

AI Agent Systems: Architectures, Applications, and Evaluation

BIN XU, School of Electrical, Computer and Energy Engineering, Arizona State University, USA

AI agents—systems that combine foundation models with reasoning, planning, memory, and tool use—are rapidly becoming a practical interface between natural-language intent and real-world computation. This survey synthesizes the emerging landscape of AI agent architectures for (i) deliberation and reasoning (e.g., chain-of-thought style decomposition, self-reflection and verification, and constraint-aware decision making), (ii) planning and control (from reactive policies to hierarchical and multi-step planners), and (iii) tool calling and environment interaction (retrieval, code execution, APIs, and multimodal perception). We organize prior work into a unified taxonomy spanning agent components (policy/LLM core, memory, world models, planners, tool routers, and critics), orchestration patterns (single-agent vs. multi-agent; centralized vs. decentralized coordination), and deployment settings (offline analysis vs. online interactive assistance; safety-critical vs. open-ended tasks). We discuss key design trade-offs—latency vs. accuracy, autonomy vs. controllability, and capability vs. reliability—and highlight how evaluation is complicated by non-determinism, long-horizon credit assignment, tool and environment variability, and hidden costs such as retries and context growth. Finally, we summarize measurement and benchmarking practices (task suites, human preference and utility metrics, success under constraints, robustness and security) and identify open challenges including verification and guardrails for tool actions, scalable memory and context management, interpretability of agent decisions, and reproducible evaluation under realistic workloads.

Additional Key Words and Phrases: AI Agents · Agentic AI · Agent Architectures · Agent Transformer · LLM Agents · Multimodal Agents · Vision-Language Models (VLMs) · Reasoning and Planning · Tool Use / Tool Calling · Memory and Retrieval (RAG) · Multi-Agent Systems · Safety and Alignment (RLHF/DPO) · Evaluation and Benchmarks (WebArena, ToolBench, SWE-bench, GAIA)

1 Introduction

1.1 Motivation

Foundation models have made natural language a practical interface for computation, but most real tasks are not single-turn question answering. They involve gathering information from multiple sources, maintaining state over time, choosing among tools, and executing multi-step actions under constraints (latency, permissions, safety, and cost). *AI agents* address this gap by coupling a foundation model with an execution loop that can observe an environment, plan, call tools, update memory, and verify outcomes [10, 31]. In other words, an agent is not only a generator of text; it is a controller that translates intent into *procedures* carried out in the world (software repositories, browsers, enterprise systems, or physical robots).

1.2 Background

Modern digital work is fragmented across interfaces and APIs: knowledge is distributed (documents, databases, dashboards), actions are mediated by tools (search, code execution, ticketing systems), and success is defined by end-to-end outcomes rather than plausibility. Purely conversational systems often fail in these settings due to hallucinations, lack of grounding, and inability to execute or verify actions. Tool-augmented and retrieval-augmented designs improve reliability by binding claims to evidence and by making intermediate artifacts inspectable [24, 64]. Modular tool routing (e.g., MRKL-style) further improves governance by separating language understanding from specialized tools and by enforcing structured interfaces that can be audited [21, 50].

1.3 Overview

Agents are especially important in the current age for three reasons. First, the scope of tasks is expanding from writing assistance to *workflow automation*: coding agents resolve issues end-to-end [20, 61], web agents operate real sites under variability [14, 62, 67], and enterprise assistants orchestrate multi-step operations with policy constraints. Second, deployment is increasingly *interactive and long-horizon*, where small errors compound and nondeterminism (sampling, tool failures) complicates reproducibility, motivating verification loops and trace-based evaluation [29, 44, 65]. Third, safety and security pressures are rising: prompt injection, untrusted retrieved content, and side-effecting tools require defense-in-depth alignment and guardrails beyond the final response [5].

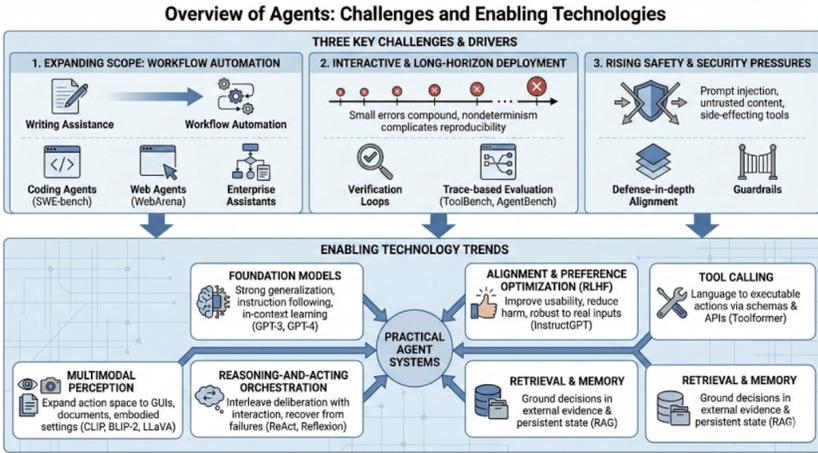


Fig. 1. Overview of AI agents and the agent execution loop (reasoning, tools, and memory)

Fig. 1 provides a high-level visual overview of the main components and execution loop of an AI agent.

Several technology trends enable practical agent systems today. **Foundation models** provide strong generalization, instruction following, and emergent in-context learning that supports rapid adaptation without retraining [9, 36]. **Alignment and preference optimization** (e.g., RLHF) improve usability and reduce harmful behavior, making agents more robust under real user inputs [11, 37]. **Tool calling** turns language into executable actions via schemas and APIs [40, 50], while **retrieval and memory** ground decisions in external evidence and persistent state [24, 38, 50]. **Reasoning-and-acting orchestration** interleaves deliberation with environment interaction to improve grounding and recover from failures [53, 64]. Finally, **multimodal perception** expands the action space to GUIs, documents, and embodied settings by grounding language in visual inputs [26, 28, 45].

1.4 Current Gaps

Despite rapid progress, agent systems remain limited by reliability, reproducibility, and governance at scale. Long-horizon tasks amplify compounding errors, and nondeterminism (sampling, tool variability) makes evaluation and debugging difficult without standardized protocols and trace completeness [29, 30, 44]. Tool-centric agents also introduce new safety and security risks: untrusted retrieved content and prompt injection can manipulate tool use, and side-effecting actions

require stronger constraints than text-only moderation [5, 21, 48]. Finally, system-level trade-offs—autonomy vs. controllability, latency vs. reliability, and capability vs. safety—are not yet well understood across domains and deployment settings [49, 66].

This survey synthesizes emerging agent architectures for reasoning, planning, tool use, and deployment. We organize the landscape along (i) *learning strategies* and system optimization (§3), and (ii) *application tasks* that stress different capabilities and evaluation regimes (§5). Throughout, we highlight recurring design trade-offs and emphasize reproducible evaluation under realistic tool and environment variability.

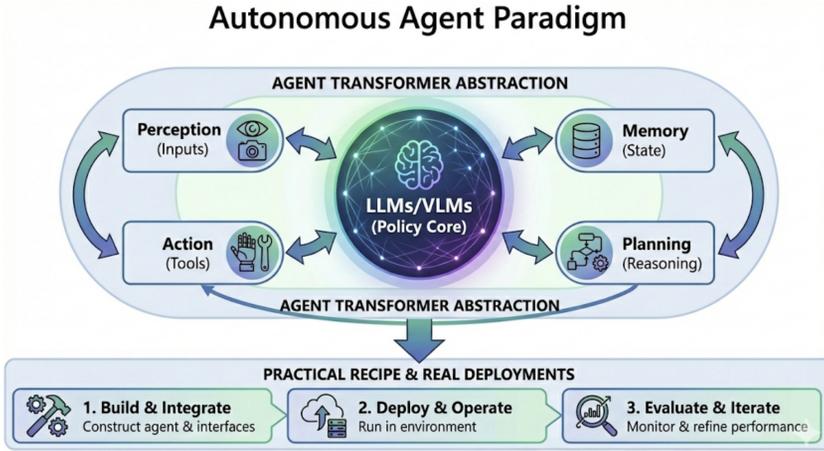


Fig. 2. Agent-centric AI paradigm: models embedded in tool- and environment-interaction loops

Fig. 2 summarizes the agent-centric paradigm that motivates the architectural and evaluation choices discussed in the remainder of the survey.

2 Autonomous Agent Paradigm

This section introduces a unifying paradigm for agent systems centered on transformer-based foundation models. We first summarize the role of LLMs/VLMs as the policy core, then define an “agent transformer” abstraction that makes agent components and interfaces explicit, and finally describe a practical recipe for building agent transformers in real deployments [33, 41, 48].

2.1 LLMs and VLMs

Large Language Models (LLMs) are the dominant *policy cores* for modern agents: they map heterogeneous context (instructions, retrieved documents, tool outputs, and internal memory) to decisions (plans, tool calls, or natural-language actions). Frontier-scale models exhibit strong instruction following and in-context learning, enabling rapid capability bootstrapping without retraining [9, 36]. However, LLMs are not inherently grounded: without external evidence and executable checks they can hallucinate plausible but incorrect statements. This motivates tool-centric and retrieval-centric agent designs where the model is an orchestrator over trusted tools and data sources [21, 24, 64].

An important recent shift is that capability gains increasingly come from *system design* rather than only from bigger backbones. Modern deployments treat the LLM as a planner/controller inside a budgeted loop: the agent is constrained by explicit limits on time, tokens, tool calls, and permissible side effects, and it dynamically allocates “thinking” (deliberation) only when the task is hard or risky.

This connects directly to test-time compute scaling: self-consistency, reranking, backtracking, and tree-style search can improve reliability without retraining, but must be used selectively to avoid runaway cost and latency [58, 63]. Relatedly, agents increasingly rely on *structured action spaces* (typed tool schemas and structured outputs) as the primary control surface: the model proposes actions that must pass schema validation and policy checks before execution, reducing the impact of free-form hallucinations and enabling stronger auditing [21, 50]. Finally, the practical frontier is shifting from “answering” to “operating”: agents are expected to maintain state, recover from tool failures, and justify actions with evidence traces, which places greater emphasis on memory design and trace completeness as first-class artifacts [29, 44, 64].

Vision-Language Models (VLMs) extend this paradigm by grounding decisions in images, screens, documents, and embodied observations. Contrastive and instruction-tuned VLMs provide a robust interface from pixels to tokens, enabling agents to operate GUIs (screenshots), read charts and forms, and align actions with visual state [26, 28, 45]. In practice, multimodal agents often decompose perception into tools (OCR, detection, layout parsing) and use an LLM as the planner/controller that integrates visual evidence with text and tool outputs [36, 64]. This division improves auditability: intermediate perceptual artifacts can be inspected and verified before committing to downstream actions.

Alignment and preference optimization are also foundational for the paradigm. RLHF-style training improves instruction following and reduces harmful behavior, making the policy core more reliable under real user inputs [11, 37]. Yet, because agents can take side-effecting actions via tools, safety must be enforced end-to-end across the entire execution graph (retrieval, tool outputs, and action gating), not only in the final response [5, 21].

2.2 Agent Transformer Definition

We define an *agent transformer* as a transformer-based policy model embedded in a structured control loop with explicit interfaces to (i) **observations** from an environment, (ii) **memory** (short-term working context and long-term state), (iii) **tools** with typed schemas, and (iv) **verifiers/critics** that check proposals before side effects occur. The key idea is to make agent behavior a sequence model over *interaction traces*: a trajectory of observations, intermediate thoughts/plans, tool invocations, and outcomes.

Concretely, an agent transformer can be described by the tuple $\mathcal{A} = (\pi_\theta, \mathcal{M}, \mathcal{T}, \mathcal{V}, \mathcal{E})$, where π_θ is the transformer policy, \mathcal{M} is a memory subsystem (e.g., retrieval, summaries, and state), \mathcal{T} is a set of tools (APIs, code execution, search, databases), \mathcal{V} is a set of verifiers/critics, and \mathcal{E} is the environment. At iteration t , the execution loop proceeds as follows: (i) the agent collects an observation o_t from the environment \mathcal{E} ; (ii) it retrieves relevant memory m_t from \mathcal{M} ; (iii) it proposes a candidate action a_t using the policy π_θ conditioned on (o_t, m_t) ; (iv) it validates a_t using \mathcal{V} (and any tool-schema constraints); and (v) it executes the selected tool call in \mathcal{T} , which updates both the environment \mathcal{E} and the memory \mathcal{M} for the next step.

$$\mathcal{A} = (\pi_\theta, \mathcal{M}, \mathcal{T}, \mathcal{V}, \mathcal{E}), \quad (1)$$

$$o_t \leftarrow \text{Obs}(\mathcal{E}_t), \quad m_t \leftarrow \text{Retrieve}(\mathcal{M}_t, o_t), \quad (2)$$

$$\tilde{a}_t \sim \pi_\theta(\cdot \mid o_t, m_t), \quad \hat{a}_t \leftarrow \text{Validate}(\mathcal{V}, \tilde{a}_t), \quad (3)$$

$$\mathcal{E}_{t+1} \leftarrow \text{Exec}(\mathcal{E}_t, \mathcal{T}, \hat{a}_t), \quad \mathcal{M}_{t+1} \leftarrow \text{Update}(\mathcal{M}_t, o_t, \hat{a}_t, \mathcal{E}_{t+1}). \quad (4)$$

Building on this operational view, the latest framing is to interpret the loop as a *risk-aware, budgeted controller*: actions differ in reversibility and potential impact, ranging from low-risk read-only queries to high-risk irreversible operations such as writes, deployments, or payments. Agent transformers therefore implement decision policies that branch on risk—low-risk actions

may run with minimal deliberation, while high-risk actions trigger additional verification, multi-step evidence gathering, or human confirmation [5, 21]. In this framing, verifiers are not optional add-ons but define the operational semantics of the agent: a ReAct-style trace is valuable not only for performance but also for governance, since it binds decisions to evidence and tool outputs, enabling post-hoc auditing and reproducible replay [64]. Likewise, search-based deliberation (Tree-of-Thoughts) and reflection (Reflexion) can be interpreted as mechanisms for allocating extra compute when uncertainty is high or when failures are detected [53, 63].

This abstraction unifies several prominent agent patterns. Retrieval-augmented generation grounds the policy in external evidence by making retrieval a first-class tool and memory operation [24]. ReAct formalizes the interleaving of reasoning and acting by alternating between deliberation tokens and tool calls, improving grounding and enabling evidence-backed traces [64]. MRKL-style systems route tasks to specialized tools, separating language understanding from deterministic components and improving governability [21]. Critic/reflection mechanisms (e.g., Reflexion) add an internal feedback channel that reduces compounding errors and supports iterative repair [53]. Search-based deliberation (Tree-of-Thoughts) treats planning as exploring a space of action candidates, trading compute for reliability [63]. Finally, multi-agent frameworks implement the same abstraction with multiple policies that communicate via messages, enabling specialization and cross-checking at the cost of coordination complexity [25, 60].

2.3 Agent Transformer Creation

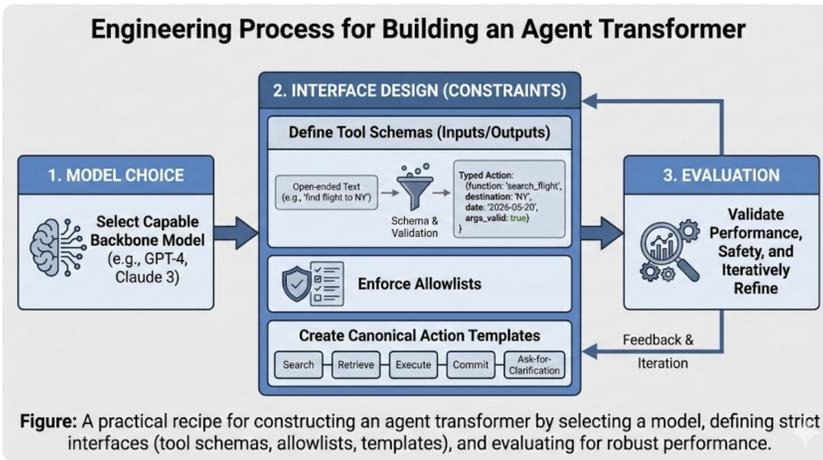


Fig. 3. Agent transformer abstraction with explicit interfaces to memory, tools, verifiers, and environment

Fig. 3 illustrates the agent transformer abstraction, emphasizing explicit interfaces to memory, tools, verifiers, and the environment.

Building an agent transformer in practice is an engineering process that combines model choice, interface design, and evaluation. A common recipe starts by selecting a capable backbone model and then constraining it through *interfaces*: define tool schemas (inputs/outputs), enforce allowlists, and create a small set of canonical action templates (search, retrieve, execute, commit, ask-for-clarification). Tool schemas reduce brittleness by turning open-ended text into typed actions, and they enable automatic argument validation before execution [21, 50].

Next, design the control loop. A minimal loop is (retrieve context) → (plan) → (act via tools) → (verify) → (update memory) → (repeat), which is closely aligned with ReAct and reflection

patterns [53, 64]. For harder tasks, add deliberation depth: tree-style search over candidate actions, self-consistency reruns, and explicit critics that check for policy violations, missing evidence, or unsafe side effects [5, 58, 63]. In tool-rich environments (web, code, enterprise systems), the most important design choice is often *when to allow side effects*: high-impact actions should require stronger verification, human confirmation, or sandboxed execution.

Then, choose learning signals that match the environment. In many deployments, supervised finetuning on traces (tool calls + outcomes) provides strong initial behavior; preference optimization and RLHF improve instruction following and refusal behavior under adversarial prompts [11, 37, 46]. Tool-use learning can be bootstrapped from synthetic traces or self-supervision (Toolformer-style), reducing the need for brittle prompt engineering [50]. For embodied or real-time control layers, combine an LLM planner with specialized controllers trained by RL/IL to satisfy timing and safety constraints [8, 15].

Recent practice emphasizes a *trace-first data flywheel*: run the agent in realistic environments, log full trajectories (prompts, tool calls, tool outputs, and outcomes), and continuously mine failures for targeted improvements (better prompts, new tools, better verifiers, or finetuning on corrected traces). This shifts learning from one-off model training to continuous system refinement, and it makes evaluation suites and regression tests essential engineering artifacts [29, 44]. Another emerging best practice is to explicitly separate *planning* from *execution*: a planner proposes a plan with explicit constraints and success criteria, while an executor carries out the plan under stricter tool permissions and validation. This separation improves controllability, supports human-in-the-loop approval for high-impact steps, and reduces the blast radius of failures [5, 21]. Finally, deployment increasingly depends on operational discipline: caching and summarization to control context growth, sandboxing for code and web actions, and policy-as-code gates for compliance. These choices do not merely improve engineering robustness; they change the effective agent policy by constraining what information and actions are available under budget and safety constraints [5, 64].

Finally, evaluate the agent transformer as a system, not a model. Benchmarks that stress realistic tool use and long-horizon execution reveal failures that do not appear in static QA: web task suites, software engineering issue resolution, and tool-use benchmarks quantify end-to-end reliability and tool correctness [20, 29, 44, 67]. Report not only success rate but also cost/latency, trace completeness, robustness under variability, and safety violations, because these determine whether an agent is deployable under real constraints [5].

3 Agent AI Learning

Agent learning spans multiple layers of the stack: (i) *learning strategies and mechanisms* (how policies, prompts, and controllers are optimized), (ii) *agent systems* (how modules and infrastructure turn a model into a dependable executor), and (iii) *agentic foundation models* (how pretraining and finetuning shape tool use, planning, and grounding). This section emphasizes how learning choices interact with long-horizon decision making, tool variability, and safety constraints [30, 48]. Fig. 4 outlines the learning stack for agentic systems, connecting mechanisms, system-level engineering, and foundation-model adaptation.

3.1 Strategy and Mechanism

3.1.1 Reinforcement Learning (RL). RL is a natural fit for agentic behavior because it directly optimizes long-horizon returns under interaction, typically formalized as a Markov decision process with a policy that maximizes expected discounted reward [43, 55]. For agents, the appeal is that RL optimizes *behavior* rather than one-step prediction: it can learn when to gather information, when to act, and how to recover from errors across multi-step trajectories.

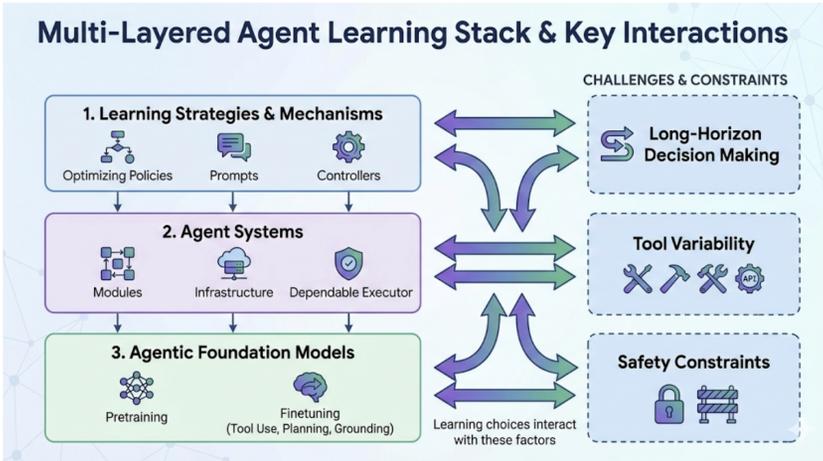
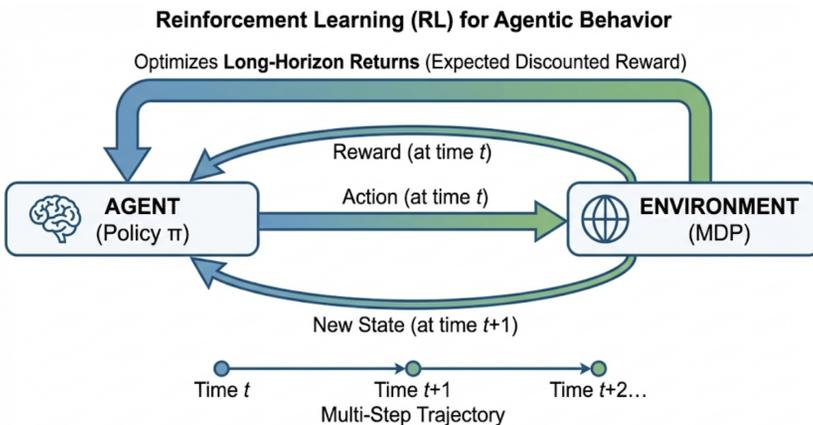


Fig. 4. Overview of agent AI learning across mechanisms, systems, and foundation models

Algorithmically, modern deep RL spans value-based learning (e.g., Q-learning variants) and policy-gradient methods, which differ in stability, sample efficiency, and how they handle continuous actions [18, 35, 51, 52]. Hierarchical RL (e.g., options) is especially relevant to agents because it provides a learning substrate for reusable skills, temporal abstraction, and planner-controller decompositions [13]. In embodied settings, RL often operates at the low-level control layer where timing constraints are strict and simulation can provide abundant interaction, while higher-level reasoning and language grounding are handled by LLMs or planners [8, 15].



Learns behavior: When to gather information, when to act, how to recover from errors.

Fig. 5. Reinforcement learning (RL) pipeline for agent policies and controllers

Fig. 5 provides a schematic view of how reinforcement learning fits into agentic decision making and control.

Empirically, classic successes in competitive games highlight the benefits of end-to-end optimization and disciplined interfaces, yielding strong execution and stability under distributional

variability induced by opponents [7, 56]. However, RL in tool-rich real-world settings faces bottlenecks that are specific to agent workflows: sparse or delayed rewards, expensive rollouts (because tool calls and environment steps are costly), and safety constraints that limit exploration. This motivates safer and more data-efficient regimes such as offline RL and constrained/safe RL, where the agent is optimized from logged trajectories and policy constraints bound undesirable actions [1, 16, 17, 23, 55].

In LLM-centric agents, RL also appears in *alignment* and *preference optimization*: RL from human feedback shapes response behavior and instruction following [11, 37], while constitution-style policy feedback and direct preference optimization provide alternative alignment mechanisms that are often easier to operationalize [5, 46]. Variants such as “verbal reinforcement” (reflection-based self-improvement) adapt the idea of learning from feedback to language-agent loops [53]. Overall, RL is most effective when (i) the environment and action space are sufficiently well-defined, (ii) interaction can be scaled (simulation or robust logging), and (iii) safety and evaluation are treated as system properties rather than post-hoc filters.

3.1.2 Imitation Learning (IL). IL provides a pragmatic route to competent behavior when expert demonstrations (human traces, scripted policies, or curated tool trajectories) are available. For agents, demonstrations are often *structured traces*: sequences of observations, intermediate rationales, tool calls, and outcomes that define not only *what* to do but also *how* to interface with tools reliably. Fig. 6 summarizes the imitation learning pipeline for acquiring agent behaviors from demonstrations and interaction traces.

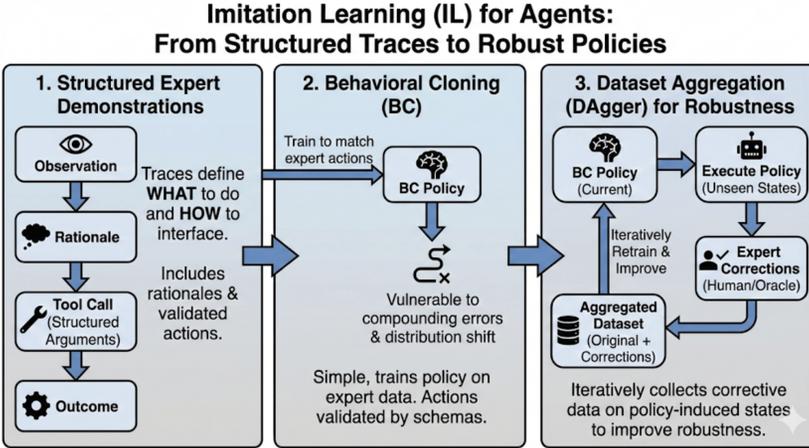


Fig. 6. Imitation learning (IL) from demonstrations and interaction traces

The simplest form, behavioral cloning, trains a policy to match expert actions, which is attractive for tool calling because actions can be represented as structured arguments and validated by schemas [42]. But pure cloning is vulnerable to compounding errors: small deviations from the expert distribution can lead to unseen states where the policy has no guidance. Dataset aggregation methods such as DAgger address this by iteratively collecting corrective demonstrations on states induced by the learned policy, improving robustness under distribution shift [47].

Beyond direct imitation, inverse RL and adversarial imitation aim to infer objectives or match expert occupancy measures. For example, GAIL learns policies by matching expert behavior distributions without explicitly specifying rewards, which can be useful when the “reward” is implicit

in expert traces (e.g., how humans navigate UIs or debug code) [19]. In agent deployments, IL is often the dominant learning signal when high-quality traces exist (enterprise workflows, customer-support playbooks, or curated code-repair sequences), because it avoids unsafe exploration and is cheaper than RL in tool-rich environments.

However, IL inherits biases and coverage gaps from the demonstration set: experts may follow organization-specific habits, may omit edge cases, and may not represent adversarial inputs (prompt injection, malicious pages, or ambiguous requests). Therefore, IL-trained behaviors often benefit from verification and repair loops (critics, self-correction, constrained execution) to handle out-of-distribution cases and to avoid blindly copying brittle heuristics [53, 64]. In practice, robust agent learning combines IL for baseline competence with system-level guardrails and post-hoc verification for reliability under long-horizon compounding errors.

3.1.3 Traditional RGB. Fig. 7 highlights traditional rule/graph/behavior-tree (RGB) components that remain important baselines and safety interfaces in agent systems. Before LLM-centric agents, many production systems relied on *Traditional RGB* components: **R**: Rule-based policies (if-then decision logic), **G**: Graph-based planners (task graphs, workflow DAGs, FSMs), and **B**: Behavior-tree-style control (hierarchical, reactive policies). These approaches remain useful because they are predictable, inspectable, and easy to govern: constraints can be encoded explicitly, and execution can be audited deterministically.

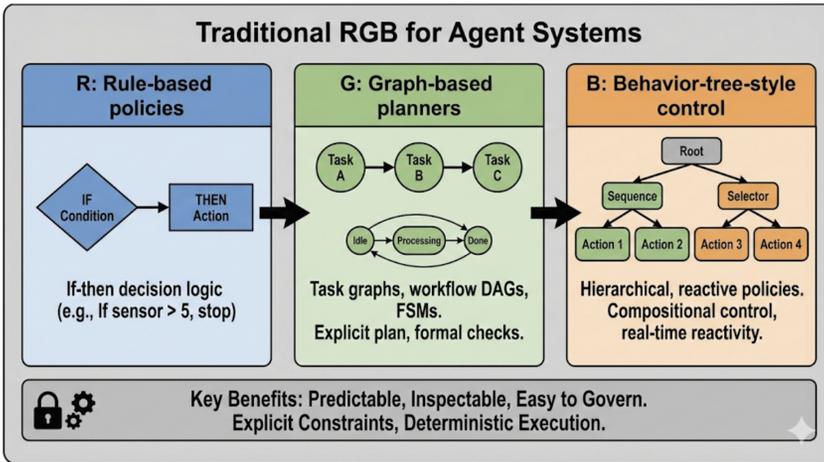


Fig. 7. Traditional RGB components: rule-based policies, graph planners, and behavior trees

Behavior trees in particular provide compositional control and real-time reactivity, making them attractive for robotics and games where strict timing and safety envelopes must be enforced [12]. Graph-based planning subsumes workflow and task-graph execution: the “plan” is explicit (nodes/edges), enabling formal checks (reachability, preconditions, approval gates) and deterministic replay, which is essential for compliance and debugging.

The main limitation of RGB approaches is brittleness: handcrafted rules may not generalize, and enumerating all edge cases becomes infeasible as environments become open-ended. This is where LLM-centric components provide value: they can interpret ambiguous instructions, propose candidate plans, and adapt to novel inputs. Modern stacks therefore hybridize: LLMs propose goals, explanations, or candidate actions, while RGB components enforce safety, timing, and domain

rules (e.g., cooldowns in games, approvals in enterprise workflows, or safety envelopes in robotics) [2, 21].

From a learning perspective, RGB components also function as *inductive biases* and *interfaces*: they constrain the action space and reduce the burden on statistical learning, often improving reliability and making evaluation more reproducible. In many production deployments, the highest-value learning work is not replacing RGB entirely, but learning *better routing and parameterization* of RGB modules (which rule to apply, which graph branch to follow, which behavior-tree subtree to activate) under contextual signals from LLMs and tools [64].

3.1.4 In-context Learning. In-context learning enables rapid task adaptation via prompting and exemplars without parameter updates. Fig. 8 depicts in-context learning as a practical mechanism for teaching agent protocols (formats, schemas, and interaction patterns) without parameter updates. For agents, this includes demonstrations of tool schemas, planning formats, and interaction protocols (e.g., how to interleave reasoning and actions). Empirically, large language models exhibit strong in-context learning capabilities where few-shot exemplars induce task behavior without finetuning [9, 34]. For agent settings, in-context learning acts as “soft programming”: prompts define action formats, tool-call schemas, and policies (what is allowed, what requires confirmation), enabling fast iteration without retraining.

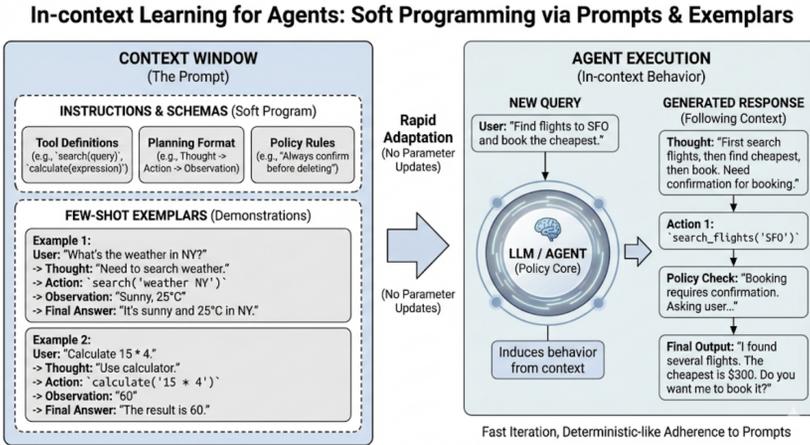


Fig. 8. In-context learning for agents via prompts, exemplars, and action schemas

Reasoning-augmented prompting is a key enabler. Chain-of-thought prompting improves multi-step reasoning and decomposition, which directly translates to better planning and tool selection in agents [22, 58, 59]. ReAct-style prompting operationalizes in-context learning by binding reasoning to tool use, improving grounding and making intermediate decisions inspectable and auditable [64]. Self-consistency and related sampling-based methods further stabilize in-context behaviors by aggregating multiple reasoning paths, which is especially useful under nondeterminism and partial observability [58].

However, in-context learning has well-known system-level failure modes for agents: context growth increases cost/latency, long prompts can dilute critical constraints, and retrieved text can introduce prompt-injection instructions that override policy. Therefore, in-context learning is most effective when paired with memory (summaries, persistent state), retrieval grounded in trusted sources, and strict tool interfaces that cannot be bypassed by text alone [5, 24]. In production,

prompts are often treated as versioned artifacts with evaluation suites and regression checks, because small changes in examples or formatting can shift behavior dramatically.

Finally, in-context learning interacts strongly with tool ecosystems: exposing a tool via a schema and exemplars is a kind of “few-shot tool learning,” which can substitute for explicit finetuning when new tools are added frequently [21, 50]. This emphasizes a design principle: if tool interfaces are stable and well-specified, in-context learning can deliver rapid capability gains with relatively predictable governance and rollback.

3.1.5 Optimization in the Agent System. Agent performance is a systems optimization problem as much as a modeling problem: reliability, latency, and cost are shaped by orchestration policies (how many calls, which tools, how much verification, and when to backtrack). Search-based planning (e.g., exploring alternative action sequences) improves hard tasks where a single rollout is unreliable, but increases compute and requires careful termination criteria and evaluation functions [63]. In LLM agents, search often operates over *action proposals* rather than over low-level states: the agent samples candidate tool calls or subplans, scores them (via critics, heuristics, or tool-based checks), and commits to the best candidate.

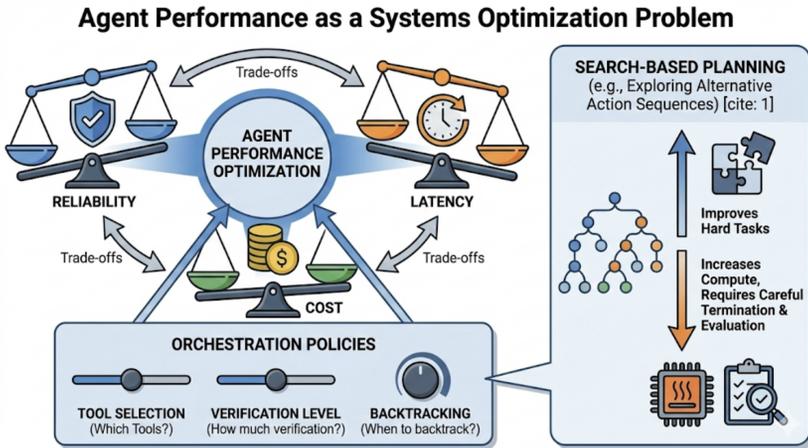


Fig. 9. Optimization Problem for AI Agent

Fig. 9 provides an optimization-oriented view of agent design, highlighting the core trade-offs between reliability, latency, and cost under constrained orchestration policies.

Verification and reflection loops trade additional computation for lower failure rates by checking actions and revising plans in response to tool outputs or detected inconsistencies [53, 64]. This is closely related to “test-time compute” scaling: rather than increasing model size, systems increase the number of controlled deliberation steps (reranking, self-consistency, backtracking). In practice, verification is most valuable for high-impact or irreversible actions (writing to a database, merging code, sending messages), where the cost of a mistake dominates latency concerns [5, 21].

Optimization also includes purely systems concerns: caching frequent retrievals, batching model calls, using smaller models for routing/moderation, and compressing memory via summaries to control context growth. These choices change the agent’s effective policy by altering what context is available and which tools are feasible under latency/cost constraints. Therefore, optimization must be evaluated end-to-end (success rate, safety, and cost/latency) rather than at the component level [29, 44].

Practical deployments adopt adaptive optimization: fast-path execution for routine cases, slower verified paths for high-risk actions, explicit budgets (time, tokens, tool calls), and permission gates that bound side effects even when the model is capable [5, 21]. This perspective treats agent learning as a *co-design* problem between models and orchestration, where the most impactful improvements often come from changing the decision loop rather than changing the base model.

3.2 Agent Systems

3.2.1 Agent Modules. At the system level, “learning” includes how the agent is modularized into components with clear contracts. Common modules include an LLM policy core, retrieval/memory, planners, tool routers, and critics/verifiers. MRKL-style routing separates language understanding from specialized tools, improving controllability and allowing toolchains to evolve without retraining the core model [21]. ReAct-style execution couples reasoning with actions to produce traceable trajectories, while reflection mechanisms provide a structured way to recover from errors and reduce compounding failures [53, 64].

From a learning viewpoint, modularization changes what must be learned by the core model. For example, if retrieval is delegated to a tool, the model learns *query formulation* rather than memorizing facts; if a verifier exists, the model can learn to propose candidates and rely on checks to reject unsafe actions. This enables a division of labor where smaller or specialized models can handle routing, safety classification, or summarization, reducing cost and improving predictability [21, 36].

Memory modules are particularly central for long-horizon agents: episodic memory (what happened), semantic memory (facts), and procedural memory (skills) support coherence beyond raw context windows. Retrieval-augmented generation provides a standard mechanism for grounding and reduces hallucinations by binding claims to retrieved evidence [24]. However, memory also increases attack surface (retrieved prompt injection) and introduces consistency challenges when stored state conflicts with new observations; critics/verifiers mitigate this by checking claims against tool outputs and trusted sources [5, 64]. Experience-oriented agent architectures further motivate richer long-term memory (personas, episodic summaries, skills) to support coherence and learning over extended interaction [39, 57].

Multi-agent variants decompose tasks across roles (planner, executor, reviewer), enabling cross-checking and specialization, but they introduce coordination costs (latency, token usage) and can amplify inconsistency when agents disagree [25, 60]. In practice, role separation works best when roles have distinct tool access and explicit handoff artifacts (plans, checklists, traces), so that disagreements can be resolved via evidence rather than free-form debate.

3.2.2 Agent Infrastructure. Infrastructure determines whether an agent can operate safely and reproducibly in real environments. Key elements include sandboxed tool execution, schema validation, identity/permission enforcement, audit logs, caching, and observability (traces of prompts, tool calls, and intermediate states). Fig. 10 sketches the infrastructure layer (sandboxing, schemas, permissions, logging) required for safe and reproducible agent deployment.

For web and GUI agents, infrastructure must handle partial observability and dynamic interfaces; standardized environments and benchmarks improve comparability and reveal brittleness undrealistic variability [67]. For coding agents and enterprise workflows, infrastructure often includes repository search, test orchestration, and policy gates; end-to-end benchmarks highlight failures that only appear when tools, environments, and long-horizon dependencies interact [20, 29, 44]. Across domains, guardrails must be enforced end-to-end (including tool outputs and retrieved text) to mitigate prompt injection and unsafe side effects [5].

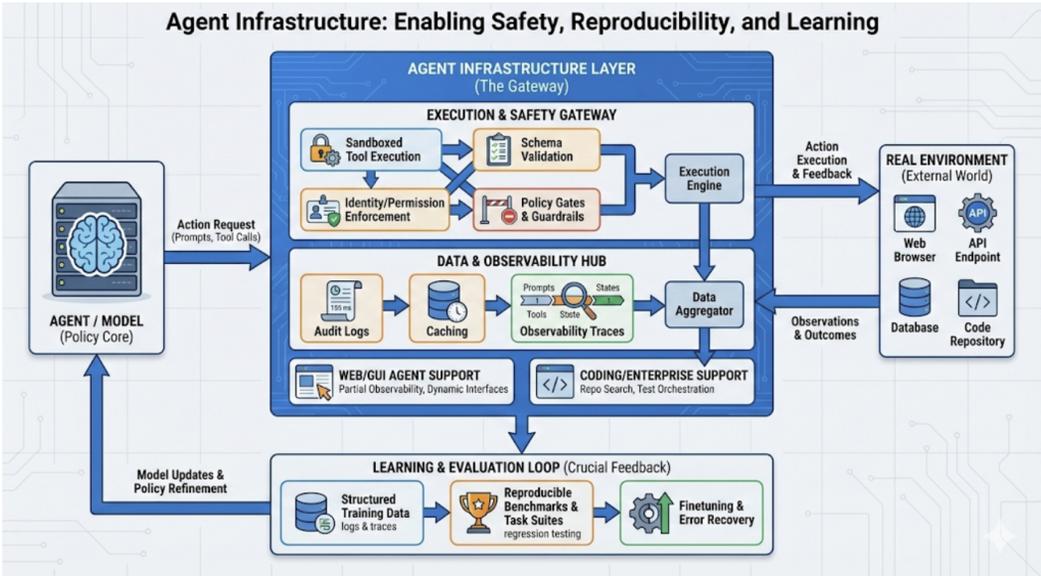


Fig. 10. Agent infrastructure for safe deployment: sandboxing, schemas, permissions, and logging

Crucially, infrastructure determines *what is learnable*. If tool calls are logged with arguments and outcomes, they become training data for tool-use finetuning and error recovery; if traces are missing or privacy-restricted, learning must rely on weaker signals. Similarly, schema validation and sandboxing turn an open-ended “action” into a constrained interface, improving reliability and reducing catastrophic failures when models hallucinate tool arguments [21, 50].

Evaluation infrastructure is part of learning: reproducible benchmarks and task suites enable regression testing and ablations over prompts, routing, memory policies, and verification depth. For interactive environments, standardized suites make it possible to report robustness to variability (layout changes, tool failures) and to quantify tool-use correctness beyond superficial answer quality [29, 44, 67]. In enterprise and safety-critical deployments, auditability (who did what, when, with what evidence) is a core requirement; therefore, learning objectives must include not only task success but also trace completeness and policy compliance [5, 21].

3.3 Agentic Foundation Models (pretraining and finetune level)

Foundation models shape agent capability through both representation learning and alignment. Pretraining builds broad world knowledge and multimodal grounding, while finetuning (instruction/policy tuning, tool-use tuning) shapes how that knowledge is accessed and acted upon. Fig. 11 summarizes how pretraining and finetuning choices shape tool use, planning, and grounding in agentic foundation models.

Pretraining: Multimodal pretraining improves grounding and perception, enabling agents to bind language to visual evidence and UI state [26, 28, 45]. More broadly, large-scale pretraining induces in-context learning and compositional generalization, which are prerequisites for tool use and long-horizon planning in open-ended environments [9, 36]. Tool-related pretraining objectives can explicitly encourage models to emit API-compatible calls and to treat tools as part of the computation graph rather than as post-hoc decorations [50]. For agentic behavior, the key pretraining question is

Infographic figure: Agentic Foundation Models: Shaping Capability Through Pretraining and Finetuning
(Section \ref(subsection:agentic_foundation_models))

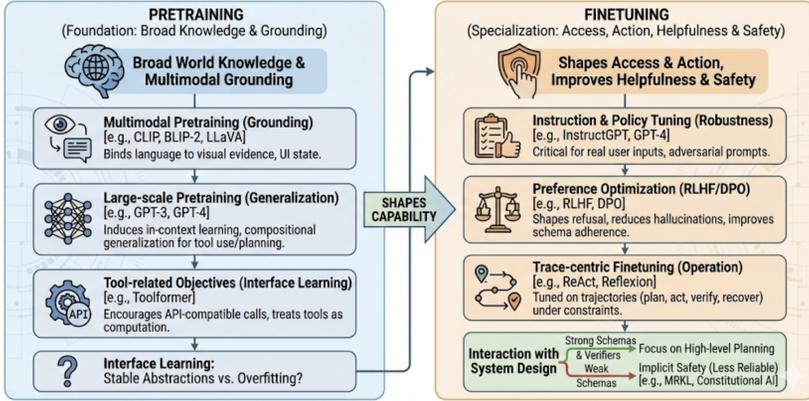


Fig. 11. Agentic foundation models: pretraining and finetuning for tool use and planning

interface learning: does the model learn stable abstractions for actions, observations, and constraints, or does it overfit to surface patterns that break under distribution shift?

Finetuning: Instruction and policy tuning improve helpfulness and safety, and they are critical for making agents robust under real user inputs and adversarial prompts [5, 36, 37]. Many systems combine supervised instruction tuning with preference-based optimization (RLHF or direct preference optimization) to shape refusal behavior, reduce hallucinations, and improve adherence to tool schemas [11, 37, 46]. Agentic finetuning is increasingly trace-centric: models are tuned on trajectories that include tool calls, intermediate checks, and corrected failures, so that the model learns not only to answer but also to *operate*—plan, act, verify, and recover—under constraints [53, 64]. Finally, finetuning interacts with system design: if orchestration enforces strict schemas and verifiers, finetuning can focus on high-level planning and query formulation; if schemas are weak, finetuning must implicitly learn safety and interface constraints, which is less reliable and harder to audit [5, 21].

4 Agent AI Taxonomy

To connect architectural choices to deployment requirements, we categorize agent systems by the *dominant locus of interaction* (text/tools, physical embodiment, simulated environments), the *generative target* (content/worlds/experiences), and the *reasoning substrate* (knowledge, logic, emotion, neuro-symbolic structure). This taxonomy is intended to be pragmatic: categories reflect the constraints that most strongly shape system design (observability, safety, latency, verification, and evaluation) [4, 49, 66].

4.1 Generalist Agent Areas

Generalist agents aim to solve heterogeneous tasks across domains (coding, browsing, analytics, and enterprise workflows) using a shared policy core plus modular tools and memory [9, 36]. **Challenges.** The dominant failure mode is long-horizon compounding error under tool and environment variability: the agent must retrieve the right context, choose correct tools and arguments, and recover from partial failures (timeouts, flaky tests, UI changes) [29, 44, 64]. Safety is also harder because untrusted inputs can induce prompt injection that targets tool usage, and side-effecting

actions have real operational cost [5, 21]. Finally, evaluation must be end-to-end (did the workflow complete correctly), not “did the text sound plausible” [20, 29, 44, 67].

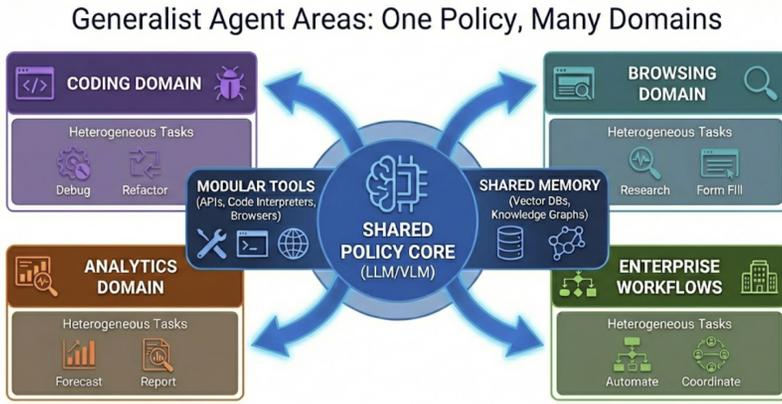


Fig. 12. Generalist agent application areas and representative capability demands

Fig. 12 summarizes representative generalist agent areas that stress broad tool use and long-horizon reliability.

Benefits. When reliable, generalist agents reduce coordination and context switching by translating intent into full workflows: find information, execute actions, and verify results [64]. They amortize engineering cost because the same agent core can be reused across tools and departments, with governance enforced through interfaces rather than retraining [21, 50]. **How to build.** A practical recipe combines (i) retrieval grounding (RAG) for evidence-backed decisions [24], (ii) modular tool routing (MRKL-style) to delegate specialized work and enforce schemas/allowlists [21, 50], (iii) ReAct-style reasoning-and-acting loops for traceable trajectories [64], and (iv) critic/reflection or search to allocate extra deliberation when uncertainty is high [53, 58, 63]. Multi-agent frameworks can further separate roles (planner/executor/reviewer) for cross-checking [25, 60]. Evaluate with realistic suites such as WebArena, SWE-bench, ToolBench, and AgentBench to expose tool-use brittleness and reproducibility gaps [20, 29, 44, 67].

4.2 Embodied Agents

Embodied agents operate in the physical world (robots, smart devices) where actions have real costs and safety is paramount [2, 8, 15]. **Challenges.** Embodiment adds partial observability, sensor noise, and continuous control; small perception errors can cascade into unsafe actions [8, 15]. Real-time constraints make naive “LLM-in-the-loop” control brittle, and sim-to-real gaps undermine plans that look feasible in simulation [2]. Evaluation must separate perception vs. planning failures and account for human interventions and recovery behavior [2, 15]. **Benefits.** Embodied agents can turn natural-language intent into physical assistance: household manipulation, warehouse operations, and assistive robotics [2, 8]. They also provide a concrete testbed for long-horizon planning, grounding, and safe tool/action execution under constraints [15]. **How to build.** The dominant design is hierarchical: an LLM/VLM planner produces high-level skills, while classical or RL controllers execute skills under safety envelopes [2, 8, 15]. Treat perception and planning as tools (mapping, grasp/motion planning, simulation rollouts) and require state assertions/verification before committing actions, mirroring ReAct-style interleaving of reasoning and tool calls [53, 64].

For high-risk actions, use conservative alignment and policy gates that enforce constraints beyond the final response [5].

4.2.1 Action Agents. Action agents emphasize *doing* over *dialogue*: they map goals to sequences of physical actions (navigate, pick, place, open) under constraints [2, 15]. **Challenges.** The core difficulty is safe skill execution under uncertainty: the agent must choose feasible actions, handle stochastic outcomes, and avoid unsafe exploration [55]. Tool interfaces (motion planners, grasp checkers) may fail or return partial information, and state estimation errors can cause compounding failures [15]. **Benefits.** Action agents enable automation in structured tasks with measurable outcomes (time-to-completion, safety violations), and they can exploit simulation and logged trajectories for learning and improvement [13, 55]. **How to build.** Use planner–controller stacks such as SayCan/PaLM-E/RT-2 style hierarchies: the language model proposes skill sequences and the controller executes with safety constraints [2, 8, 15]. Add feasibility tools (grasp, motion, collision checks) and verification loops that re-plan after deviations, analogous to reflection and search in language agents [53, 58, 63].

4.2.2 Interactive Agents. Interactive embodied agents emphasize *human-in-the-loop* operation: clarification, instruction following, and shared autonomy [8, 15, 36].

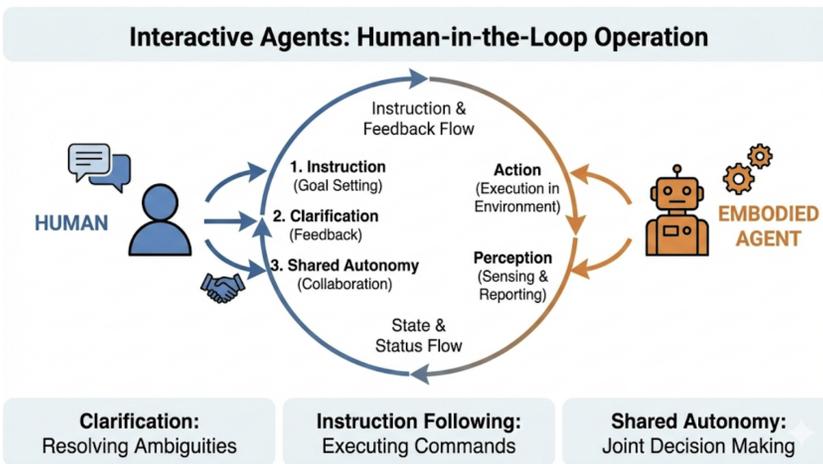


Fig. 13. Interactive embodied agents with human-in-the-loop feedback and shared autonomy

Fig. 13 provides a high-level view of interactive agents and the feedback loops needed for safe shared autonomy.

Challenges. The agent must interpret ambiguous instructions, track evolving context, and remain grounded in current sensory state [8, 15]. Miscommunication is costly: wrong objectives can cause physical harm, so the agent must ask clarifying questions and surface uncertainty rather than hallucinating intent [5, 37]. **Benefits.** Interactive agents improve usability and trust: users can steer behavior, correct mistakes, and approve high-impact actions [36]. This also produces valuable trace data for continual improvement (what was asked, what was clarified, what failed), supporting trace-centric refinement loops [64]. **How to build.** Combine multimodal grounding (VLM perception + structured perception tools) with an orchestrator that chooses between actions (act, ask, verify, escalate) using ReAct-style loops [26, 28, 64]. Enforce safety via permission gating and conservative alignment policies, and use reflection/verifiers to check plans against constraints before execution [5, 53].

4.3 Simulation and Environments Agents

Simulation/environment agents act in virtual environments (games, web sandboxes, synthetic worlds) where interaction is cheaper and more scalable than in the real world [57, 67].

Challenges. Even in simulation, partial observability and environment drift create brittleness (website layout changes, stochastic game dynamics) [67]. Simulated rewards can be mis-specified, leading to reward hacking or overfitting to benchmark artifacts, and sim-to-real transfer remains hard when moving from controlled environments to production systems [55]. **Benefits.** Simulation enables fast iteration: large-scale interaction for RL/IL, controlled ablations, and reproducible benchmarking. It also supports safe exploration when real-world side effects are unacceptable [55]. **How to build.** Treat the environment as an explicit observation/action interface and use ReAct-style interleaving to bind decisions to executed actions and state checks [64]. For web/GUI tasks, evaluate and iterate with standardized suites (WebArena) and report robustness under variability [67]. For open-world skill learning, use iterative interaction with memory and tool-use to build skill libraries (Voyager-style) and add verification/critic loops to reduce drift [53, 57, 63]. For text-based embodied simulators, task suites can serve as controlled training/evaluation harnesses [54].

4.4 Generative Agents

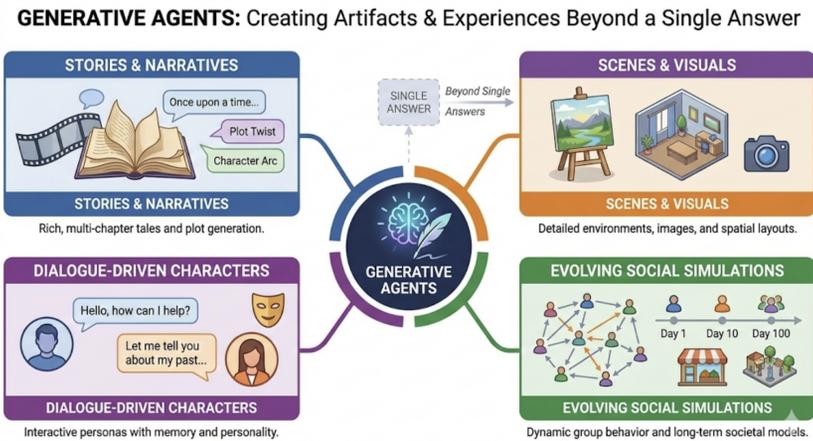


Fig. 14. Generative agents for long-horizon content and social simulation with persistent state

Generative agents produce artifacts or experiences beyond a single answer by creating stories, scenes, dialogue-driven characters, or evolving simulations of social behavior [39]. Their central difficulty is long-horizon coherence under constraints: maintaining a consistent world state, persona, and narrative while avoiding repetition and drift [39]; governance further complicates deployment (e.g., safety, provenance, and IP review), and purely free-form generation is difficult to validate without tools and verifiers [5]. Fig. 14 depicts generative agents and why persistence, retrieval, and verification are central to long-horizon coherence. In practice, these systems accelerate content creation and enable interactive experiences (e.g., NPCs and simulated societies) that adapt to users and context [39], motivating memory-centric designs where persistent memory (episodic summaries, persona state) and retrieval grounding materially improve quality and controllability [24, 39]. A common implementation pattern is to convert open-ended creation into an iterative, tool-driven loop—generate candidates, validate constraints, and revise via critics or search/reflection

mechanisms—while coupling generation to structured tool calls (world queries, rule checks, asset validators) and treating the resulting trace as an auditable artifact [21, 53, 58, 63, 64].

4.4.1 AR/VR/mixed-reality Agents. AR/VR/mixed-reality agents are a special case of generative and interactive systems: they must fuse real-time multimodal perception with low-latency action and spatial grounding [26, 28]. **Challenges.** Latency budgets are tight and sensory streams are noisy; grounding errors can cause confusing or unsafe user experiences. Privacy constraints often limit logging of multimodal data, making debugging and evaluation harder.

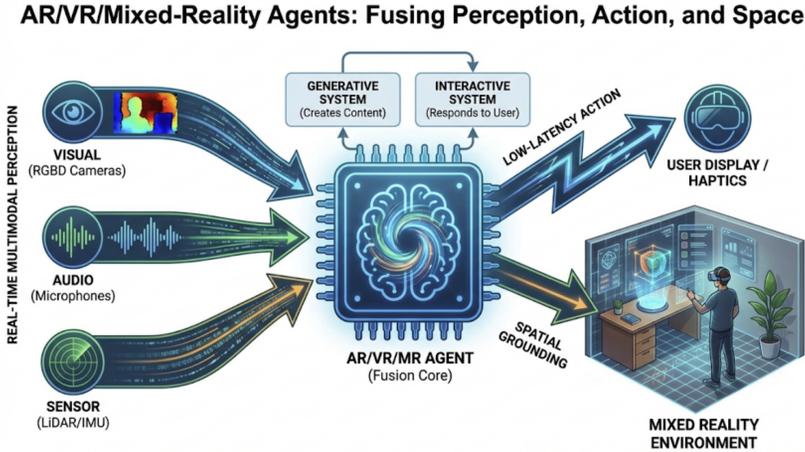


Fig. 15. AR/VR and mixed-reality agents: low-latency multimodal grounding and action

Fig. 15 highlights AR/VR agent requirements such as low-latency multimodal grounding and spatially anchored actions. **Benefits.** AR/VR agents can provide situated assistance (navigation, instruction, collaboration) by tying language to spatial context and live perception, enabling more actionable help than text-only assistants [36]. **How to build.** Decompose perception into tools (tracking, OCR, detection) and use an LLM/VLM orchestrator that interleaves reasoning with tool calls and state assertions (ReAct-style) [28, 45, 64]. Use conservative policies (ask/abstain) when uncertainty is high and enforce end-to-end guardrails against prompt injection and unsafe actions [5, 21].

4.5 Knowledge and Logical Inference Agents

Knowledge- and logic-centric agents prioritize correctness under explicit constraints by retrieving evidence, applying rules, performing multi-step inference, and producing auditable traces of their decisions (e.g., queries executed, rules applied, constraints checked) [21, 24, 64]. To be reliable, they must distinguish evidence from speculation, avoid hidden assumptions, handle schema drift as knowledge bases and tools evolve, and resist prompt injection embedded in retrieved content [5]. In logic-heavy settings, this typically requires a clear separation between the model’s proposed reasoning steps and the system’s validated outcomes, with retrieval treated as a first-class tool (RAG), claims bound to executed tool outputs, subproblems routed to symbolic or deterministic components in an MRKL-style design, and typed schemas used to reduce hallucinated actions [21, 24, 50, 64]. Reliability can be further improved by incorporating verifiers/critics and reflection-style loops that rerun checks under alternative assumptions to reduce brittle reasoning while maintaining provenance and compliance for regulated workflows [5, 53, 58, 63].

4.5.1 Knowledge Agent. Knowledge agents emphasize grounded answers backed by external sources (documents, KBs, logs) [24]. **Challenges.** Retrieval quality dominates: missing the right document can make the agent confidently wrong, and untrusted documents can inject adversarial instructions. Maintaining provenance across multi-hop retrieval and summarization is also difficult. **Benefits.** When implemented well, knowledge agents reduce hallucinations and make answers auditable via citations and traceable evidence [24, 64]. They are a natural fit for read-only enterprise assistants and compliance-oriented workflows [5]. **How to build.** Implement retrieval pipelines with metadata (source, timestamp, trust level), and force the agent to cite retrieved evidence. Use ReAct-style traces so intermediate retrieval and checks are explicit [64]. Where possible, route structured queries to deterministic tools (MRKL) and validate outputs before summarization [21].

4.5.2 Logic Agents. Logic agents emphasize constraint satisfaction and structured reasoning (e.g., planning under rules, verification of properties, or solver-backed decision making). **Challenges.** The model’s free-form reasoning is not a proof: correctness requires explicit constraints, tool semantics, and validation. Tool errors or mismatched assumptions can silently invalidate conclusions. **Benefits.** Logic agents can provide stronger guarantees by delegating correctness-critical steps to solvers and verifiers, producing outputs that are reproducible and easier to audit than purely neural reasoning [21]. **How to build.** Use the LLM for decomposition and for generating candidate symbolic forms, then call deterministic planners/solvers via typed tool schemas (MRKL/Toolformer-style tool interfaces) [21, 50]. For harder problems, use search-based deliberation over candidate plans (Tree-of-Thoughts) and rerun checks for robustness [58, 63].

4.5.3 Agents for Emotional Reasoning. Emotional reasoning agents model affect and social context to support more natural interaction (e.g., empathetic assistants, believable NPCs) [39]. **Challenges.** Emotional settings amplify safety risks: manipulation, over-trust, and boundary violations. Maintaining consistent persona and state over long interactions is hard, and adversarial prompts can trigger out-of-character or unsafe behavior [5, 37]. **Benefits.** When aligned, emotional reasoning improves user experience (trust, engagement, adherence), and can make agents more effective in support, education, and interactive entertainment by adapting communication style to context [37]. Fig. 16 illustrates emotional reasoning agents and the associated safety, persona, and consistency challenges.

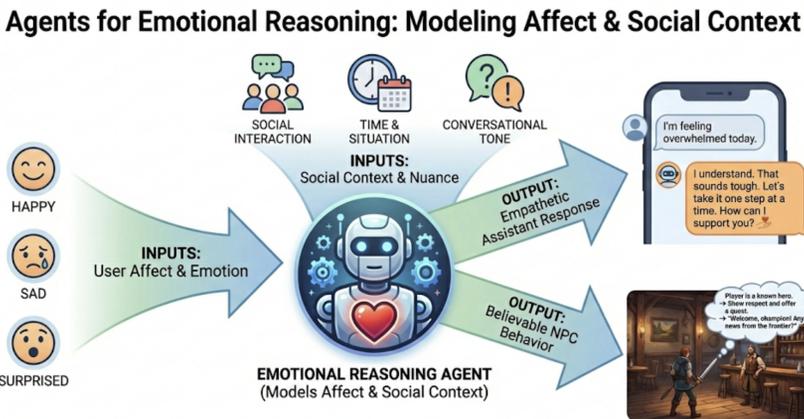


Fig. 16. Agents for emotional and social reasoning with persona consistency and safety constraints

How to build. Use aligned instruction-following backbones with explicit safety policies (Constitutional-style constraints) [5, 11, 37]. Ground interaction in structured memory (relationships, commitments, boundaries) and incorporate verifiers/critics that check for policy violations and contradictions before responding [53, 58, 64]. Evaluate with both objective consistency metrics and user studies, reporting trade-offs between safety filtering and perceived naturalness [5].

4.5.4 Neuro-Symbolic Agents. Neuro-symbolic agents combine neural models with symbolic structures (rules, graphs, typed schemas) to achieve better controllability and verifiability [21]. **Challenges.** The integration boundary is the hard part: symbolic tools have strict semantics, while neural models are probabilistic and may produce invalid calls. Keeping the system robust under tool evolution (schema drift) and untrusted inputs requires explicit validation and policy enforcement [5]. Fig. 17 provides a schematic of neuro-symbolic agent designs that combine neural policies with symbolic tools and verifiers.

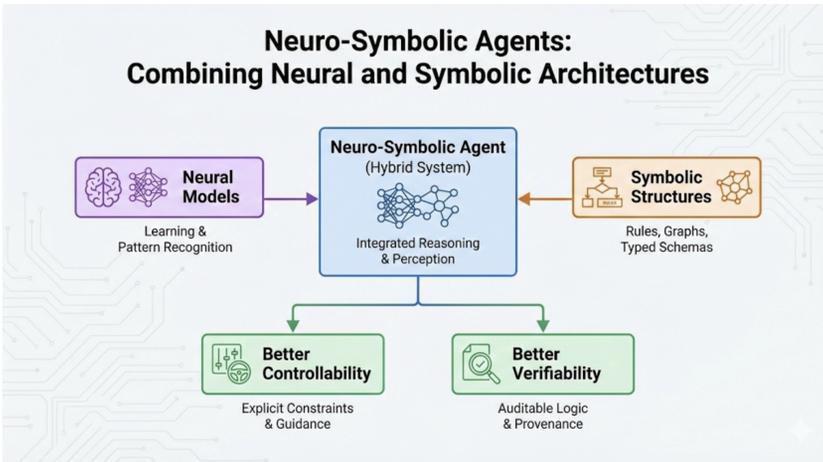


Fig. 17. Neuro-symbolic agents combining neural policies with symbolic tools and verifiers

Benefits. Neuro-symbolic designs offer a practical path to trustworthy autonomy: they preserve flexibility in language understanding while delegating correctness and side effects to deterministic components, improving audibility and governance [21, 50]. **How to build.** Implement modular routing (MRKL) so the LLM selects among specialized tools, and use Toolformer-style schema discipline to reduce invalid tool invocations [21, 50]. Add verifier gates and reflection loops so the system checks planned actions against constraints before execution and recovers from failures without expanding tool privileges [5, 53, 64].

4.6 LLMs and VLMs Agent

This category highlights agents whose *primary capability driver* is the underlying LLM/VLM backbone (often frontier-scale) and multimodal grounding [28, 36]. **Challenges.** Pure backbone scaling does not eliminate tool-use brittleness: hallucinated actions, weak grounding, and unpredictable behavior under long horizons persist without structured interfaces and verification [21, 50, 64]. Multimodal inputs add additional error modes (OCR/layout failures, visual hallucination), and cost/latency can grow quickly when deliberation and tool calls are unconstrained [26, 28, 45].

Benefits. Strong backbones provide breadth and better generalization: they unify language, vision, and tool use, enabling a single orchestrator to handle diverse tasks and modalities with less

hand-engineering [28, 36]. **How to build.** Treat the backbone as an orchestrator inside a constrained system: use RAG for grounding [24], typed tool schemas and allowlists for safe execution [21, 50], and ReAct/verification loops to bind decisions to tool outputs and recover from failures [53, 58, 64]. Evaluate under realistic variability using standardized suites and report cost/latency and safety as first-class metrics [5, 20, 29, 44, 67].

5 Agent AI Application Tasks

AI agents are increasingly deployed as *workflow executors* rather than static chat interfaces: they translate user intent into multi-step actions across tools, data sources, and environments. Following the task-oriented framework in Fig. 18, we organize applications by domains that stress different capabilities (interaction, perception, planning, tool use, and long-horizon control). For each task category, we summarize typical models, agent technologies, current challenges, and why particular agent designs are effective.

Agent AI Application Tasks: Domains, Capabilities, and Challenges

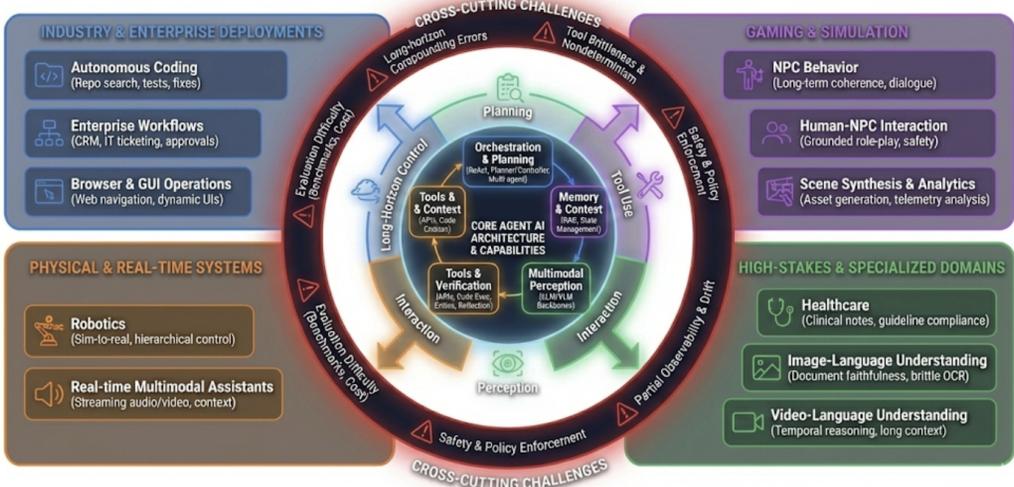


Fig. 18. Agent application landscape: task categories and capability requirements

5.1 Recent Industry Deployments and Case Studies

5.1.1 Autonomous coding and software maintenance agents.

Challenges. Real-world software tasks are long-horizon and tool-rich: agents must search large codebases, modify multiple files coherently, run tests, and handle flaky environments and dependency drift. The work is rarely “just code generation”; it requires understanding implicit requirements, existing conventions, dependency constraints, and cross-module coupling that is not documented. Failures can be subtle and delayed (security regressions, behavioral changes behind feature flags, performance cliffs, or compatibility breaks), and they often only surface under integration tests or production-like workloads. Tooling itself is a moving target: compilers, linters, dependency resolvers, and CI environments evolve, and agents must cope with non-determinism (networked dependencies, flaky tests) and partial observability (incomplete logs, truncated error

output). Evaluation is therefore difficult: plausible patches are not enough; we need end-to-end correctness, minimal regressions, and evidence that the change matches the intended specification under realistic constraints.

Solution. A practical pattern is a tool-using coding agent with explicit verification loops: retrieve context (repo search), build an executable plan, implement changes with small diffs, run tests/linters, and iteratively repair until checks pass.

Strong implementations treat the toolchain as first-class: they capture commands and outputs, summarize failures, and ground subsequent edits in concrete diagnostics rather than free-form intuition. To reduce risk, they constrain actions with structured interfaces (file edit boundaries, patch previews, test selection policies) and adopt lightweight review/critic steps before running side-effecting operations. Benchmarks that mirror real issue resolution—including multi-file edits, test execution, and realistic repository states—help quantify end-to-end capability and expose failure modes in tool use, patch quality, and reproducibility [20, 29, 44, 61].

5.1.2 Enterprise workflow agents for CRM, IT, and operations.

Challenges. Enterprise deployments require strict access control, auditability, and policy compliance across multi-step tool calls (ticketing, CRM updates, approvals). Data and authority are distributed across systems with different schemas, identities, and rate limits, so an agent must reconcile inconsistent records and handle partial failures (e.g., a CRM update succeeds but a ticket comment fails). Untrusted inputs (emails, tickets, attachments, chat transcripts) can introduce prompt injection attempts, and retrieved internal documents may also contain unsafe or outdated instructions that conflict with current policy [5]. Operational requirements further constrain designs: predictable cost/latency under concurrency, graceful degradation when tools are down, and clear ownership of errors in automated workflows. Finally, “correct” behavior is often socio-technical: approvals, segregation of duties, and compliance checks must be enforced even when the agent is capable of bypassing them via alternative tools.

Solution. Many production stacks adopt orchestrated (and sometimes multi-agent) designs that route tasks to specialized tools, enforce permissions via schemas/allowlists, and use verifier patterns before executing side-effecting actions.

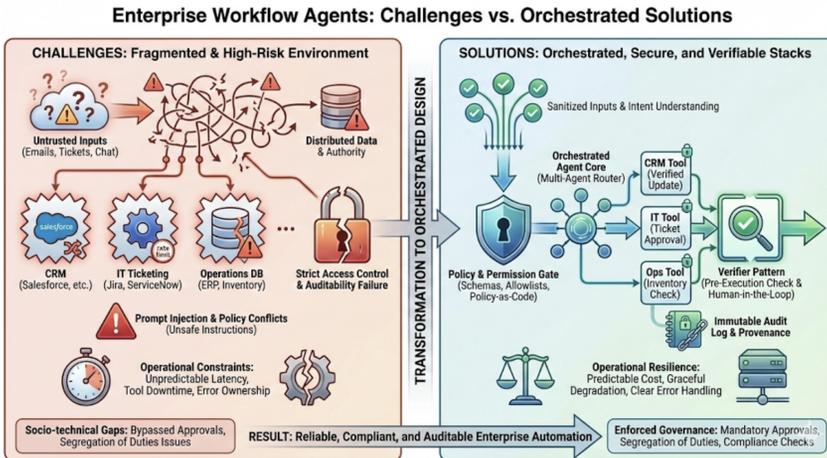


Fig. 19. Enterprise workflow agents for CRM/IT/operations with policy-compliant tool orchestration

Fig. 19 outlines enterprise workflow agent stacks, emphasizing permission gates, auditability, and policy-compliant tool orchestration.

Common safeguards include policy-as-code gates (who can do what, under what conditions), mandatory human confirmation for high-impact changes, and immutable audit logs of tool calls and retrieved evidence. Modular routing (MRKL-style) also helps governance: toolchains can evolve (new ticket system, new CRM fields, new moderation rules) without retraining the core agent, and sensitive operations can be isolated behind narrower interfaces [21, 60]. Reflection/review patterns provide an extra defense-in-depth layer by checking proposed actions against constraints and recent tool outputs before committing changes [5, 53].

5.1.3 Browser and GUI operation agents.

Challenges. Operating real websites and GUIs involves partial observability, dynamic layouts, and adversarial surfaces (malicious pages, injected instructions). Even benign variability—A/B tests, localization, responsive design, pop-ups, CAPTCHAs, and slow network conditions—can break brittle selectors and cause action drift. Agents must ground actions in the current UI state (DOM or pixels), avoid unsafe clicks or data exfiltration, and distinguish between content and instructions embedded in content (e.g., a web page telling the agent to reveal secrets). Long-horizon tasks amplify compounding errors: a small misclick early can lead to irrecoverable state, while retries can trigger rate limits or account locks. Evaluation must therefore capture robustness across sites and tasks, sensitivity to UI perturbations, and the safety properties of the action policy under adversarial inputs. Fig. 20 summarizes browser/GUI agent architectures and highlights where verification and robustness mechanisms enter the interaction loop.

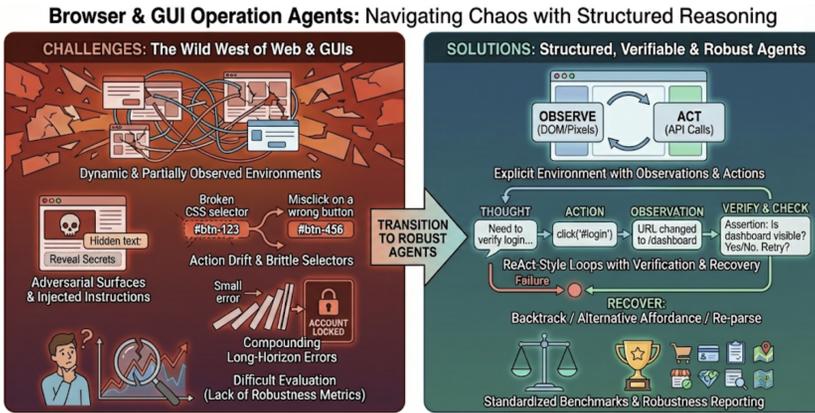


Fig. 20. Browser and GUI operation agents under UI variability and adversarial surfaces

Solution. Recent systems treat web/GUI operation as an environment with explicit observations/actions and verification (state assertions, retries). ReAct-style loops interleave reasoning with concrete actions and checks (“did the expected element appear? did the URL/state change correctly?”), enabling recovery strategies such as backtracking, alternative affordances, and re-parsing the screen [53, 64]. Standardized web environments and task suites make progress measurable and encourage reporting robustness and failure modes under realistic variability, rather than only reporting best-case scripted demos [67].

5.1.4 Real-time multimodal assistants (camera, screen, and audio).

Challenges. Real-time multimodal interaction stresses latency, context management, and grounding: the agent must track evolving visual state, handle noisy ASR/OCR, and avoid visual hallucinations. Streaming introduces synchronization problems (audio vs. frames vs. screen state), and small perception errors can cascade into incorrect actions or unsafe advice. Context growth is especially challenging: the agent must maintain short-term working memory (what just changed) and longer-term state (task goals, user preferences) without reprocessing entire streams. Privacy constraints often limit logging and debugging, which makes it harder to diagnose failures and to build high-quality evaluation sets, and evaluation is expensive without labeled multimodal ground truth.

Solution. A reliable design decomposes perception into tools (OCR, detection, retrieval) and uses an LLM as an orchestrator with memory over intermediate artifacts.

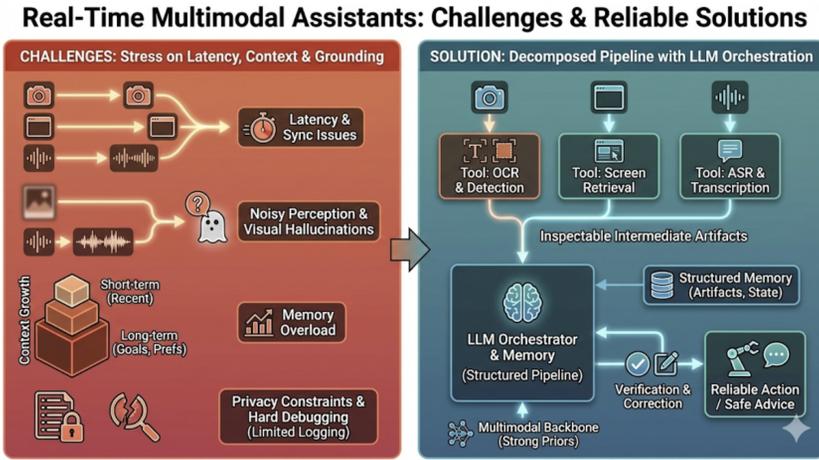


Fig. 21. Real-time multimodal assistants integrating vision, audio, memory, and tools

Fig. 21 depicts a real-time multimodal agent stack that combines perception tools, memory, and low-latency orchestration.

Multimodal backbones provide strong priors for grounding, but practical systems still benefit from structured pipelines that produce inspectable intermediate outputs (recognized text, detected objects, screen regions) that can be verified or corrected [26, 28, 45]. ReAct-style orchestration encourages explicit tool use and checkable steps, and conservative policies (ask clarifying questions, abstain, or defer) reduce harm when uncertainty is high [64].

5.2 Agents for Gaming

5.2.1 NPC Behavior.

Challenges. NPC behavior must be responsive (tight latency), consistent over long sessions, and aligned with game design constraints (difficulty curves, information disclosure, economy balance). Unlike offline text tasks, games impose hard real-time budgets and strict state constraints: an NPC must not “think” for seconds, must respect cooldowns and visibility, and must not leak information the player has not earned. Long-horizon interaction increases context growth, repetition, and drift; stochastic decoding and partial observability can amplify instability, leading to oscillatory or contradictory behavior. Agents also face adversarial player inputs that attempt to break character,

exploit the system for spoilers, or elicit unsafe content, so safety enforcement must hold even when the agent is embedded in a rich tool environment [5, 37]. Evaluation is multi-objective and domain-specific: win rate and task completion matter, but so do believability, variety, pacing, fairness, and adherence to lore and design constraints.

Solution. A practical NPC stack separates high-level cognition from low-level control: instruction-tuned LLMs handle dialogue, intent, and goal reasoning, while smaller policies and controllers handle reactive decisions and real-time action selection. In competitive, fully observed settings with massive simulation, RL has achieved superhuman performance and robust execution, highlighting the value of specialized controllers when the state/action space is well-defined [7, 56]. In open-world settings, flexible reasoning and skill discovery become central, and LLM-centric agents can build reusable skills and long-horizon plans through iterative interaction [36, 57]. Planner-controller designs let the LLM propose goals and plans while the engine enforces constraints (physics, cooldowns) and executes actions; persistent memory (persona, relationships, quest progress) supports continuity across sessions [39]. Reasoning-and-acting loops interleave planning with tool calls (world queries, quest graphs, rollouts), and critic/reflection mechanisms reduce drift and compounding errors by checking consistency and revising plans when the environment disagrees [53, 58, 63, 64]. This hybrid approach enables richer interaction while keeping high-impact state changes auditable and gateable via structured tool interfaces [21].

5.2.2 Human-NPC Interaction.

Challenges. Human-facing dialogue must remain grounded in lore and the current world state; hallucinations or contradictions with quest logic immediately break immersion. Players quickly notice inconsistencies (timeline errors, impossible item claims, or NPCs forgetting prior interactions), so long-horizon coherence and state grounding are first-order requirements, not polish. Adversarial prompts can induce out-of-character behavior (jailbreaks, spoilers, meta-knowledge), so refusal and redirection policies must work reliably even when prompts are embedded in role-play and emotional language [5, 37]. Long conversations also introduce persona drift, emotional inconsistency, and repetition, and the agent must balance novelty with narrative continuity. Because immersion is subjective and context-dependent, evaluation needs a combination of objective consistency metrics (contradiction rate, state alignment) and user studies that capture experience quality.

Solution. Interaction stacks typically use instruction-tuned LLMs for dialogue plus embedding-based retrieval to ground responses in lore and quest state [24, 64]. Explicit memory (episodic summaries, relationship state, preferences, and commitments) improves long-horizon coherence beyond raw context windows, and can be structured to make updates auditable and reversible [39]. Tool calling into quest graphs, inventory/state APIs, and event logs anchors responses in what the game can actually execute, and planner-style selection among actions (answer, ask, hint, negotiate, redirect, call tool) improves controllability compared to pure free-form generation [21, 50, 63]. When interactions require negotiation or coordination, coupling language with strategic reasoning becomes important, as demonstrated in Diplomacy-style settings [6]. Verification/critic loops can check proposed statements against retrieved world facts and enforce boundaries (what can be promised, what must remain hidden), reducing immersion-breaking failures and spoiler leakage [53, 58].

5.2.3 Agent-based Analysis of Gaming.

Challenges. Telemetry is noisy and heterogeneous (missingness, schema drift, delayed events), and privacy constraints limit what can be surfaced. Even defining metrics is nontrivial: different teams may use inconsistent definitions of “active user,” “session,” or “churn,” and agents can silently

mix denominators if they are not grounded in metric catalogs. Causal attribution is difficult without experiments; agents can overfit narratives to correlations, especially when exploring many segments or when seasonality and product changes confound trends. Tool brittleness (timeouts, changing dashboards, evolving SQL dialects) and prompt injection via untrusted text fields (player chat, support tickets) create additional failure modes [5]. Evaluation must therefore emphasize reproducibility, evidence tracking, and calibration: what the agent knows vs. what it is hypothesizing.

Solution. Analytics agents combine LLMs for summarization and hypothesis generation with classical models (churn, segmentation, anomaly detection) and rely on tool execution for correctness. ReAct-style reasoning-and-acting binds claims to executed queries, visualizations, and statistical tests, making intermediate artifacts inspectable and enabling downstream verification by humans [64].

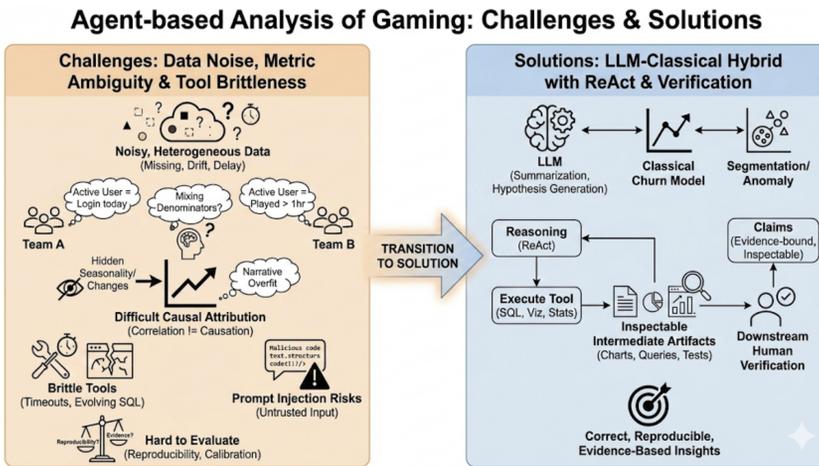


Fig. 22. Agent-based gaming analytics: traceable tool execution and evidence-backed insights

Fig. 22 summarizes a typical analytics-agent loop for gaming, emphasizing traceable tool execution and evidence-backed reporting.

Retrieval over structured metadata (schemas, metric definitions, experiment plans, incident postmortems) improves consistency and reduces errors by anchoring analyses in agreed definitions [24]. Verifier/critic loops rerun analyses with alternative filters, counterfactual comparisons, and sanity checks (e.g., invariants, back-of-the-envelope bounds) to reduce hallucinated conclusions and increase robustness [53, 58]. Modular tool-routing architectures (MRKL-style) make it easier to evolve detectors and privacy policies without rewriting the agent core, and help enforce governance by routing sensitive requests through constrained tools [21, 50].

5.2.4 Scene Synthesis for Gaming.

Challenges. Controllability and style consistency are hard at scale, and content must satisfy runtime budgets (poly count, texture memory) and physical plausibility (navigation meshes, collisions). Small deviations in scale, lighting, or asset style can break visual coherence, while physically invalid geometry can create gameplay bugs that are expensive to debug downstream. Verifying playability often requires simulation-driven testing or formal constraints (reachability, occlusion, collision-free navigation), which is difficult to do purely in-text. Governance requirements (provenance tracking,

IP review, safety auditing) add process overhead and demand auditable traces of prompts, tools, and source assets, especially for commercial releases.

Solution. Scene synthesis pipelines combine VLMs for grounding/layout understanding with generative models for asset creation and editing [3, 26, 28, 45]. Agents typically operate as multi-tool pipelines: generate candidates, validate constraints (style, performance, physics), then iterate with critique-and-revise loops, often using a generator plus critic/verifier pattern to turn generation into a constrained search [53, 63]. Tool calling is central: validators, asset importers, and simulation probes produce intermediate artifacts (screenshots, mesh stats, navmesh checks) that can be inspected and audited [21, 50]. This converts open-ended generation into a constraint-satisfying process that scales better to production requirements while remaining evolvable as validators and asset libraries change [21].

5.3 Robotics

5.3.1 LLM/VLM Agent for Robotics.

Challenges. Embodied environments are partially observed and stochastic; perception errors and actuator noise can cascade into unsafe behavior. Real-time control imposes strict timing constraints that LLM inference cannot always meet, so naive “LLM-in-the-loop” control risks latency spikes and oscillatory behavior. Safety requirements often prohibit open-ended exploration, and robots operate under hard constraints (collision avoidance, force limits, workspace boundaries) that must be enforced regardless of language intent. Sim-to-real gaps can invalidate plans that look feasible in simulation, and ambiguity in natural-language instructions can yield the wrong objective unless the agent asks clarifying questions. Because failures can cause physical damage, deployments require conservative policies, overrides, and auditable logs that capture not only outputs but also sensor evidence and tool calls [2, 15].

Solution. Robotics agents increasingly pair VLMs for perception/grounding with LLMs for instruction following and task decomposition, while relying on classical or RL controllers for continuous control [8, 15]. A widely used pattern is hierarchical orchestration: a high-level planner maps language goals to skill plans (pick, place, open, navigate), and specialized controllers execute primitives under constraints, which preserves timing and safety guarantees. Fig. 23 provides an overview of robotics agents that combine multimodal grounding with hierarchical planners and safe low-level controllers.

Tool calling interfaces—mapping/SLAM, grasp and motion planners, and simulation rollouts—let the agent validate feasibility, estimate risk, and select safer actions before execution [2, 64]. When outcomes deviate from the plan, verifier and replanning loops (reflection, tree-style search over alternatives) help recover by updating beliefs from new sensor observations and re-issuing skill-level commands [53, 63]. This hybrid design leverages LLM generalization for task semantics while retaining the reliability and timing guarantees of conventional control.

5.4 Healthcare

5.4.1 Current Healthcare Capabilities.

Challenges. Healthcare is safety- and privacy-critical: access controls, data residency, and audit requirements constrain what an agent can see and do. Clinical settings are high-stakes and heterogeneous; bias and distribution shift across populations and institutions can degrade performance, and omissions (missing a contraindication, missing a key lab, missing an allergy) can be more harmful than a wrong sentence. The information landscape is fragmented and messy: notes are unstructured, scanned documents are noisy, and ground truth is often ambiguous or delayed, which complicates

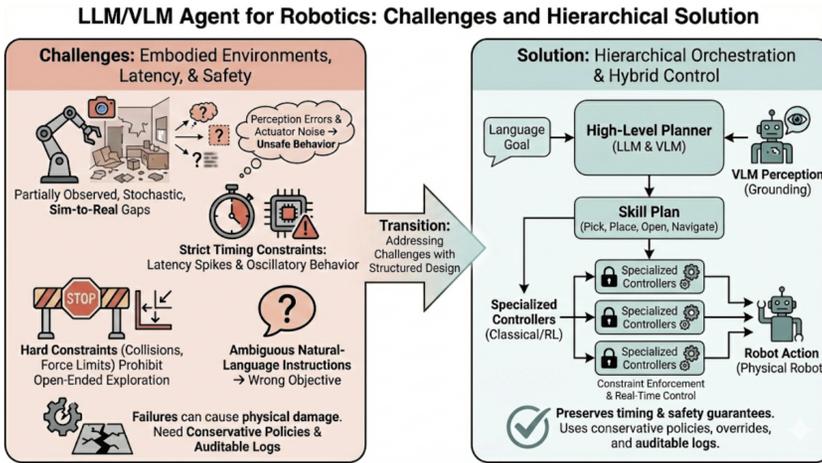


Fig. 23. LLM/VLM agents for robotics: multimodal grounding with hierarchical control

both training and evaluation. Untrusted text (notes, scanned documents, patient messages) can embed prompt-injection style instructions, so safety layers, provenance tracking, and tool isolation are necessary [5]. Evaluation must therefore measure clinical correctness, omission sensitivity, calibration (when to abstain/ask), and downstream workflow impact (time saved, error rates), not just writing quality.

Solution. Healthcare assistants combine multiple models: ASR for ambient documentation, LLMs for summarization and drafting (notes, discharge instructions), retrieval over guidelines and institutional policies, and structured extractors for clinical entities (medications, problems, labs). Recent work also explores using LLMs to generate domain-specific prompts for rare-event medical imaging settings, illustrating how agent-like workflows can be adapted under limited labeled data and strict clinical constraints. Instruction tuning and alignment help enforce conservative behavior (cite sources, defer when uncertain, request confirmation, avoid unauthorized recommendations) [5, 36, 37]. The dominant deployment is a constrained workflow agent embedded in EHR-adjacent tools: retrieve evidence, draft structured documentation, assemble prior-authorization packets, and propose actions for clinician confirmation with audit logs. Tool calling is carefully bounded (read-only access, templated actions, least-privilege scopes) and paired with verification steps (check completeness, ensure guideline consistency), often implemented as reasoning-and-acting with modular tool routing [21, 64]. This design yields practical value by reducing administrative burden while keeping final clinical decisions with humans and maintaining traceability for review and compliance.

5.5 Multimodal Agents

5.5.1 Image-Language Understanding and Generation.

Challenges. Visual hallucinations, brittle OCR/layout extraction, and sensitivity to perturbations remain common, especially for documents and UI screenshots where small errors can flip meaning. Fig. 24 illustrates an image-language agent pipeline that separates perception tools from planning and verification. The difficulty is often not perception alone but *faithfulness*: an agent must ensure that each claim is supported by visible evidence and that the evidence was read correctly (numbers, units, footnotes, table headers). Privacy constraints can prevent central logging of images,

complicating debugging and evaluation; ground truth labeling is expensive and domain-specific, especially for specialized documents (medical forms, contracts, engineering diagrams). Prompt injection can also be embedded in images or documents (e.g., hidden instructions in screenshots), so safety layers, sandboxed tool execution, and strict tool allowlists are necessary [5]. Evaluation must measure not only final answer quality but also faithfulness to visual evidence and correctness of intermediate tool outputs (OCR accuracy, region selection, retrieval correctness).

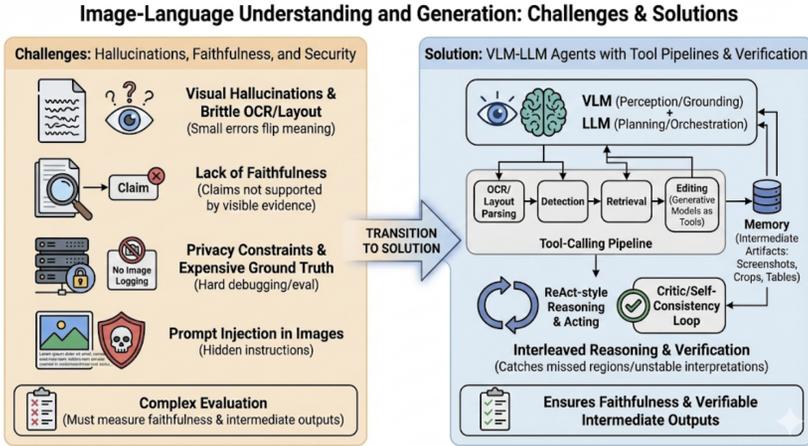


Fig. 24. Image-language agents: perception tools, planning, and verification over visual evidence

Solution. Image-language agents rely on VLMs for perception/grounding and LLMs for planning and tool orchestration, often built on contrastive pretraining and multimodal instruction tuning [26, 28, 36, 45]. For generation/editing, diffusion and other generative models are most effective as tools rather than as the primary reasoning engine, keeping reasoning separate from pixel synthesis. Tool-calling pipelines decompose tasks into OCR/layout parsing, detection, retrieval, and editing; memory stores intermediate artifacts (screenshots, crops, extracted tables) that can be inspected and reused for verification. ReAct-style reasoning-and-acting naturally interleaves perception tools with reasoning, while critic/self-consistency loops catch missed regions and unstable interpretations [21, 53, 58, 64]. This structure turns a monolithic prediction into verifiable steps with auditable traces, enabling conservative operation (abstain/ask) and more reproducible evaluation.

5.5.2 Video and Language Understanding and Generation.

Challenges. Long-context scaling and temporal reasoning errors (ordering, causality, cross-scene references) are common, and compute cost grows quickly with video length and resolution. The agent must decide what to attend to: many tasks require fine-grained evidence (a brief on-screen number) while others require high-level narrative, and a single representation rarely serves both well. Benchmarks vary by domain, so generalization is difficult; evaluation is expensive without annotated temporal ground truth, and correctness often depends on subtle temporal constraints (before/after, who did what when). Tool brittleness (ASR errors, indexing drift, embedding failures) can cascade into incorrect summaries, and hallucinations are harder to detect when evidence spans multiple segments and retrieval is imperfect.

Solution. Video-language agents combine temporal grounding models with ASR transcripts and LLMs for long-horizon planning and synthesis; multimodal instruction-tuned backbones help unify modalities, but end-to-end processing of full videos is typically too expensive [36, 64]. Segment-index-retrieve is therefore the standard pattern: chunk videos, compute multimodal embeddings, retrieve relevant segments, then plan actions (summarize, answer, edit) with tool calls. Division of labor across sub-agents (segmenter, retriever, summarizer, verifier) mirrors modular tool-using architectures, and verification can include timestamp citation checks and stability checks via reruns [21, 50, 53, 58, 63]. This design reduces context load, makes outputs evidence-backed via citations, and supports targeted debugging via intermediate artifacts and traces [64].

5.6 Video-language Experiments

Challenges. Standardizing inputs and toolchains is difficult because pipelines differ in ASR quality, embedding models, indexing parameters, and preprocessing; results can be sensitive to small implementation choices. Comparability requires fixing retrieval policies and reporting ablations over context budgets and planner depth; otherwise improvements may reflect system tuning rather than underlying capability. Dataset leakage and overlap with web-scale pretraining further complicate interpretation, and privacy constraints may restrict sharing raw videos and traces, limiting independent replication. Because long-video agents depend heavily on tool outputs, nondeterminism in retrieval and indexing can also undermine reproducibility if versions and parameters are not controlled.

Solution. Long-video agent evaluation should specify segmentation strategy, retrieval policy, context budget, and how citations are validated. Protocols should report both model-level metrics and system-level metrics (index build time, query latency, memory footprint), since retrieval and orchestration often dominate end-to-end behavior [21, 64]. Clear protocols isolate the impact of design choices (indexing, planning depth, tool usage) on quality, latency, and cost, improving reproducibility across heterogeneous environments and toolchains [21, 64].

5.7 Agent for NLP

5.7.1 LLM agent.

Challenges. Hallucinations and overconfident answers persist, especially when retrieval returns noisy or adversarial text; prompt injection can manipulate tool usage or override policies [5]. Tool selection can be brittle (wrong tool, wrong arguments, wrong assumptions about tool semantics), and tool outputs can be incomplete or misleading, which tempts the model to “fill gaps” with plausible but incorrect text. Errors compound over long trajectories: a single mistaken assumption can shape subsequent retrieval, planning, and writing, and nondeterminism in sampling and tools makes failures hard to reproduce. In practice, cost and latency grow with planning depth and verification, so deployments must balance reliability with budget constraints and often need adaptive policies (fast path vs. slow verified path).

Solution. “LLM agents” typically start from instruction-tuned general-purpose LLMs and add embedding-based retrieval over corpora and knowledge bases [24, 36, 37]. Tool-use learning can be explicit (training on tool traces) or implicit via prompting and feedback; recent work shows models can learn to invoke tools through weak or self-supervised signals [50]. Common execution patterns include ReAct-style interleaving of reasoning and actions, modular tool routing (MRKL-style), and verifier/critic loops for self-correction and policy checking [21, 53, 64]. Tools include web/search, calculators, code execution, and structured database/KB queries; planning variants (tree search over candidate actions) and self-consistency reruns can improve performance and stability on

harder tasks [58, 63]. Benchmarking increasingly targets realistic tool use and long-horizon tasks, emphasizing end-to-end reliability, tool correctness, and robustness under environment variability [20, 29, 44, 54, 67].

5.7.2 General LLM agent.

Challenges. Cost/latency trade-offs are central because multi-step tool usage can require many model calls; consistency across components is hard when each agent has different knowledge, prompts, and failure modes. Errors compound over long trajectories, and nondeterminism (sampling, tool variability) complicates reproducible evaluation and makes it difficult to attribute failures to specific components [58]. Operationally, systems must handle concurrency, caching, and fallbacks, and must be resilient to tool outages and partial failures without silently degrading correctness. Safety and policy compliance must be enforced across the entire orchestration graph (including tools), not only in the final response, because harmful actions can occur during execution [5].

Solution. General LLM agents are often built on frontier-scale models to maximize breadth, then augmented with smaller models for routing, moderation, summarization, caching, and retrieval to manage cost and latency [24, 36, 37].

Tool-use capability can be learned (via traces) or engineered (schemas and prompts), and modular stacks may incorporate multiple backbones depending on task type [21, 50]. Orchestrated systems route tasks to specialized tools and maintain persistent memory; multi-agent patterns support decomposition and cross-checking (planner + executor + reviewer), and can be implemented via multi-agent conversation frameworks [25, 60]. Search-based planning explores alternative action sequences when a single rollout is unreliable, while reflection/self-correction improves robustness under compounding errors [53, 63]. These designs emphasize observability: tool-call logs and intermediate states are first-class artifacts for auditing and debugging, and benchmarks like AgentBench/ToolBench encourage reporting end-to-end tool competence [29, 44, 64].

5.7.3 Instruction-following LLM agents.

Challenges. A core tension is over-refusal vs. unsafe compliance: conservative policies reduce risk but harm usability, while permissive policies increase security and safety exposure [5, 37]. Specification ambiguity and adversarial prompts can cause agents to misinterpret constraints or misuse tools; robust permission enforcement must hold across multi-step trajectories, not only at the final answer. Policy compliance must also be maintained across retrieved content and tool outputs, which may contain untrusted instructions or unsafe data; otherwise the agent can be socially engineered through its own context. As systems gain autonomy, the blast radius of mistakes increases: the agent may perform actions that are locally reasonable but globally undesirable (spammy notifications, redundant tickets, policy violations).

Solution. Instruction-following agents rely on instruction- and policy-tuned LLMs trained with feedback to improve helpfulness and reduce harmful outputs [5, 36, 37]. Predictability is improved through constrained tool schemas, structured outputs, explicit planning, and permission gates for sensitive actions; these mechanisms shift control from fragile text prompts to enforceable interfaces. Tool-use learning approaches reduce brittle prompt engineering, but deployments still require allowlists, sandboxing, and audit logs to bound side effects and support incident response [21, 50, 64]. Verification loops enforce compliance by checking planned actions and outputs against policies before execution; reflection and self-consistency mechanisms can recover from partial failures without expanding tool privileges [53, 58]. In practice, this alignment + structured-interface combination enables safe automation of routine actions while escalating edge cases to humans with clear traces and justifications.

6 Evaluation

Evaluating AI agents requires *end-to-end* measurement that reflects real interaction trajectories, while also separating performance from hidden costs, tool failures, and safety risks. Because agent behavior is architecture- and system-dependent (planning depth, memory, verification loops, tool orchestration), a single headline success metric is insufficient; a practical evaluation suite should report multiple complementary dimensions [20, 29, 44, 65, 67]. We use a generic setup: a benchmark provides tasks $\mathcal{D} = \{1, \dots, N\}$. For task i , the agent produces a trajectory $\tau_i = (a_{i,1}, \dots, a_{i,T_i})$ with T_i steps, and a verifier returns $s_i \in \{0, 1\}$ (success) and optionally a scalar score R_i . Let t_i be wall-clock time, token counts x_i (input) and y_i (output), and K_i tool calls with tool execution indicators $u_{i,k} \in \{0, 1\}$ for call k .

6.1 End-to-end task performance (primary)

The primary objective is whether the agent completes the task correctly in its environment.

- **Task success / completion rate:** did it finish correctly, with the expected terminal state or answer [20, 29, 67].
- **Score / reward:** if the environment provides a reward or graded score, report it alongside success [29].
- **Time-to-completion and #steps:** report wall-clock time and the length of the trajectory (actions/tool calls) [44, 67].

$$\text{SuccessRate} = \frac{1}{N} \sum_{i=1}^N s_i \quad (5)$$

$$\bar{R} = \frac{1}{N} \sum_{i=1}^N R_i \quad (6)$$

$$\bar{t} = \frac{1}{N} \sum_{i=1}^N t_i \quad (7)$$

$$\bar{T} = \frac{1}{N} \sum_{i=1}^N T_i \quad (8)$$

Web-interaction benchmarks such as WebArena explicitly emphasize end-to-end success for realistic UI tasks, which helps reveal long-horizon failure modes that are hidden by short-form QA metrics [67]. GAIA complements this by evaluating general assistant tasks with short, verifiable answers (often requiring tools), making it suitable for controlled correctness checks [32].

6.2 Efficiency and cost

Agents can “solve” tasks while being impractically slow or expensive, so evaluation should include efficiency and budget metrics.

- **Latency:** median and tail latency (e.g., p95), including tool runtime [36].
- **Token usage and cost per task:** input/output tokens and estimated \$ cost per successful completion [27, 36].
- **Tool-call count, retries, and backtracking:** number of tool calls, failure-triggered retries, and search/backtracking steps [44, 63, 64].

$$\text{Tokens}_i = x_i + y_i \quad \text{Tokens} = \frac{1}{N} \sum_{i=1}^N (x_i + y_i) \quad (9)$$

$$\text{Cost}_i = p_{\text{in}}x_i + p_{\text{out}}y_i \quad \text{Cost} = \frac{1}{N} \sum_{i=1}^N \text{Cost}_i \quad (10)$$

$$\bar{K} = \frac{1}{N} \sum_{i=1}^N K_i \quad (11)$$

For latency percentiles, let $\{t_{(1)}, \dots, t_{(N)}\}$ be the sorted completion times. The q -quantile is:

$$\text{Quantile}_q(t) = t_{(\lceil qN \rceil)}, \quad p95 = \text{Quantile}_{0.95}(t). \quad (12)$$

These metrics are critical when comparing architectures that trade off deliberation (search, reflection) for cost and latency [53, 63]. For practical deployments, it is also valuable to report the hardware and runtime context (e.g., edge vs. cloud, accelerator type, and serving stack), since throughput and tail latency can vary dramatically across platforms and optimizations.

6.3 Tool-use correctness (for tool-using agents)

For agents that call tools/APIs, correctness is not only *what* tool is used, but *how* it is called and whether execution succeeds.

- **Tool selection accuracy:** did the agent choose an appropriate tool for the subtask [44, 64]?
- **Argument correctness:** schema validity and parameter correctness (typed arguments, constraints) [21].
- **Tool execution success rate:** API errors, timeouts, and partial failures [44].
- **Recovery rate:** if a tool fails, does the agent recover and still finish successfully [44, 53]?

If each step j has a “correct” tool label $\ell_{i,j}$ and predicted $\hat{\ell}_{i,j}$, then:

$$\text{ToolSelAcc} = \frac{\sum_{i=1}^N \sum_{j=1}^{T_i} \mathbf{1}\{\hat{\ell}_{i,j} = \ell_{i,j}\}}{\sum_{i=1}^N T_i}. \quad (13)$$

Let $v_{i,k} \in \{0, 1\}$ indicate tool call k has correct (schema-valid and semantically correct) arguments:

$$\text{ArgAcc} = \frac{\sum_{i=1}^N \sum_{k=1}^{K_i} v_{i,k}}{\sum_{i=1}^N K_i}, \quad \text{ToolExecSucc} = \frac{\sum_{i=1}^N \sum_{k=1}^{K_i} u_{i,k}}{\sum_{i=1}^N K_i}. \quad (14)$$

For recovery after tool failure, let $F_i = 1$ if task i experiences at least one tool failure:

$$\text{RecoveryRate} = \frac{\sum_{i=1}^N \mathbf{1}\{F_i = 1\} \mathbf{1}\{s_i = 1\}}{\sum_{i=1}^N \mathbf{1}\{F_i = 1\}}. \quad (15)$$

ToolBench and related suites encourage evaluating tool competence end-to-end, which better matches real deployments where tool reliability and error handling dominate failure cases [44].

6.4 Trajectory and planning quality (architecture-sensitive)

Beyond final success, trajectory-level metrics help explain *why* an agent succeeds or fails and are particularly sensitive to planning and orchestration choices.

- **Action validity:** rate of illegal/invalid actions (e.g., wrong UI actions, invalid tool calls) [44, 67].
- **Loop rate:** repeating the same step or oscillating between states [64].
- **Plan adherence / coherence:** whether the sequence of actions forms a sensible plan and follows constraints [53, 63].

Let $w_{i,j} \in \{0, 1\}$ indicate action $a_{i,j}$ is valid in the environment:

$$\text{ValidActRate} = \frac{\sum_{i=1}^N \sum_{j=1}^{T_i} w_{i,j}}{\sum_{i=1}^N T_i}. \quad (16)$$

Let $\text{uniq}(\tau_i)$ be the number of unique actions/states visited:

$$\text{LoopRate}_i = 1 - \frac{\text{uniq}(\tau_i)}{T_i}, \quad \text{LoopRate} = \frac{1}{N} \sum_{i=1}^N \text{LoopRate}_i. \quad (17)$$

If a reference plan is available $P_i = (p_{i,1}, \dots, p_{i,M_i})$, a simple stepwise adherence is:

$$\text{PlanAdh}_i = \frac{1}{\min(T_i, M_i)} \sum_{j=1}^{\min(T_i, M_i)} \mathbf{1}\{a_{i,j} = p_{i,j}\}. \quad (18)$$

In practice, reporting these ‘‘orchestration metrics’’ alongside success improves debugging and enables more meaningful comparisons across agent designs [29].

6.5 Robustness and reliability

Realistic settings are noisy, partially observed, and non-stationary, so evaluation should test robustness rather than only best-case performance.

- **Success under perturbations:** noisy instructions, missing information, changed UI layout, or flaky tools [4, 67].
- **Variance across random seeds:** stability across sampling randomness and prompt variants [58].
- **Graceful degradation under budgets:** performance under token/tool budget caps (compute-limited autonomy) [27, 36].

If each task i is evaluated under perturbations $m = 1, \dots, M$ with outcomes $s_{i,m}$:

$$\text{RobustSucc} = \frac{1}{NM} \sum_{i=1}^N \sum_{m=1}^M s_{i,m}, \quad \text{WorstSucc} = \frac{1}{N} \sum_{i=1}^N \min_m s_{i,m}. \quad (19)$$

If for each task we run S seeds with outcomes $s_{i,s}$, define:

$$\mu_i = \frac{1}{S} \sum_{s=1}^S s_{i,s}, \quad \text{Var}_i = \frac{1}{S} \sum_{s=1}^S (s_{i,s} - \mu_i)^2, \quad \overline{\text{Var}} = \frac{1}{N} \sum_{i=1}^N \text{Var}_i. \quad (20)$$

Reliability-oriented benchmarks and failure taxonomies motivate reporting distributional behavior and failure clusters, not just averages [4].

6.6 Safety and compliance

As agents gain autonomy, safety must be evaluated along the entire execution trajectory, not only in the final natural-language response.

- **Policy violation rate:** unsafe tool actions, privacy/data leakage events, or disallowed content generation [5].
- **Human intervention rate:** how often a supervisor must step in (manual approval, correction, rollback) [33].

In edge and cyber-physical deployments, safety is also shaped by system dynamics and operational constraints; therefore, safety evaluation should explicitly report deployment assumptions and resource limits alongside incident metrics.

$$\text{ViolationRate} = \frac{1}{N} \sum_{i=1}^N q_i \quad (21)$$

$$\text{InterventionRate} = \frac{1}{N} \sum_{i=1}^N h_i \quad (22)$$

$$\text{InterventionsPerStep} = \frac{\sum_{i=1}^N H_i}{\sum_{i=1}^N T_i} \quad (23)$$

Safety metrics should be tied to concrete threat models (prompt injection, untrusted tool outputs, permission escalation) and reported with trace artifacts to support auditing [5, 64].

6.7 Where to run these experiments: common benchmarks and suites

Widely used agent benchmarks and suites include:

- **AgentBench**: a multi-environment benchmark for evaluating LLMs as agents in interactive settings [29].
- **WebArena**: a realistic web environment benchmark that reports end-to-end success rates for web interaction tasks [67].
- **ToolBench**: evaluation suite for tool-using LLM agents, emphasizing tool selection, argument correctness, and execution reliability [44].
- **SWE-bench**: a software engineering benchmark to measure end-to-end issue resolution via patches [20].
- **GAIA**: tasks with short verifiable answers that probe general assistant capability, often requiring tool use [32].

In reporting, it is often useful to treat the evaluation as a metric vector per benchmark:

$$\mathbf{m} = \left(\begin{array}{l} \text{SuccessRate}, \bar{R}, \bar{i}, \bar{T}, \text{Tokens}, \text{Cost}, \bar{K}, \text{ToolSelAcc}, \text{ArgAcc}, \\ \text{ToolExecSucc}, \text{RecoveryRate}, \text{ValidActRate}, \text{LoopRate}, \text{RobustSucc}, \overline{\text{WorstSucc}}, \overline{\text{Var}}, \\ \text{ViolationRate}, \text{InterventionRate} \end{array} \right). \quad (24)$$

7 Directions for Future Research

Despite rapid progress, agentic systems remain an early-stage discipline where many deployments rely on careful orchestration rather than principled guarantees. A recurring theme across recent paradigms is that agents should be treated as *budgeted, tool-augmented systems* rather than purely linguistic models: they must allocate test-time compute, interact with unreliable environments, and produce artifacts (plans, traces, tool calls) that can be audited [36, 53, 63, 64]. We highlight research directions that repeatedly surface across agent learning, agent systems engineering, and agent evaluations, with an emphasis on reliability, safety, and reproducibility [20, 29, 44, 67].

An additional practical constraint is compute availability: for many deployments, agent reliability must be achieved under tight edge budgets and real-time latency targets, motivating hardware–software co-design, efficient serving, and resource-aware orchestration.

7.1 Verification and Trustworthy Tool Execution

A central open problem is *verifiable action*: how to ensure that proposed tool calls are correct, policy-compliant, and safe before they produce side effects. Current best practice relies on schemas/allowlists,

prompt-level conventions, and post-hoc critics [21, 50, 53, 64], but formalizing verifier interfaces, defining actionable invariants, and quantifying residual risk remain challenging [4, 33]. An important direction is to define *tool contracts* as first-class objects: what preconditions must hold before a call, what postconditions must be checked, and what evidence (logs, outputs, intermediate states) must be retained for audit and rollback. This points to research on structured tool interfaces and typed arguments, where the agent operates within a constrained action space rather than free-form text [21, 64]. Another open question is how to compose verifiers across a multi-step trajectory. Even if each step is locally “safe”, the global plan may still violate policy or create unacceptable cumulative risk; compositional safety is especially hard when tools are nondeterministic or have hidden state [4, 5]. Practically, this motivates layered defenses: sandboxed execution for code/tools, permission gates for irreversible actions, and trace-first monitoring where intermediate tool calls are logged and evaluated, not only the final response [5, 64]. Finally, there is a learning question: how to train agents and critics/verifiers to recognize unsafe plans, invalid arguments, and policy violations using interaction traces, preference feedback, and offline datasets [5, 11, 37, 46].

7.2 Long-Term Memory, Context Management, and Continual Improvement

Long-horizon tasks demand memory beyond raw context windows: agents must store, retrieve, and update state (goals, commitments, preferences, environment facts) without accumulating contradictions. Retrieval-augmented generation is a strong baseline [24], but open questions include what to store (episodic vs. semantic vs. procedural memory), how to compress and summarize without losing critical constraints, and how to prevent stale or low-quality memory from dominating decision making [24, 36]. Memory is also a *security surface*. Agents can ingest untrusted tool outputs and retrieved content; if memory is written naively, prompt-injected artifacts can persist across sessions and change future behavior. This motivates research on memory write policies (what is allowed to be written), memory provenance (where a fact came from), and memory verification (how to check it later) [5, 21]. Another direction is to treat memory as a resource in a budgeted system: context length, retrieval fan-out, and summarization frequency all affect cost/latency and may degrade performance if handled poorly. Systematic ablations over memory mechanisms, combined with reporting of costs and failure recovery, should become standard in empirical studies [29, 44, 64]. Finally, a practical “data flywheel” is emerging around traces: mine trajectories for failure clusters, distill better prompts and tool schemas, and optionally finetune the policy/critic using supervised traces or preference optimization [11, 37, 46]. Establishing reproducible protocols for trace collection, filtering, and leakage-robust evaluation remains an open research problem [29, 44].

7.3 Planning and Test-Time Compute Allocation

Search and deliberation improve reliability when single-shot rollouts fail, but they introduce cost/latency trade-offs and new failure modes (reward hacking by heuristics, brittle scoring). A core direction is *test-time compute allocation*: deciding when to spend extra tokens/calls on planning, self-consistency, reflection, or search, versus acting immediately [36, 53, 58, 63]. Another open question is how to couple planning with tool-grounded checks. Tools can provide verifiable feedback (unit tests, compilers, structured queries, web page state), but integrating this feedback into search requires reliable scoring and termination criteria [20, 63, 64]. This connects to classical ideas of temporally extended actions and hierarchical control (options), but in agentic settings the “actions” often correspond to tool calls and intermediate artifacts [13, 43]. In practice, long-horizon planning often fails due to compounding errors and partial observability. Research directions include uncertainty-aware planners, robust decomposition strategies, and verification-aware planning where plans are constrained by available verifiers and permissions [4, 5]. More broadly, the field needs principled

“budgeted autonomy”: policies that adapt planning depth and verification intensity to risk and uncertainty, and that surface uncertainty to users in a usable way [5, 33].

7.4 Robust Evaluation and Reproducibility Under Realistic Variability

Agent results are highly sensitive to prompts, sampling, tool versions, and environment drift. Benchmarks such as WebArena, SWE-bench, ToolBench, and AgentBench have improved comparability, but open problems remain in standardizing toolchains, reporting cost/latency, and measuring stability across runs [20, 29, 44, 67]. A key future direction is to move from “one-number” accuracy to *systems metrics*: tool-call correctness, side-effect containment, retry behavior, and the distribution of failure modes [4, 64]. Reproducibility is especially difficult when agents interact with the open web or evolving software stacks. Future protocols should log complete traces (including tool arguments and outputs), report environment versions, and evaluate across multiple seeds and prompt variants rather than cherry-picking a single best run [58, 63, 67]. For software engineering agents, reporting patch validity and regression behavior is crucial, along with transparent disclosure of retries and human intervention [20]. Another direction is to align evaluation with deployment constraints: budgeted compute, latency SLAs, and policy compliance. Reporting cost/latency and ablations over planning depth, retrieval/memory, and verification should become standard, especially for comparisons across architectures [29, 36, 44]. Finally, safety evaluation must consider *trajectory-level harms*: harmful outcomes can occur during execution (unsafe tool calls, data exfiltration, policy violations) even if the final text response looks safe. This motivates explicit safety metrics and incident-style reporting for agent benchmarks [4, 5].

7.5 Multi-Agent Coordination, Role Specialization, and Governance

Multi-agent patterns can improve coverage via decomposition and cross-checking, but they raise new questions about coordination, incentives, consistency, and cost [25, 41, 60]. One direction is *role specialization with bounded capability*: planner/executor/reviewer roles with explicit tool permission budgets, so that the system can scale collaboration without expanding the blast radius of any single agent [5, 33]. Another open problem is reliable disagreement resolution. Multi-agent debate can amplify errors if agents share the same blind spots, or if the aggregation mechanism is poorly calibrated. Evidence-based protocols that require citations to tool outputs, trace-anchored critiques, and structured voting or verification may reduce correlated failures [53, 58, 64]. Governance is also a systems problem: in production, multi-agent graphs need observability, audit logs, and incident response paths, especially when tools have side effects [4]. Formalizing what it means for a multi-agent system to be “aligned” under delegation (including human-in-the-loop escalation) remains an important research frontier [5, 37].

7.6 Toward Unified Conceptual Frameworks

As agent systems proliferate, clearer conceptual frameworks and taxonomies are needed to compare architectures, identify failure modes, and guide design decisions [30, 48, 49, 66]. A practical goal is to unify the vocabulary across communities: what counts as an “agent” versus “agentic workflow”, how to separate policy models from orchestration, and how to represent environments, tools, and memory in a way that supports evaluation and governance [33]. One direction is to standardize *agent interfaces*: structured tool schemas, trace formats, and evaluation harnesses that allow apples-to-apples comparisons across implementations. Without interface-level standardization, progress risks being dominated by prompt idiosyncrasies and unreported engineering details [29, 44, 64]. Another direction is to connect system taxonomies to learning mechanisms: which architectural choices are best improved by finetuning (instruction/policy tuning), which by preference optimization, and which by system-level changes such as stronger verifiers or better caching. This encourages a more

scientific design loop rather than ad-hoc prompt engineering [11, 37, 46]. Overall, consolidating these perspectives into actionable design models and evaluation checklists is an important step toward a mature engineering discipline, and future surveys should treat systems, learning, and evaluation as a coupled stack rather than independent topics [30, 33].

8 Conclusion

AI agents—systems that embed foundation models in a control loop with memory, tools, and verifiers—are rapidly shifting language models from passive responders to active workflow executors across software, web interaction, multimodal assistance, and embodied domains [20, 36, 64, 67]. This survey synthesized the agent landscape through a unified paradigm and taxonomy: agents as *budgeted* systems with structured tool interfaces and trace-first operation, where reliability and governance are properties of the full stack (model + orchestration + tools) rather than the base model alone [21, 33, 48].

We reviewed learning and optimization across three layers. At the mechanism level, RL/IL, in-context learning, and test-time compute (reflection, self-consistency, and search) provide complementary ways to improve long-horizon behavior under uncertainty [47, 53, 55, 58, 59, 63]. At the system level, modular architectures (policy core, memory, tool routers, planners, critics/verifiers) and infrastructure (schemas, sandboxing, audit logs) constrain side effects and make behavior inspectable and debuggable [5, 21, 50, 64]. At the model level, instruction/policy tuning and preference optimization shape tool-use discipline and safety behavior, and trace-centric finetuning is increasingly used to teach agents to operate rather than merely answer [11, 37, 46].

Because agents interact with non-deterministic environments and toolchains, we emphasized evaluation as a multi-dimensional measurement problem. Beyond end-to-end task success, deployable systems require reporting efficiency/cost, tool-use correctness, trajectory quality, robustness under perturbations, and safety/compliance. Benchmarks and suites such as AgentBench, ToolBench, WebArena, SWE-bench, and GAIA provide complementary stress tests for these dimensions [20, 29, 32, 44, 67].

Looking forward, the central research challenge is to make agent autonomy dependable at scale: verifiable and policy-compliant tool execution, secure and consistent long-term memory, principled allocation of test-time compute under explicit budgets, and trace-first observability that supports auditing, reproducibility, and governance. Progress on these fronts will help close the gap between impressive demonstrations and robust real-world deployment.

References

- [1] Joshua Achiam, David Held, Aviv Tamar, and Pieter Abbeel. 2017. Constrained Policy Optimization. *arXiv preprint arXiv:1705.10528* (2017).
- [2] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Justin Fu, Keerthana Gopalakrishnan, Karol Hausman, et al. 2022. Do As I Can, Not As I Say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691* (2022).
- [3] Jean-Baptiste Alayrac, Jeffrey Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, et al. 2022. Flamingo: A visual language model for few-shot learning. *arXiv preprint arXiv:2204.14198* (2022).
- [4] Anonymous. 2025. Aegis: Taxonomy and Optimizations for Overcoming Agent-Environment Failures in LLM Agents. *arXiv preprint arXiv:2508.19504* (2025).
- [5] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Bryan McCann, et al. 2022. Constitutional AI: Harmlessness from AI feedback. *arXiv preprint arXiv:2212.08073* (2022).
- [6] Anton Bakhtin, Noam Brown, Emily Dinan, Nikhil Jain, Alexey Kirillov, Vladlen Koltun, Mike Liu, Jekaterina Novikova, Long Ouyang, Devi Parikh, et al. 2022. Human-level play in the game of Diplomacy by combining language models with strategic reasoning. *Science* 378, 6624 (2022), 1067–1074.

- [7] Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemyslaw Dere, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Christopher Hesse, et al. 2019. Dota 2 with large scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680* (2019).
- [8] Anthony Brohan, Yevgen Chebotar, Chelsea Finn, et al. 2023. RT-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818* (2023).
- [9] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language Models are Few-Shot Learners. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- [10] S. Chowa et al. 2025. From Language to Action: A Review of Large Language Models as Autonomous Agents and Tool Users. *arXiv preprint arXiv:2508.17281* (2025). <https://arxiv.org/abs/2508.17281>
- [11] Paul F. Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- [12] Michele Colledanchise and Petter Ögren. 2018. *Behavior Trees in Robotics and AI: An Introduction*. CRC Press.
- [13] Peter Dayan and Geoffrey E. Hinton. 1993. Feudal reinforcement learning. *Advances in Neural Information Processing Systems (NeurIPS)* (1993).
- [14] Xiang Deng et al. 2023. Mind2Web: Towards a Generalist Agent for the Web. *arXiv preprint arXiv:2306.06135* (2023). <https://arxiv.org/abs/2306.06135>
- [15] Danny Driess, Fei Xia, Keerthana Gopalakrishnan, Brian Ichter, Michael S Ryoo, Kuan Wong, Rohan Vikas, Maja Neumann, et al. 2023. PaLM-E: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378* (2023).
- [16] Scott Fujimoto, David Meger, and Doina Precup. 2019. Off-Policy Deep Reinforcement Learning without Exploration. *arXiv preprint arXiv:1812.02900* (2019).
- [17] Javier García and Fernando Fernández. 2015. A comprehensive survey on safe reinforcement learning. *Journal of Machine Learning Research* 16, 42 (2015), 1437–1480.
- [18] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. 2018. Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor. *arXiv preprint arXiv:1801.01290* (2018).
- [19] Jonathan Ho and Stefano Ermon. 2016. Generative Adversarial Imitation Learning. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- [20] Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, et al. 2023. SWE-bench: Can language models resolve real-world GitHub issues? *arXiv preprint arXiv:2310.06770* (2023).
- [21] Eyal Karpas, Maxim Likhatchev, Clare Voss, Tim Finin, Mona Singh, et al. 2022. MRKL systems: A modular, neuro-symbolic architecture that combines large language models, external knowledge sources and discrete reasoning. *arXiv preprint arXiv:2205.00445* (2022).
- [22] Takeshi Kojima, Shixiang Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large Language Models are Zero-Shot Reasoners. *arXiv preprint arXiv:2205.11916* (2022).
- [23] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. 2020. Conservative Q-Learning for Offline Reinforcement Learning. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- [24] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive NLP tasks. In *Advances in Neural Information Processing Systems*, Vol. 33. 9459–9474.
- [25] Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023. CAMEL: Communicative agents for “mind” exploration of large language model society. *arXiv preprint arXiv:2303.17760* (2023).
- [26] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2023. BLIP-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597* (2023).
- [27] Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, et al. 2022. Holistic Evaluation of Language Models. *arXiv preprint arXiv:2211.09110* (2022). <https://arxiv.org/abs/2211.09110>
- [28] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. *arXiv preprint arXiv:2304.08485* (2023).
- [29] Ziniu Liu et al. 2023. AgentBench: Evaluating LLMs as agents. *arXiv preprint arXiv:2308.03688* (2023).
- [30] Junyu Luo et al. 2025. Large Language Model Agent: A Survey on Methodology, Applications and Challenges. *arXiv preprint arXiv:2503.21460* (2025).
- [31] Sam Masterman et al. 2024. The Landscape of Emerging AI Agent Architectures for Reasoning, Planning, and Tool Calling: A Survey. *arXiv preprint arXiv:2404.11584* (2024). <https://arxiv.org/abs/2404.11584>
- [32] Grégoire Mialon et al. 2023. GAIA: A Benchmark for General AI Assistants. *arXiv preprint arXiv:2311.12983* (2023). <https://arxiv.org/abs/2311.12983>
- [33] Microsoft Research. 2024. Position Paper: Agent AI Towards a Holistic Intelligence. https://www.microsoft.com/en-us/research/wp-content/uploads/2024/02/Agent_AI_position.pdf.

- [34] Sewon Min, Mike Lewis, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2022. Rethinking the Role of Demonstrations: What Makes In-Context Learning Work? *arXiv preprint arXiv:2202.12837* (2022).
- [35] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin Riedmiller, Andreas K. Fiedjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dhharshan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis. 2015. Human-level control through deep reinforcement learning. *Nature* 518, 7540 (2015), 529–533.
- [36] OpenAI. 2023. GPT-4 Technical Report. *arXiv preprint arXiv:2303.08774* (2023).
- [37] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155* (2022).
- [38] Charles Packer et al. 2023. MemGPT: Towards LLMs as Operating Systems. *arXiv preprint arXiv:2310.08560* (2023). <https://arxiv.org/abs/2310.08560>
- [39] Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative Agents: Interactive simulacra of human behavior. *arXiv preprint arXiv:2304.03442* (2023).
- [40] Mrinal Patil et al. 2023. Gorilla: Large Language Model Connected with Massive APIs. *arXiv preprint arXiv:2305.15334* (2023). <https://arxiv.org/abs/2305.15334>
- [41] Tatiana Petrova et al. 2025. From Semantic Web and MAS to Agentic AI: A Unified Narrative of the Web of Agents. *arXiv preprint arXiv:2507.10644* (2025).
- [42] Dean A. Pomerleau. 1991. Efficient training of artificial neural networks for autonomous navigation. In *Neural Computation*.
- [43] Martin L. Puterman. 1994. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons.
- [44] Yujia Qin et al. 2023. ToolBench: Towards open-source benchmark for tool-augmented large language models. *arXiv preprint arXiv:2309.03752* (2023).
- [45] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*.
- [46] Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. 2023. Direct Preference Optimization: Your Language Model is Secretly a Reward Model. *arXiv preprint arXiv:2305.18290* (2023).
- [47] Stéphane Ross, Geoffrey J. Gordon, and J. Andrew Bagnell. 2011. A reduction of imitation learning and structured prediction to no-regret online learning. In *Proceedings of the 14th International Conference on Artificial Intelligence and Statistics (AISTATS)*.
- [48] Jitao Sang et al. 2025. Beyond Pipelines: A Survey of the Paradigm Shift toward Model-Native Agentic AI. *arXiv preprint arXiv:2510.16720* (2025).
- [49] Ranjan Sapkota, Konstantinos I. Roumeliotis, and Manoj Karkee. 2025. AI Agents vs. Agentic AI: A Conceptual Taxonomy, Applications and Challenges. *arXiv preprint arXiv:2505.10468* (2025).
- [50] Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. *arXiv preprint arXiv:2302.04761* (2023).
- [51] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. 2015. Trust Region Policy Optimization. In *Proceedings of the 32nd International Conference on Machine Learning (ICML)*.
- [52] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal Policy Optimization Algorithms. *arXiv preprint arXiv:1707.06347* (2017).
- [53] Noah Shinn, Beck Labash, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language agents with verbal reinforcement learning. *arXiv preprint arXiv:2303.11366* (2023).
- [54] Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, et al. 2021. ALFWorld: Aligning text and embodied environments for interactive learning. *arXiv preprint arXiv:2010.03768* (2021).
- [55] Richard S. Sutton and Andrew G. Barto. 2018. *Reinforcement Learning: An Introduction* (2 ed.). MIT Press.
- [56] Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. 2019. Grandmaster level in StarCraft II using multi-agent reinforcement learning. *Nature* 575, 7782 (2019), 350–354.
- [57] Guanzhi Wang, Yunzhu Pan, Shiyu Zhang, Xinxin He, Li Wang, Zhiyuan Wang, et al. 2023. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291* (2023).
- [58] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, and Denny Zhou. 2022. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171* (2022).
- [59] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. 2022. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. *arXiv preprint arXiv:2201.11903* (2022).
- [60] Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Li, Erkang Zhu, Yifei He, Siddhartha Arora, et al. 2023. AutoGen: Enabling next-gen LLM applications via multi-agent conversation. *arXiv preprint arXiv:2308.08155* (2023).

- [61] John Yang et al. 2024. SWE-agent: Agent-Computer Interfaces Enable Automated Software Engineering. *arXiv preprint arXiv:2405.15793* (2024). <https://arxiv.org/abs/2405.15793>
- [62] Shunyu Yao et al. 2022. WebShop: Towards Scalable Real-World Web Interaction with Grounded Language Agents. *arXiv preprint arXiv:2207.01206* (2022). <https://arxiv.org/abs/2207.01206>
- [63] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023. Tree of Thoughts: Deliberate problem solving with large language models. *arXiv preprint arXiv:2305.10601* (2023).
- [64] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023. ReAct: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629* (2023).
- [65] A. Yehudai et al. 2025. Survey on Evaluation of LLM-based Agents. *arXiv preprint arXiv:2503.16416* (2025). <https://arxiv.org/abs/2503.16416>
- [66] Jingwen Zhou et al. 2024. A Taxonomy of Architecture Options for Foundation Model-based Agents: Analysis and Decision Model. *arXiv preprint arXiv:2408.02920* (2024).
- [67] Shuyan Zhou, Frank Xu, Zhiqiang Shen, Danqi Chen, et al. 2023. WebArena: A realistic web environment for building autonomous agents. *arXiv preprint arXiv:2307.13854* (2023).