and Shiyu Chang. 2026f. How well do agen-	language models: a survey . <i>arXiv preprint</i>	990
935	tic skills work in the wild: Benchmarking llm skill us-	<i>arXiv:2302.07842</i> .	991
936	age in realistic settings . <i>Preprint</i> , arXiv:2604.04323.	Microsoft. 2024. Prompt flow: Build high-quality llm	992
937	Zhengxi Lu, Zhiyuan Yao, Jinyang Wu, Chengcheng	apps from prototyping to production .	993
938	Han, Qi Gu, Xunliang Cai, Weiming Lu, Jun Xiao,	Bowen Yu Feifan Song Hangyu Li Haiyang Yu	994
939	Yueting Zhuang, and Yongliang Shen. 2026a. Skill0:	Zhoujun Li Fei Huang Yongbin Li Minghao Li,	995
940	In-context agentic reinforcement learning for skill	Yingxiu Zhao. 2023. Api-bank: A comprehensive	996
941	internalization . <i>arXiv preprint arXiv:2604.02268</i> .	benchmark for tool-augmented llms . <i>arXiv preprint</i>	997
942	Zijian Lu, Yiping Zuo, Yupeng Nie, Xin He, Weibei	<i>arXiv:2304.08244</i> .	998
943	Fan, Lianyong Qi, and Shi Jin. 2026b. Contractskill:	Jingwei Ni, Yihao Liu, Xinpeng Liu, Yutao Sun,	999
944	Repairable contract-based skills for multimodal web	Mengyu Zhou, Pengyu Cheng, Dexin Wang, Er-	1000
945	agents . <i>Preprint</i> , arXiv:2603.20340.	chao Zhao, Xiaoxi Jiang, and Guanjun Jiang.	1001
946	Junyu Luo, Weizhi Zhang, Ye Yuan, Yusheng Zhao, Jun-	2026. Trace2skill: Distill trajectory-local lessons	1002
947	wei Yang, Yiyang Gu, Bohan Wu, Binqi Chen, Ziyue	into transferable agent skills . <i>arXiv preprint</i>	1003
948	Qiao, Qingqing Long, Rongcheng Tu, Xiao Luo, Wei	<i>arXiv:2603.25158</i> .	1004
949	Ju, Zhiping Xiao, Yifan Wang, Meng Xiao, Chenwu	Zheng Nie, Ruolin Shen, Xinlei Yu, Bo Yin, Jiangning	1005
950	Liu, Jinyang Yuan, Shichang Zhang, and 7 others.	Zhang, and Xiaobin Hu. 2026. Skillgraph: Self-	1006
951	2025. Large language model agent: A survey on	evolving multi-agent collaboration with multimodal	1007
952	methodology, applications and challenges . <i>Preprint</i> ,	graph topology . <i>arXiv preprint arXiv:2604.17503</i> .	1008
953	arXiv:2503.21460.		

1009	Siru Ouyang, Jun Yan, Yanfei Chen, Rujun Han, Zifeng Wang, Bhavana Dalvi Mishra, Rui Meng, Chun-Liang Li, Yizhu Jiao, Kaiwen Zha, and 1 others. 2026. Skillos: Learning skill curation for self-evolving agents. <i>arXiv preprint arXiv:2605.06614</i> .	1065
1010		1066
1011		1067
1012		
1013		
1014	Charles Packer, Vivian Fang, Shishir G Patil, Kevin Lin, Sarah Wooders, and Joseph E Gonzalez. 2023a. Memgpt: Towards llms as operating systems. <i>arXiv preprint arXiv:2310.08560</i> .	1068
1015		1069
1016		1070
1017		1071
1018	Charles Packer and 1 others. 2023b. Memgpt: Towards llms as operating systems. <i>arXiv preprint arXiv:2310.08560</i> .	1072
1019		1073
1020		1074
1021	Joon Sung Park, Joseph C O'Brien, Carrie J Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In <i>Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology (UIST)</i> , pages 1–22.	1075
1022		1076
1023		1077
1024		1078
1025		1079
1026		1080
1027	Shishir G Patil, Tianjun Zhang, Xin Wang, and Joseph E Gonzalez. 2024. Gorilla: Large language model connected with massive apis. In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> , volume 37.	1081
1028		1082
1029		1083
1030		
1031	Hongji Pu, Xinyuan Song, and Liang Zhao. 2026. Skillops: Managing llm agent skill libraries as self-maintaining software ecosystems. <i>Preprint</i> , arXiv:2605.13716.	1084
1032		1085
1033		1086
1034		1087
1035	Denys Pushkin and Emmanuel Abbe. 2026. Lead: Breaking the no-recovery bottleneck in long-horizon reasoning. <i>Preprint</i> , arXiv:2603.06870.	1088
1036		1089
1037		1090
1038	Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen, Yusheng Su, Xin Cong, and 1 others. 2024. Chatdev: Communicative agents for software development. In <i>Proceedings of the 62nd annual meeting of the association for computational linguistics (volume 1: Long papers)</i> , pages 15174–15186.	1091
1039		1092
1040		1093
1041		1094
1042		1095
1043		1096
1044		1097
1045	Yubin Qu, Yi Liu, Tongcheng Geng, Gelei Deng, Yuekang Li, Leo Yu Zhang, Ying Zhang, and Lei Ma. 2026. Supply-chain poisoning attacks against llm coding agent skill ecosystems. <i>arXiv preprint arXiv:2604.03081</i> .	1098
1046		1099
1047		1100
1048		1101
1049		1102
1050	Christopher Rawles, Alice Li, Daniel Rodriguez, Oriana Riva, and Timothy P Lillicrap. 2023. Androidinthewild: A large-scale dataset for android device control. In <i>Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track</i> .	1103
1051		1104
1052		1105
1053		1106
1054		1107
1055		
1056	Sampriti Saha and Pranav Hemanth. 2026. Skilldex: A package manager and registry for agent skill packages with hierarchical scope-based distribution. <i>arXiv preprint arXiv:2604.16911</i> .	1108
1057		1109
1058		1110
1059		1111
1060	Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. <i>arXiv preprint arXiv:2302.04761</i> .	1112
1061		1113
1062		1114
1063		1115
1064		1116
		1117
		1118
	David Schmotz, Sahar Abdelnabi, and Maksym Andriushchenko. Agent skills enable a new class of realistic and trivially simple prompt injections.	
	David Schmotz, Sahar Abdelnabi, and Maksym Andriushchenko. 2025. Agent skills enable a new class of realistic and trivially simple prompt injections. <i>arXiv preprint arXiv:2510.26328</i> .	
	David Schmotz, Luca Beurer-Kellner, Sahar Abdelnabi, and Maksym Andriushchenko. 2026. Skill-inject: Measuring agent vulnerability to skill file attacks. <i>arXiv preprint arXiv:2602.20156</i> .	
	Shuaike Shen, Wenduo Cheng, Mingqian Ma, Al-stair Turcan, Martin Jinye Zhang, and Jian Ma. 2026. Skillfoundry: Building self-evolving agent skill libraries from heterogeneous scientific resources. <i>arXiv preprint arXiv:2604.03964</i> .	
	Haochen Shi, Xingdi Yuan, and Bang Liu. 2026. Evolving programmatic skill networks. <i>Preprint</i> , arXiv:2601.03509.	
	Tianmin Shu, Caiming Xiong, and Richard Socher. 2017. Hierarchical and interpretable skill acquisition in multi-task reinforcement learning. <i>Preprint</i> , arXiv:1712.07294.	
	Yifan Song and 1 others. 2403. Taskbench: Evaluating tool use and task planning capabilities in llm agents. <i>arXiv preprint arXiv:2403.09115</i> .	
	Weihang Su, Jianming Long, Qingyao Ai, Yichen Tang, Changyue Wang, Yiteng Tu, and Yiqun Liu. 2026. Skill retrieval augmentation for agentic ai. <i>arXiv preprint arXiv:2604.24594</i> .	
	Theodore Summers, Shunyu Yao, Karthik R Narasimhan, and Thomas L. Griffiths. 2024. Cognitive architectures for language agents. <i>Transactions on Machine Learning Research</i> . Survey Certification, Featured Certification.	
	Yiqun Sun, Pengfei Wei, and Lawrence B. Hsieh. 2026. Don't retrieve, navigate: Distilling enterprise knowledge into navigable agent skills for qa and rag. <i>Preprint</i> , arXiv:2604.14572.	
	Didac Suris, Ruoshi Gupta, and Carl Vondrick. 2023. Vipergpt: Visual grounding via python execution. In <i>IEEE International Conference on Computer Vision (ICCV)</i> .	
	Romal Thoppilan, Daniel De Freitas, and 1 others. 2022. Lamda: Language models for dialog applications. <i>arXiv preprint arXiv:2201.08239</i> .	
	Guiyao Tie, Jiawen Shi, Pan Zhou, and Lichao Sun. 2026. Badskill: Backdoor attacks on agent skills via model-in-skill poisoning. <i>Preprint</i> , arXiv:2604.09378.	
	Songjun Tu, Chengdong Xu, Qichao Zhang, Yaocheng Zhang, Xiangyuan Lan, Linjing Li, Dong Li, and Dongbin Zhao. 2026. Dynamic dual-granularity skill bank for agentic rl. <i>Preprint</i> , arXiv:2603.28716.	

1119	Georgios Tzifas and Hamidreza Kasaei. 2024. Life-long robot library learning: Bootstrapping composable and generalizable skills for embodied control with language models. In <i>2024 IEEE International Conference on Robotics and Automation (ICRA)</i> , pages 515–522. IEEE.	1173
1120		1174
1121		1175
1122		1176
1123		1177
1124		1178
1125	Yash Vishe, Rohan Surana, Xunyi Jiang, Zihan Huang, Xintong Li, Nikki Lijing Kuang, Tong Yu, Ryan A. Rossi, Jingbo Shang, Julian McAuley, and Junda Wu. 2026. <i>Skill-r1: Agent skill evolution via reinforcement learning</i> . <i>Preprint</i> , arXiv:2605.09359.	1179
1126		1180
1127		1181
1128		1182
1129		
1130	Chenxi Wang, Zhuoyun Yu, Xin Xie, Wuguannan Yao, Runnan Fang, Shuofei Qiao, Kexin Cao, Guozhou Zheng, Xiang Qi, Peng Zhang, and 1 others. 2026a. Skillx: Automatically constructing skill knowledge bases for agents. <i>arXiv preprint arXiv:2604.04804</i> .	1183
1131		1184
1132		1185
1133		1186
1134		
1135	Fali Wang, Chenglin Weng, Xianren Zhang, Siyuan Hong, Hui Liu, and Suhang Wang. 2026b. <i>Graph-skill: Documentation-guided hierarchical retrieval-augmented coding for complex graph reasoning</i> . <i>Preprint</i> , arXiv:2603.06620.	1187
1136		1188
1137		1189
1138		1190
1139		1191
1140	Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. 2023. Voyager: An open-ended embodied agent with large language models. <i>arXiv preprint arXiv:2305.16291</i> .	1192
1141		1193
1142		1194
1143		1195
1144		
1145	Hao Wang, Guozhi Wang, Han Xiao, Yufeng Zhou, Yue Pan, Jichao Wang, Ke Xu, Yafei Wen, Xiaohu Ruan, and Xiaoxin Chen. 2026c. Skill-conditioned self-distillation for multi-turn llm agents. <i>arXiv preprint arXiv:2604.10674</i> .	1196
1146		1197
1147		1198
1148		1199
1149		1200
1150	Jiayu Wang, Yifei Ming, Zixuan Ke, Shafiq Joty, Aws Albarghouthi, and Frederic Sala. 2026d. <i>Skillorchestra: Learning to route agents via skill transfer</i> . <i>Preprint</i> , arXiv:2602.19672.	1201
1151		1202
1152		1203
1153		1204
1154	Jiayu Wang, Yifei Ming, Zixuan Ke, Shafiq Joty, Aws Albarghouthi, and Frederic Sala. 2026e. Skillorchestra: Learning to route agents via skill transfer. <i>arXiv preprint arXiv:2602.19672</i> .	1205
1155		1206
1156		1207
1157		1208
1158	Jiongxiao Wang, Qiaojing Yan, Yawei Wang, Yijun Tian, Soumya Smruti Mishra, Zhichao Xu, Megha Gandhi, Panpan Xu, and Lin Lee Cheong. 2025a. Reinforcement learning for self-improving agent with skill library. <i>arXiv preprint arXiv:2512.17102</i> .	1209
1159		1210
1160		1211
1161		1212
1162		1213
1163	Qianli Wang, Boyang Ma, Minghui Xu, and Yue Zhang. 2026f. When skills lie: Hidden-comment injection in llm agents. <i>arXiv preprint arXiv:2602.10498</i> .	1214
1164		1215
1165		1216
1166	Ruoyao Wang, Peter Jansen, Marc-Alexandre Côté, and Prithviraj Ammanabrolu. 2022. <i>ScienceWorld: Is your agent smarter than a 5th grader?</i> In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 11279–11298, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	1217
1167		1218
1168		1219
1169		1220
1170		1221
1171		1222
1172		1223
		1224
		1225
		1226
		1227
	Xinyu Jessica Wang, Haoyue Bai, Yiyu Sun, Haorui Wang, Shuibai Zhang, Wenjie Hu, Mya Schroder, Bilge Mutlu, Dawn Song, and Robert D Nowak. 2026g. <i>The long-horizon task mirage? diagnosing where and why agentic systems break</i> . <i>Preprint</i> , arXiv:2604.11978.	1228
		1229
	Zimu Wang, Yuling Shi, Mengfan Li, Zijun Liu, Jie M. Zhang, Chengcheng Wan, and Xiaodong Gu. 2026h. <i>Effiskill: Agent skill based automated code efficiency optimization</i> . <i>Preprint</i> , arXiv:2603.27850.	1230
		1231
	Zora Zhiruo Wang, Apurva Gandhi, Graham Neubig, and Daniel Fried. 2025b. Inducing programmatic skills for agentic tasks. <i>arXiv preprint arXiv:2504.06821</i> .	1232
		1233
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> .	1234
		1235
	Hongbo Wen, Ying Li, Hanzhi Liu, Chaofan Shou, Yanju Chen, Yuan Tian, and Yu Feng. 2026. <i>Semia: Auditing agent skills via constraint-guided representation synthesis</i> . <i>Preprint</i> , arXiv:2605.00314.	1236
		1237
	Zikai Alex Wen. 2026. <i>Toward user comprehension supports for LLM agent skill specifications</i> . In <i>First Workshop on Agent Skills</i> .	1238
		1239
	Niklas Wretblad, Fredrik Gordh Riseby, Rahul Biswas, Amin Ahmadi, and Oskar Holmström. 2024. <i>Understanding the effects of noise in text-to-sql: An examination of the bird-bench benchmark</i> . <i>Preprint</i> , arXiv:2402.12243.	1240
		1241
	Haolin Wu, Yuecheng Liu, Junyi Dong, Heng Zhang, Sitong Mao, Hesheng Wang, Weigang Wu, and Shunbo Zhou. 2025. <i>Ascent: Autonomous skill learning toward complex embodied tasks with foundation models</i> . In <i>2025 IEEE International Conference on Robotics and Automation (ICRA)</i> , pages 16752–16758.	1242
		1243
	Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Riley Zhang, Jian Liu, and 1 others. 2023. Autogen: Enabling next-gen llm applications via multi-agent conversation. <i>arXiv preprint arXiv:2308.08155</i> .	1244
		1245
	Xiyang Wu, Zongxia Li, Guangyao Shi, Alexander Duffy, Tyler Marques, Matthew Lyle Olson, Tianyi Zhou, and Dinesh Manocha. 2026. <i>Co-evolving llm decision and skill bank agents for long-horizon tasks</i> . <i>Preprint</i> , arXiv:2604.20987.	1246
		1247
	Chunqiu Steven Xia, Yinlin Deng, Soren Dunn, and Lingming Zhang. 2025. <i>Demystifying llm-based software engineering agents</i> . <i>Proc. ACM Softw. Eng.</i> , 2(FSE).	1248
		1249
	Peng Xia, Jianwen Chen, Hanyang Wang, Jiaqi Liu, Kaide Zeng, Yu Wang, Siwei Han, Yiyang Zhou, Xujiang Zhao, Haifeng Chen, and 1 others.	1250
		1251

1228	2026a. Skillrl: Evolving agents via recursive skill-	exploration and iterative feedback. <i>Preprint</i> ,	1282
1229	augmented reinforcement learning. <i>arXiv preprint</i>	arXiv:2506.04287.	1283
1230	<i>arXiv:2602.08234</i> .		
1231	Tianle Xia, Lingxiang Hu, Yiding Sun, Ming Xu, Lan	Yutao Yang, Junsong Li, Qianjun Pan, Bihao Zhan, Yux-	1284
1232	Xu, Siying Wang, Wei Xu, and Jie Jiang. 2026b.	uan Cai, Lin Du, Jie Zhou, Kai Chen, Qin Chen,	1285
1233	Grasp: Graph-structured skill compositions for llm	Xin Li, Bo Zhang, and Liang He. 2026. Autoskill:	1286
1234	agents . <i>arXiv preprint arXiv:2604.17870</i> .	Experience-driven lifelong learning via skill self-	1287
		evolution . <i>Preprint</i> , arXiv:2603.01145.	1288
1235	Tianle Xia, Lingxiang Hu, Yiding Sun, Ming Xu, Lan	Zonghan Yang, Shengjie Wang, Kelin Fu, Wenyang He,	1289
1236	Xu, Siying Wang, Wei Xu, and Jie Jiang. 2026c.	Weimin Xiong, Yibo Liu, Yibo Miao, Bofei Gao,	1290
1237	Grasp: Graph-structured skill compositions for llm	Yejie Wang, Yingwei Ma, and 1 others. 2025b. Kimi-	1291
1238	agents . <i>Preprint</i> , arXiv:2604.17870.	dev: Agentless training as skill prior for swe-agents .	1292
		<i>arXiv preprint arXiv:2509.23045</i> .	1293
1239	Wenjie Xiao, Xuehai Tang, Biyu Zhou, Songlin Hu,	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak	1294
1240	and Jizhong Han. 2026. Routeguard: Internal-signal	Shafran, Karthik Narasimhan, and Yuan Cao. 2023.	1295
1241	detection of skill poisoning in llm agents . <i>arXiv</i>	React: Synergizing reasoning and acting in language	1296
1242	<i>preprint arXiv:2604.22888</i> .	models . In <i>International Conference on Learning</i>	1297
1243	Senwei Xie, Yuntian Zhang, Ruiping Wang, and Xilin	Representations (ICLR) .	1298
1244	Chen. 2026. Uni-skill: Building self-evolving skill		
1245	repository for generalizable robotic manipulation .	Haoran Ye, Xuning He, Vincent Arak, Haonan Dong,	1299
1246	<i>Preprint</i> , arXiv:2603.02623.	and Guojie Song. 2026. Meta context engineer-	1300
1247	Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan	ing via agentic skill evolution . <i>arXiv preprint</i>	1301
1248	Li, Siheng Zhao, Ruisheng Cao, Toh Jing Hua, Zhou-	<i>arXiv:2601.21557</i> .	1302
1249	jun Cheng, Dongchan Shin, Fangyu Lei, Yitao Liu,	Simon Yu, Gang Li, Weiyang Shi, and Peng Qi.	1303
1250	Yiheng Xu, Shuyan Zhou, Silvio Savarese, Caim-	2025. Polyskill: Learning generalizable skills	1304
1251	ing Xiong, Victor Zhong, and Tao Yu. 2024. Os-	through polymorphic abstraction . <i>arXiv preprint</i>	1305
1252	world: Benchmarking multimodal agents for open-	<i>arXiv:2510.15863</i> .	1306
1253	ended tasks in real computer environments . <i>Preprint</i> ,		
1254	arXiv:2404.07972.	Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga,	1307
1255	Hanwen Xing, Haomin Zhuang, Xuandong Zhao, Yue	Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingn-	1308
1256	Huang, Zhenheng Tang, and Xiangliang Zhang. 2026.	ing Yao, Shanelle Roman, Zilin Zhang, and Dragomir	1309
1257	Recipes for agents: Understanding skills and their	Radev. 2018. Spider: A large-scale human-labeled	1310
1258	open questions . <i>Preprint, ResearchGate</i> . doi, 10.	dataset for complex and cross-domain semantic pars-	1311
1259	Linyu Li Ganghong Huang Jianfeng Liu Honglin Qiao	ing and text-to-SQL task . In <i>Proceedings of the 2018</i>	1312
1260	Xingyan Liu, Xiyue Luo. 2026. Forging domain-	<i>Conference on Empirical Methods in Natural Lan-</i>	1313
1261	specific, self-evolving agent skills in cloud technical	<i>guage Processing</i> , pages 3911–3921, Brussels, Bel-	1314
1262	support . <i>arXiv preprint arXiv:2604.08618</i> .	gium. Association for Computational Linguistics.	1315
1263	Mengda Xu, Zhenjia Xu, Cheng Chi, Manuela Veloso,	Haoqi Yuan, Chi Zhang, Hongcheng Wang, Feiyang	1316
1264	and Shuran Song. 2023. Xskill: Cross embodiment	Xie, Penglin Cai, Hao Dong, and Zongqing Lu.	1317
1265	skill discovery . <i>Preprint</i> , arXiv:2307.09955.	2023. Skill reinforcement learning and planning	1318
1266	Yangjie Xu, Lujun Li, Lama Sleem, Niccolo Gentile,	for open-world long-horizon tasks . <i>arXiv preprint</i>	1319
1267	Yewei Song, Yiqun Wang, Siming Ji, Wenbo Wu, and	<i>arXiv:2303.16563</i> .	1320
1268	Radu State. 2026a. Agent skill framework: Perspec-	Siyu Yuan, Kaitao Song, Jiangjie Chen, Xu Tan,	1321
1269	tives on the potential of small language models in in-	Yongliang Shen, Ren Kan, Dongsheng Li, and De-	1322
1270	dustrial environments . <i>Preprint</i> , arXiv:2602.16653.	qing Yang. 2024. Easytool: Enhancing llm-based	1323
1271	Zishan Xu, Yifu Guo, Yuquan Lu, Fengyu Yang,	agents with concise tool instruction . <i>Preprint</i> ,	1324
1272	Zhiyuan Yao, Jiaye Lin, Ruyi Gong, and Lihua Cai.	arXiv:2401.06201.	1325
1273	2026b. Skillervo: An experience learning framework	Renos Zabounidis, Yue Wu, Simon Stepputtis, Woojun	1326
1274	with reinforcement learning for skill evolution .	Kim, Yuanzhi Li, Tom Mitchell, and Katia Sycara.	1327
1275	John Yang and 1 others. 2024. Swe-agent: Agent-	2026. Scalar: Learning and composing skills through	1328
1276	computer interfaces for resolving real-world soft-	llm guided symbolic planning and deep rl grounding .	1329
1277	ware engineering problems . <i>arXiv preprint</i>	<i>Preprint</i> , arXiv:2603.09036.	1330
1278	<i>arXiv:2405.01115</i> .	Guijia Zhang, Shu Yang, Xilin Gong, and Di Wang.	1331
1279	Yongjin Yang, Sinjae Kang, Juyong Lee, Dongjun	2026a. Stars: Skill-triggered audit for request-	1332
1280	Lee, Se-Young Yun, and Kimin Lee. 2025a. Auto-	conditioned invocation safety in agent systems . <i>arXiv</i>	1333
1281	ated skill discovery for language agents through	<i>preprint arXiv:2604.10286</i> .	1334

1439
1440
1441
1442
1443
1444
1445
1446
1447
1448

1449
1450
1451
1452
1453

1454
1455
1456
1457
1458
1459
1460

1461

1462
1463
1464
1465
1466
1467
1468

1469

1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482

1483
1484
1485

A.4.3 Robotics and Embodied Agents

In embodied domains, cross-embodiment discovery enables robotic systems to transfer learned skills across heterogeneous hardware boundaries natively (Xu et al., 2023). Specialized frameworks leverage agentic vision-language-action models to execute contact-rich physical tasks (Liu et al., 2026a), while self-evolving capability repositories drive generalizable robotic manipulation across open-world environments (Xie et al., 2026).

A.4.4 Open-World and Gaming Environments

In unconstrained gaming environments like Minecraft, frameworks empower agents with survival capabilities through expansive, open-world skill libraries (Liu et al., 2025; Yuan et al., 2023).

A.4.5 Collaborative and Meta-Skills

Meta-skills are advanced routines operating directly over the capability space to diagnose deficiencies and rewrite existing files to continuously self-improve (Yang et al., 2026; Zhang et al., 2026b; Zhou et al., 2026b; Yu et al., 2025; Xingyan Liu, 2026).

B Skill Ecosystems and Marketplaces

The rapid maturation of the agent skills abstraction has catalyzed the development of specialized, domain-specific ecosystems and public marketplaces. These platforms package procedural expertise to solve complex, long-horizon tasks across diverse computer environments and interactive simulators.

B.1 Desktop Operating Systems

Desktop operating systems represent a mature application environment for agentic skills, enabling models to execute open-ended administrative tasks across complex file systems and local software applications (Xie et al., 2024; Chen et al., 2026d; Liu et al., 2026c). In these setups, agents must coordinate visual observations and accessibility tree states to execute deterministic mouse clicks and keyboard typing routines. State-of-the-art computer-use agents are evaluated on real Ubuntu environments via the OSWorld benchmark, which measures task completion rates across diverse software environments.

B.2 Mobile Operating Systems

Mobile operating systems present a highly dynamic, touch-centric environment where agents

must automate complex task sequences on real smartphones (Rawles et al., 2023). To navigate these devices, mobile GUI agents translate natural language user intents into sequences of swipe, scroll, tap, and back buttons (Li et al., 2026b). These mobile automation capabilities are validated on large-scale crowdsourced datasets like AndroidInTheWild, testing the agent’s ability to operate diverse apps under dynamic layout variations (Rawles et al., 2023).

B.3 Web Navigation

Web navigation represents a massive, open-world ecosystem where agents leverage specialized DOM-parsing and HTML-filtering skills to browse websites (Deng et al., 2023). In these scenarios, generalist web agents are prompted to scroll, click, and input text across real-world, complex e-commerce platforms and search engines (He et al., 2024). Web-use capabilities are validated on robust benchmarks like WebArena, testing the model’s robustness under stateful, dynamic web page modifications (Zhou et al., 2024).

B.4 Software Engineering

Software engineering is the most demanding domain for agentic computing, requiring the coordination of file navigation, code parsing, compiler interaction, and patch generation (Jimenez et al., 2024). Rather than relying on simple, stateless code completion, software agents use packaged IDE skills to resolve real-world software issues autonomously (Yang et al., 2024). These capabilities are rigorously evaluated on repository-level codebases like SWE-bench, measuring functional correctness via execution-based test verification (Jimenez et al., 2024).

B.5 Scientific Exploration

Scientific exploration represents an emerging frontier where agents deploy specialized skills to plan experiments and reason about physical and chemical laws (Wang et al., 2022). In these settings, agents are integrated with interactive text environments and symbolic physics simulators to perform hypothesis testing and analyze experimental outcomes (Wang et al., 2022). Grounding language reasoning in interactive simulations has been shown to improve the success rate of scientific exploration agents, paving the way for autonomous laboratory assistants (Wang et al., 2022).

1486
1487
1488
1489
1490
1491
1492
1493
1494
1495

1496

1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507

1508

1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520

1521

1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533