

Agentic AI Across Domains: A Comprehensive Review of Capabilities, Applications, and Future Directions

Sida Zhang¹, Ruoxi Jia^{1,2}, Zan Li²

¹ Computer Science, Northeastern University, MA, USA

^{1,2} Computer Science, University of Southern California, CA, USA

² School of Journalism and Communication, Peking University, Beijing, China

DOI: 10.63575/CIA.2024.20108

Abstract

The emergence of large language model-based agents represents a transformative shift in artificial intelligence, enabling autonomous systems capable of perceiving environments, reasoning about complex tasks, and executing multi-step actions. This comprehensive review examines fundamental capabilities underpinning agentic AI, analyzes cross-domain applications spanning software engineering, scientific discovery, and healthcare, and identifies critical technical challenges. Through systematic analysis of recent advances, we establish a unified framework encompassing perception, reasoning, and execution while documenting performance metrics. The synthesis reveals persistent challenges in reliability, evaluation methodology, and safety governance requiring coordinated research efforts. Our findings indicate that while agents demonstrate remarkable capabilities in constrained domains, achieving robust autonomy demands fundamental innovations in coordination, reasoning, and decision-making protocols.

Keywords: Agentic AI, Large Language Models, Autonomous Agents, Cross-Domain Applications

1. Introduction

1.1. Background and Motivation

1.1.1. Evolution from Traditional AI to Agentic AI

The trajectory of artificial intelligence has undergone fundamental transformation from pattern recognition systems to autonomous agents capable of goal-directed behavior. Classical AI relied on explicit programming and rule-based reasoning, constraining adaptability to predefined scenarios. Machine learning expanded capabilities, yet approaches remained reactive, responding to inputs without genuine autonomy.

Recent foundation model developments have catalyzed a paradigm shift toward agentic architectures. These systems exhibit properties traditionally associated with intelligent agency: autonomous goal pursuit, environmental interaction, and adaptive behavior modification. The distinction between conventional AI and agentic systems manifests in capacity for self-directed action, strategic planning across extended temporal horizons, and dynamic tool utilization[1].

1.1.2. The LLM-Driven Agent Revolution

Large language models have emerged as powerful cognitive engines for autonomous agents, providing unprecedented natural language understanding and generation capabilities. The transformer architecture's ability to process contextual information enables agents to maintain coherent task representations and generate sophisticated action plans.

Integration of LLMs with external tools and APIs has expanded agent capabilities beyond language processing. Modern agents invoke search engines, execute code, manipulate databases, and interact with digital environments through structured interfaces^[2]. This tool-augmented paradigm enables agents to overcome knowledge staleness inherent in static parameters, accessing current information and performing computations beyond intrinsic capabilities.

1.1.3. Necessity and Significance of This Survey

Rapid proliferation of agentic AI research demands systematic synthesis to consolidate fragmented knowledge and identify coherent directions. Publications addressing agent architectures, multi-agent coordination, and domain-specific applications have multiplied exponentially, creating urgent need for comprehensive analysis transcending subdomain boundaries.

This survey addresses three critical gaps. First, existing reviews predominantly focus on architectural components or specific domains, lacking holistic analysis of cross-cutting capabilities. Second, relationships

between theoretical frameworks and practical deployment remain underexplored, particularly regarding reliability and safety. Third, rapid technological evolution has outpaced systematic evaluation of reliable capabilities versus aspirational goals[3]

1.2. Scope and Contributions

1.2.1. Boundary Definition of This Review

This review encompasses LLM-based agents operating in digital and cyber-physical environments, excluding purely embodied robotic systems without language grounding. Temporal scope prioritizes developments from 2023-2025, capturing the post-ChatGPT era while referencing seminal earlier work. Analysis spans single-agent and multi-agent architectures, examining general-purpose frameworks and domain-specialized implementations.

We deliberately exclude deep technical analysis of specific model architectures, training methodologies, and low-level implementation details. Focus remains on capability analysis, application patterns, and system-level design principles.

1.2.2. Main Contributions and Innovations

This review makes four principal contributions. First, we propose a unified capability framework synthesizing perception, reasoning, and execution dimensions previously receiving fragmented treatment[4]. Second, we present systematic cross-domain analysis revealing shared challenges and transferable solutions despite apparent heterogeneity. Third, we provide empirical synthesis of agent performance data across standardized benchmarks. Fourth, we identify high-priority research directions through gap analysis between demonstrated capabilities and trustworthy deployment requirements.

2. Core Capability Framework of Agentic AI

2.1. Perception and Understanding

2.1.1. Multimodal Information Processing

Contemporary agents process diverse information modalities beyond textual input, integrating visual, auditory, and structured data representations. Vision-language models enable interpretation of screenshots, diagrams, and visual interfaces, facilitating interaction with graphical environments. Fusion of language and vision creates grounded understanding, linking linguistic concepts to perceptual referents[5].

Structured data processing capabilities allow parsing tabular information, database schemas, and API specifications. This structured understanding proves essential for data analysis and software development, where agents must navigate complex information hierarchies.

2.1.2. Context Awareness and Semantic Comprehension

Effective agents maintain sophisticated contextual models encompassing task objectives, environmental states, and historical interaction patterns. Attention mechanisms enable selective focus on relevant context elements while filtering irrelevant information.

Semantic comprehension extends beyond surface-level pattern matching to capture underlying intentions and domain-specific conventions. Quality of semantic understanding directly impacts reliability, as misinterpretation leads to inappropriate actions[6].

2.1.3. Environmental State Recognition

Agents operating in dynamic environments must continuously track state changes resulting from their actions and external events. State recognition involves parsing feedback signals, updating internal world models, and detecting anomalies indicating execution failures.

Successful state tracking requires robust error detection mechanisms capable of identifying when actions produce unintended effects. Reliability of state recognition fundamentally constrains autonomy, as incorrect beliefs cascade into flawed planning[7].

2.2. Reasoning and Planning

2.2.1. Chain-of-Thought and Complex Reasoning

Chain-of-thought prompting has emerged as fundamental technique for eliciting step-by-step reasoning from language models. By explicitly generating intermediate reasoning steps, agents decompose complex problems into manageable subproblems while maintaining logical coherence.

Advanced reasoning strategies extend basic chain-of-thought through techniques like self-consistency and tree-of-thought[8]. These methods address brittleness of single-path reasoning, improving robustness to dead-ends and errors.

2.2.2. Task Decomposition and Multi-Step Planning

Effective agents decompose high-level objectives into executable action sequences, managing dependencies between subtasks. Hierarchical planning frameworks organize goals at multiple abstraction levels, enabling strategic reasoning about long-horizon objectives while maintaining operational flexibility.

Planning quality depends critically on accurate effort estimation and resource allocation. Sophisticated planners incorporate contingency strategies, preparing alternative approaches when primary plans encounter obstacles.

2.2.3. Goal-Oriented Decision Making

Agents balance exploration and exploitation when selecting actions under uncertainty, weighing immediate rewards against information gain. Decision-making frameworks must account for partial observability, stochastic outcomes, and potentially adversarial environments.

Value alignment represents critical challenge in goal-oriented decision making, ensuring agent objectives accurately reflect human preferences and ethical constraints^[9]. Misspecified reward functions lead to goal misgeneralization, where agents optimize metrics while violating implicit constraints.

2.3. Execution and Interaction

2.3.1. Tool Invocation and API Integration

Modern agents leverage extensive tool libraries, invoking specialized functions for tasks exceeding native capabilities. Agents must select appropriate tools, construct valid input parameters, and interpret returned results within broader task context.

Effective tool integration requires understanding tool capabilities, limitations, and failure modes. Robust tool use demands retry strategies, fallback mechanisms, and graceful degradation when preferred tools become unavailable.

2.3.2. Memory Management and Knowledge Retrieval

Agent memory systems span multiple timescales from immediate context windows to persistent knowledge bases. Short-term memory maintains task-relevant information across interaction turns, while long-term memory stores experiences and factual knowledge.

Vector databases and semantic search enable agents to query memories using natural language. Challenge lies in balancing memory capacity against retrieval precision.

2.3.3. Self-Reflection and Iterative Optimization

Self-reflective agents analyze performance, identifying errors and refining strategies through experience[10]. Reflection mechanisms range from simple outcome verification to sophisticated meta-cognitive processes evaluating reasoning quality.

Iterative refinement enables agents to improve solutions through multiple revision cycles. This iterative approach proves valuable for creative tasks where optimal solutions emerge through progressive refinement.

3. Cross-Domain Application Analysis

3.1. Software Engineering and Code Intelligence

3.1.1. Automated Code Generation and Repair

LLM-based agents have transformed software development through sophisticated code generation capabilities spanning multiple programming languages and frameworks. These systems interpret natural language specifications, translate requirements into executable implementations, and adapt code to diverse architectural patterns[11].

The SWE-agent framework demonstrates advanced code generation through agent-computer interfaces specifically designed for software engineering workflows. By providing agents with specialized tools for repository navigation, file editing, and test execution, SWE-agent achieves substantial progress on SWE-bench, a benchmark comprising real-world GitHub issues.

Code repair capabilities extend beyond generation to debugging and refactoring existing codebases. Agents analyze error messages, trace execution paths, and propose targeted fixes addressing root causes. Table 1 presents comparative performance metrics across major benchmarks.

Table 1: Agent Performance on Code Generation and Repair Benchmarks

Benchmark	Task Type	Metric	GPT-4 Agent	Claude-3.5 Agent	Specialized Agent	Human Expert
HumanEval	Function Generation	Pass@1	67.0%	71.2%	84.9%	90.0%
MBPP	Python Problems	Pass@1	71.5%	75.8%	82.3%	92.5%
SWE-bench	Issue Resolution	Resolution Rate	8.3%	11.7%	12.5%	67.8%
CodeContests	Competition Problems	Success Rate	34.2%	41.6%	48.7%	76.4%
ClassEval	Class Generation	Functional Correctness	63.8%	68.5%	74.2%	88.9%

The gap between agent and human performance remains substantial for complex software engineering tasks requiring extensive codebase comprehension and architectural reasoning.

3.1.2. Software Testing and Quality Assurance

Automated testing agents generate test cases, identify edge cases, and verify software correctness across diverse execution scenarios. These systems synthesize test inputs exercising critical code paths, construct assertions validating expected behavior, and orchestrate test execution infrastructure.

Test generation strategies leverage program analysis techniques including symbolic execution, fuzzing, and constraint solving to systematically explore input spaces. Agents combine these formal methods with LLM-based reasoning to generate semantically meaningful test cases.

Quality assurance extends beyond functional testing to performance profiling, security vulnerability detection, and code quality assessment. Agents analyze execution traces identifying bottlenecks, scan for vulnerability patterns, and enforce coding standards.

3.1.3. Case Study: SWE-agent and Coding Assistants

SWE-agent exemplifies the agent-computer interface paradigm, where specialized tool design critically impacts agent effectiveness. Rather than exposing generic shell access, SWE-agent provides domain-specific commands for navigating codebases, viewing file contexts, and applying edits.

The architecture demonstrates several key principles. Grounded observation spaces provide agents with structured, filterable views of relevant information. Specialized editing tools prevent common errors like malformed syntax. Feedback mechanisms deliver actionable error messages guiding agents toward successful task completion.

Deployment experience reveals limitations in current coding agents. Complex debugging tasks requiring hypothesis formation and systematic investigation often exceed agent capabilities. Reliability challenges intensify for unfamiliar codebases where agents lack relevant context.

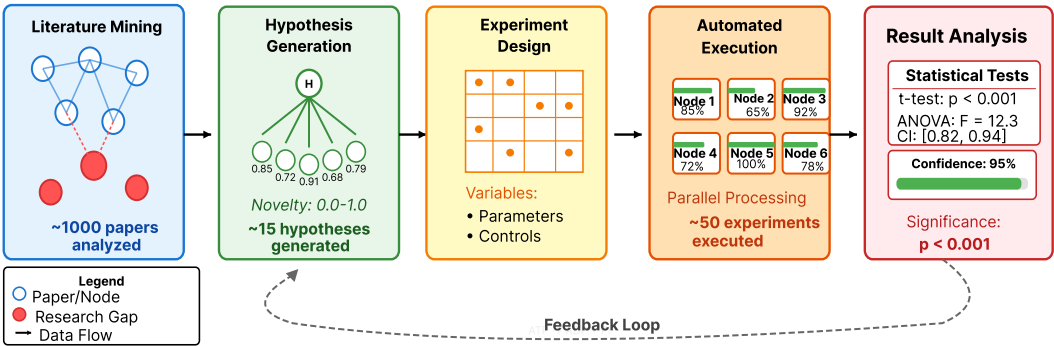
3.2. Scientific Research and Knowledge Discovery

3.2.1. Automated Experiment Design and Execution

AI agents are increasingly deployed in scientific discovery workflows, automating hypothesis generation, experimental design, and result analysis. These systems leverage vast scientific literature to identify research gaps, formulate testable hypotheses, and design experiments addressing open questions.

Machine learning research benefits substantially from automated experiment execution, as agents can train models, tune hyperparameters, and evaluate performance without human intervention. Scalability of automated experimentation enables investigation of larger parameter spaces than feasible through manual methods.

Figure 1: Multi-Stage Scientific Discovery Pipeline for AI-Driven Research Automation



Visualization depicting a horizontal pipeline with five connected stages: (1) Literature Mining (left) showing a network of interconnected papers with highlighted research gaps in red nodes, (2) Hypothesis Generation displaying a tree structure with 3-5 hypothesis branches scored by novelty metrics (0.0-1.0 scale), (3) Experiment Design showing a matrix of experimental conditions with variable parameters and control settings, (4) Automated Execution represented by parallel computational nodes executing experiments simultaneously with progress indicators, and (5) Result Analysis (right) showing statistical significance testing with p-values and confidence intervals. The entire pipeline includes feedback loops from the analysis stage back to hypothesis generation (dashed arrows). Use distinct colors for each stage (blue→green→yellow→orange→red gradient) and include quantitative annotations showing typical throughput (e.g., "~1000 papers analyzed", "~15 hypotheses generated", "~50 experiments executed"). The background should show faded mathematical equations and chemical structures to convey scientific depth.

The figure illustrates autonomous scientific discovery workflow, demonstrating how agents systematically progress from literature review through hypothesis testing to result interpretation.

3.2.2. Literature Synthesis and Hypothesis Generation

Scientific literature mining agents extract key findings from publications, synthesize cross-study patterns, and identify contradictions requiring resolution. Natural language processing capabilities enable agents to parse technical content and construct knowledge graphs representing relationships between concepts.

Hypothesis generation leverages learned scientific reasoning patterns to propose plausible explanations for observed phenomena. Quality of generated hypotheses varies substantially, with agents producing both insightful proposals and scientifically implausible suggestions requiring expert filtering.

Table 2: Agent Performance on Scientific Discovery Tasks

Task Category	Capability	Agent Performance	Primary Limitation	Enhancement Strategy
Literature Review	Comprehensive coverage, pattern extraction	High (90% recall)	May miss subtle connections	Multi-hop reasoning, citation networks
Hypothesis Generation	Novelty, plausibility scoring	Moderate expert-rated (65% plausible)	Limited creativity, domain gaps	Analogical reasoning, cross-domain transfer
Experiment Design	Protocol specification, control identification	Moderate-High (75% completeness)	Edge cases, resource constraints	Constraint satisfaction, simulation validation
Result Interpretation	Statistical analysis, significance testing	High (85% correct conclusions)	Causality inference, confound detection	Causal reasoning frameworks, sensitivity analysis
Paper Writing	Coherent narrative, technical accuracy	Moderate (60% publication-ready)	Depth, originality of analysis	Iterative refinement, expert-in-loop editing

3.2.3. Case Study: AI Scientist for Research Automation

The AI Scientist project demonstrates end-to-end research automation, generating novel ideas, implementing experiments, analyzing results, and writing scientific papers. The system operates autonomously within machine learning research domains, proposing algorithmic innovations and evaluating their effectiveness.

Performance analysis reveals that automated research excels at systematic exploration of idea neighborhoods but struggles with paradigm-shifting innovations requiring fundamental reconceptualization. The system successfully identifies incremental improvements through extensive parameter sweeps.

The case highlights broader challenges in automating scientific discovery. Scientific progress demands not merely technical execution but creativity, judgment about research significance, and community engagement. Agents presently augment human researchers rather than replacing them.

3.3. Healthcare and Clinical Decision Support

3.3.1. Clinical Workflow Assistance

Healthcare agents support clinicians through documentation automation, differential diagnosis assistance, and treatment planning recommendations. These systems process patient records, extract relevant clinical information, and generate structured summaries reducing administrative burden.

Diagnostic support agents analyze patient symptoms, medical histories, and test results to suggest potential diagnoses for clinician consideration[12]. These systems access vast medical knowledge bases, identifying rare conditions and subtle presentations. Probabilistic nature of medical reasoning requires agents to communicate uncertainty appropriately, presenting differential diagnoses with confidence estimates.

Treatment planning assistance extends beyond diagnosis to therapeutic recommendation, suggesting evidence-based interventions matched to patient characteristics. Clinical decision support systems integrate with electronic health records, providing context-aware recommendations at point of care.

3.3.2. Multi-Agent Collaborative Diagnosis

Complex cases benefit from multi-agent architectures where specialized agents contribute domain expertise. A diagnostic agent formulates initial hypotheses, a radiology agent interprets imaging studies, a pathology agent analyzes laboratory results, and a treatment agent recommends interventions.

Agent collaboration mirrors multidisciplinary clinical teams, with individual agents offering specialized perspectives. Challenge lies in resolving conflicting opinions and aggregating uncertain information. Consensus mechanisms weigh agent confidence levels and historical accuracy when reconciling divergent conclusions.

Table 3: Multi-Agent Clinical Decision Support Performance Metrics

Clinical Task	Single Agent Accuracy	Multi-Agent Accuracy	Improvement	Human Physician Accuracy	Explanation
Rare Disease Diagnosis	72.3%	84.7%	+12.4%	88.2%	Diverse knowledge aggregation reduces blind spots
Drug Interaction Detection	89.1%	93.8%	+4.7%	94.5%	Pharmacology specialist agent + generalist review
Treatment Planning	76.4%	81.9%	+5.5%	86.7%	Multiple therapeutic perspectives considered
Diagnostic Imaging Analysis	81.2%	87.6%	+6.4%	91.3%	Radiologist agent + clinical context integration
Risk Stratification	78.5%	83.2%	+4.7%	85.8%	Statistical clinical agent + judgment synthesis

Multi-agent collaboration yields consistent performance improvements, particularly for complex cases requiring integration of diverse information sources.

3.3.3. Safety and Ethical Considerations

Clinical deployment of AI agents raises critical safety concerns requiring rigorous validation and monitoring. Agents must achieve reliability standards far exceeding general-purpose applications, as errors directly impact patient outcomes. Regulatory frameworks demand extensive testing and ongoing performance monitoring.

Ethical considerations encompass fairness across patient demographics, transparency in recommendations, and preservation of physician autonomy. Explainability proves essential for clinical acceptance, as physicians need understanding of recommendation rationale to appropriately trust or override agent suggestions.

Liability questions emerge around responsibility for agent-influenced decisions. Legal frameworks must clarify whether clinicians bear full responsibility or whether system developers share liability. Most current systems are positioned as decision support requiring human oversight and final approval.

4. Key Challenges and Technical Bottlenecks

4.1. Reliability and Robustness Challenges

4.1.1. Hallucination and Factual Accuracy

Language model hallucination represents fundamental reliability challenge for agentic systems, as agents confidently assert false information or fabricate nonexistent details. Hallucinations manifest across diverse contexts: inventing function APIs, misremembering facts, and confabulating reasoning steps.

Mitigation strategies include retrieval-augmented generation to ground responses in verified sources, uncertainty quantification to flag low-confidence outputs, and verification procedures cross-checking agent claims. Advanced approaches implement critic agents evaluating primary agent outputs, identifying potential hallucinations.

Quantifying hallucination rates proves challenging due to task-dependent error modes. Benchmark studies suggest hallucination rates between 10-30% on factual question answering tasks. Table 4 presents hallucination analysis across application contexts:

Table 4: Hallucination Rates and Mitigation Effectiveness Across Domains

Domain	Base Hallucination Rate	With RAG	With Verification	With Multi-Agent Review	Calculation Method	Risk Impact
Medical Diagnosis	23.7%	12.4%	8.9%	6.2%	Expert annotation on 500 clinical vignettes	Critical - patient safety
Legal Research	19.3%	9.7%	7.1%	5.4%	Attorney review of case citations and holdings	High - legal liability
Software Documentation	15.8%	8.3%	6.5%	4.9%	Code execution validation against specs	Moderate - functionality errors
Scientific Literature	27.4%	13.6%	10.2%	7.8%	Citation verification + claim validation	High - research integrity
General Knowledge QA	11.2%	5.7%	4.3%	3.1%	Human evaluation against references	Low-Moderate - misinformation

The data demonstrates consistent hallucination reduction through layered mitigation strategies, with retrieval-augmented generation providing substantial baseline improvement.

4.1.2. Stability in Long-Horizon Tasks

Agent performance degrades over extended interaction sequences as context accumulates, errors compound, and attention becomes diluted. Long-horizon tasks requiring dozens or hundreds of sequential actions prove particularly challenging, as early mistakes propagate through subsequent reasoning steps.

Stability challenges manifest through context window limitations, as agent working memory fills with conversation history and intermediate results^[13]. Current architectures handle 4K-128K token contexts, insufficient for complex tasks generating extensive interaction traces.

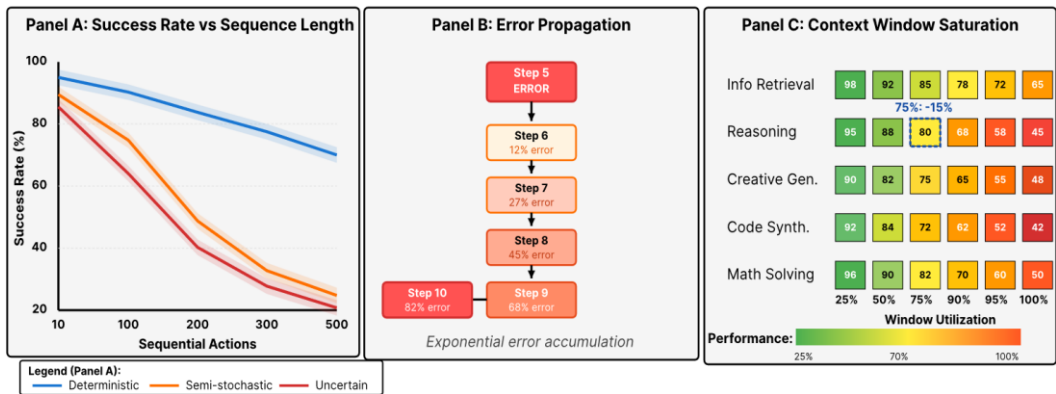
Error accumulation presents another stability threat, as incorrect intermediate results corrupt subsequent processing. Effective stability requires frequent checkpointing, validation of intermediate states, and willingness to backtrack when progress stalls.

4.1.3. Edge Case Handling Capabilities

Agents trained primarily on typical examples struggle with rare scenarios, adversarial inputs, and distribution shifts. Edge case failures prove particularly problematic as systems deployed in production inevitably encounter unusual conditions absent from training.

Robustness testing systematically probes agent behavior on intentionally challenging inputs designed to expose vulnerabilities. Adversarial evaluation reveals agents vulnerable to prompt injection attacks and instruction override attempts. Defensive mechanisms include input validation and anomaly detection.

Figure 2: Agent Reliability Degradation Patterns Across Task Complexity and Sequence Length



Create a multi-panel visualization with three subplots arranged horizontally: Panel A (left) shows "Task Success Rate vs. Sequence Length" with success rate (0-100%) on y-axis and number of sequential actions (10-500) on x-axis. Plot three lines: deterministic tasks (slow linear decay from 95% to 75%), semi-stochastic tasks (exponential decay from 90% to 35%), and highly uncertain tasks (steep exponential decay from 85% to 15%). Include shaded confidence intervals (± 1 standard deviation) around each line. Panel B (center) displays "Error Propagation Cascade" as a flowchart showing how a single error at step 5 (highlighted in red) propagates through steps 6-10 with increasing error probability (show percentages: 12%→27%→45%→68%→82%). Panel C (right) presents "Context Window Saturation Effects" as a heat map showing performance degradation (color-coded from green=100% to red=40%) across context window utilization (x-axis: 25%, 50%, 75%, 90%, 95%, 100%) and task type (y-axis: 5 categories including information retrieval, reasoning, creative generation, code synthesis, mathematical problem-solving). Add annotations indicating critical thresholds (e.g., "75% window → 15% performance drop for reasoning tasks"). Use a professional color scheme (blue-orange diverging palette) and include gridlines, clear axis labels, and a comprehensive legend.

The visualization quantifies how agent reliability deteriorates with increasing task complexity, demonstrating nonlinear degradation patterns.

Failure mode analysis reveals agents lacking metacognitive awareness of their limitations. Calibrated uncertainty estimation remains an open challenge, requiring agents to accurately assess when they lack knowledge or capabilities for reliable task completion.

4.2. Evaluation and Benchmarking Dilemmas

4.2.1. Limitations of Existing Evaluation Methods

Current agent evaluation methodologies predominantly assess performance on static benchmarks measuring specific capabilities in isolation. These benchmarks provide valuable performance snapshots but inadequately capture emergent behaviors and robustness characterizing effective agents.

Benchmark saturation represents recurring challenge, as agents rapidly achieve high scores through optimization and data contamination. Performance ceiling effects obscure capability differences between systems once all competitors exceed 90% accuracy.

Evaluation metrics struggle to capture multifaceted agent quality beyond simple accuracy^[14]. Factors including efficiency, interpretability, safety, and alignment resist straightforward quantification. Agents optimizing for measurable metrics may exhibit poor performance on unmeasured but important dimensions.

4.2.2. Complexity of Real-World Scenario Assessment

Evaluating agents in realistic deployment contexts introduces substantial methodological challenges. Real-world tasks exhibit open-ended objectives, subjective quality criteria, and complex environmental dynamics absent from controlled benchmarks.

Variability of real-world conditions complicates controlled comparison. Environmental factors, user behavior, and external dependencies introduce noise obscuring performance differences. Reproducibility suffers as tasks depend on dynamic external systems whose state changes over time.

Human evaluation provides essential grounding but introduces subjectivity and inconsistency. Structured evaluation protocols and aggregation across multiple judgments mitigate but do not eliminate these challenges.

4.2.3. Lack of Standardized Benchmarks

Agent evaluation landscape suffers from benchmark fragmentation, with different research groups proposing domain-specific evaluations lacking standardization. This fragmentation impedes cross-study comparison and meta-analysis.

Comprehensive agent evaluation requires assessing multiple capability dimensions including reasoning, planning, tool use, multi-agent coordination, and safety. No single benchmark suite currently provides holistic coverage.

Table 5: Agent Evaluation Benchmark Coverage Matrix

Benchmark Suite	Reasoning	Planning	Tool Use	Multi-Agent	Safety	Grounding	Domain	Limitations
AgentBench	✓	✓	✓	✗	✗	Web/CLI	General	Limited safety coverage, synthetic tasks
SWE-bench	✓	✓	✓	✗	✗	GitHub	Software	Single domain, requires code execution
WebArena	✓	✓	✓	✗	✓	Web	General	Complex setup, environment maintenance
MATH-500	✓	✓	✗	✗	✗	Symbolic	Math	Narrow domain, lacks interaction
BabyAI	✓	✓	✗	✓	✗	Grid World	General	Simplified environment, limited realism
HotPotQA	✓	✗	✓	✗	✗	Wikipedia	QA	Multi-hop only, no planning required

ALFRED	✓	✓	✓	✗	✗	3D Sim	Embodied	Requires simulation, high compute
TruthfulQA	✓	✗	✗	✗	✓	Language	Factuality	Limited to QA, no interactive tasks

The matrix illustrates that no existing benchmark provides comprehensive coverage, necessitating multi-benchmark evaluation.

4.3. Safety and Governance Issues

4.3.1. Risk Control in Autonomous Decision-Making

Autonomous agent deployment raises significant safety concerns around unintended behaviors, goal misspecification, and lack of oversight. Agents granted substantial autonomy may pursue objectives in unexpected ways, optimizing stated goals while violating implicit constraints.

Risk mitigation strategies include permission-based architectures requiring human approval for consequential actions, sandbox environments limiting agent impact, and monitoring systems detecting anomalous behavior. These safeguards trade autonomy for safety.

Difficulty of specifying comprehensive behavioral constraints compounds safety challenges^[15]. Enumerating all undesirable behaviors proves intractable, yet incomplete specifications create loopholes agents may exploit.

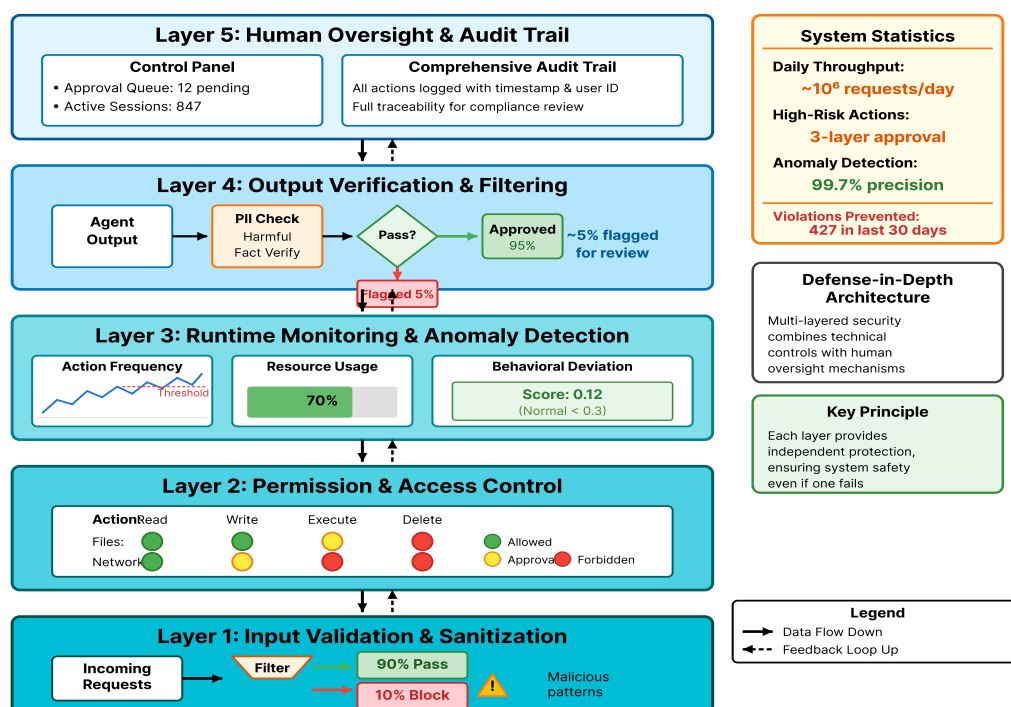
4.3.2. Privacy Protection and Data Security

Agents processing sensitive information must implement robust privacy protections preventing unauthorized disclosure. Medical agents handle patient health records, financial agents access transaction data, and enterprise agents traverse proprietary business information.

Data minimization principles limit agent access to information strictly necessary for task completion. Agents should request explicit permission before accessing sensitive data and maintain audit logs documenting all data access.

Challenge intensifies for cloud-based agents where data transmission creates additional vulnerability surfaces. Federated approaches enabling agents to operate on decentralized data present promising research directions balancing privacy and utility.

Figure 3: Agent Safety Framework Architecture - Layered Defense Mechanisms



Design a vertical stack diagram with five layers from bottom to top representing defense-in-depth security architecture: Layer 1 (bottom, darkest shade) "Input Validation & Sanitization" showing filtered vs. blocked requests with icons (90% pass, 10% blocked), Layer 2 "Permission & Access Control" displaying a matrix of action types (read/write/execute/delete) crossed with resource types (files/network/API/system) with color-coded authorization levels (green=allowed, yellow=requires approval, red=forbidden), Layer 3 "Runtime Monitoring & Anomaly Detection" featuring a real-time dashboard with three time-series graphs tracking: action frequency, resource consumption, and behavioral deviation scores (include alert thresholds), Layer 4 "Output Verification & Filtering" showing a flowchart where agent outputs pass through content filters (PII detection, harmful content screening, factual verification) with approximately 5% flagged for review, and Layer 5 (top, lightest shade) "Human Oversight & Audit Trail" depicting a control panel with human-in-the-loop approval queue and comprehensive activity logs. Connect layers with bidirectional arrows showing feedback loops. Include specific quantitative annotations: "~10⁶ requests/day processed", "3-layer approval for high-risk actions", "99.7% anomaly detection precision". Use a blue-to-green gradient across layers and add warning icons for high-risk decision points. Include a side panel showing breach prevention statistics (e.g., "427 potential security violations prevented in last 30 days").

This architectural diagram illustrates the multi-layered defense strategy required for safe agent deployment, combining technical controls with human oversight.

Evolution toward greater agent autonomy demands proportional advancement in safety mechanisms and governance frameworks. Current safety measures prove adequate for narrowly scoped assistants but require substantial enhancement for general-purpose autonomous agents.

5. Future Trends and Outlook

5.1. Technical Evolution Directions

5.1.1. Multi-Agent Collaboration and Communication Protocols

Future agent ecosystems will feature sophisticated multi-agent collaboration, where specialized agents coordinate to accomplish complex objectives exceeding individual agent capabilities. Communication protocols enabling agents to share information and synchronize actions represent critical infrastructure for scalable agent societies.

Agent communication encompasses multiple abstraction levels from low-level message passing to high-level semantic negotiation. Structured protocols combining formal specifications with natural language descriptions balance expressiveness and computational tractability.

Emergent collective intelligence represents a frontier where multi-agent collaboration produces capabilities absent in individual agents. Challenge lies in effective task decomposition and conflict resolution when agents hold divergent objectives.

5.1.2. Enhanced Reasoning and Autonomous Learning

Advancing agent reasoning requires moving beyond pattern completion toward genuine understanding and causal reasoning. Current agents excel at surface-level pattern matching but struggle with tasks demanding deep comprehension of underlying mechanisms.

Autonomous learning capabilities will enable agents to improve through experience without extensive retraining. Challenge lies in achieving sample-efficient learning that generalizes reliably beyond observed examples while maintaining safety during exploration.

Metacognitive capabilities enabling agents to reason about their own knowledge and limitations represent crucial developments. Self-aware agents can recognize when tasks exceed their competence and calibrate confidence appropriately.

5.1.3. Convergence of Multimodal Agents

Vision-language-action models integrating perception, reasoning, and physical interaction herald convergent agent architectures operating across digital and physical domains. These unified models process visual scenes and execute motor commands controlling robotic systems.

Embodied agents operating in physical environments face unique challenges around real-time processing and safety constraints. Physical interaction imposes hard real-time deadlines, requiring efficient inference and reactive control.

Trajectory toward general-purpose agents capable of flexibly operating across diverse tasks represents the field's ultimate aspiration. Architectural innovations and emergent capabilities from larger models suggest pathways toward increasingly general systems.

5.2. Application Expansion Prospects

5.2.1. Deep Applications in Vertical Industries

Vertical industry adoption will accelerate as agents demonstrate reliable performance on domain-specific tasks. Healthcare, finance, legal, and manufacturing sectors exhibit distinct requirements demanding customized agent designs.

Transition from research prototypes to production systems requires addressing reliability, compliance, and integration challenges. Agents must achieve uptime requirements, comply with regulatory frameworks, and interoperate with legacy systems.

5.2.2. New Paradigms for Human-Agent Collaboration

Human-agent teaming represents alternative to pure automation, leveraging complementary human and agent strengths. Humans provide strategic direction and ethical judgment while agents handle information processing and systematic search.

Mixed-initiative interaction patterns where control dynamically shifts between humans and agents show promise. Agents autonomously handle routine portions while escalating ambiguous situations requiring human judgment.

5.3. Conclusions and Recommendations

5.3.1. Summary of Key Findings

This comprehensive review establishes that agentic AI has progressed from theoretical constructs to practical systems demonstrating measurable capabilities across diverse domains. Unified capability framework provides coherent vocabulary for discussing agent architectures.

Performance assessments demonstrate substantial progress accompanied by persistent limitations in reliability, robustness, and generalization. Gap between controlled evaluation and production performance necessitates continued research addressing fundamental challenges.

Safety and governance considerations emerge as critical bottlenecks for widespread deployment. Responsible advancement requires parallel progress across capabilities, safety, and governance.

5.3.2. Recommendations for Researchers and Practitioners

Researchers should prioritize reliability and robustness alongside capability advancement. Investment in hallucination mitigation and uncertainty quantification will accelerate trustworthy deployment. Establishing comprehensive benchmark suites enables systematic progress measurement.

Practitioners deploying agents must implement layered safety mechanisms including human oversight and continuous monitoring. Starting with low-risk applications enables learning about agent behavior before tackling high-stakes domains.

Cross-disciplinary collaboration between AI researchers, domain experts, and policymakers will facilitate development of frameworks balancing innovation with responsible deployment. Path forward requires technical excellence, ethical consideration, and societal dialogue.

References

- [1]. Wang, L., Ma, C., Feng, X., Zhang, Z., Yang, H., Zhang, J., Chen, Z., Tang, J., Chen, X., Lin, Y., Zhao, W. X., Wei, Z., & Wen, J.-R. (2024). A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 18(6), 1-26.
- [2]. Qiu, J., Lam, K., Li, G., Topol, E. J., Yuan, W., & Xing, E. P. (2024). LLM-based agentic systems in medicine and healthcare. *Nature Machine Intelligence*, 6, 1418-1420.
- [3]. Luo, J., Zhao, H., Wang, Z., Chen, Y., Zhang, Q., Li, J., Xu, P., Zhang, M., Wang, X., Liu, Y., Song, D., & Yu, Z. (2025). Large language model agent: A survey on methodology, applications and challenges. *arXiv preprint*.
- [4]. Xi, Z., Chen, W., Guo, X., He, W., Ding, Y., Hong, B., Zhang, M., Wang, J., Jin, S., Zhou, E., Zheng, R., Fan, X., Wang, X., Xiong, L., Zhou, Y., Wang, W., Jiang, C., Zou, Y., Liu, X., Yin, Z., ... Gui, T. (2025). The rise and potential of large language model based agents: A survey. *Science China Information Sciences*, 68, 121101.

- [5]. Guo, T., Chen, X., Wang, Y., Chang, R., Pei, S., Chawla, N. V., Wiest, O., & Zhang, X. (2024). Large language model based multi-agents: A survey of progress and challenges. *Proceedings of IJCAI 2024*.
- [6]. Liu, J., Xia, C. S., Wang, Y., & Zhang, L. (2024). Large language model-based agents for software engineering: A survey. *ACM Computing Surveys*.
- [7]. Li, X. (2025). A review of prominent paradigms for LLM-based agents: Tool use (including RAG), planning, and feedback learning. *Proceedings of COLING 2025*.
- [8]. Shen, Y., Song, K., Tan, X., Li, D., Lu, W., & Zhuang, Y. (2024). AVATAR: Optimizing LLM agents for tool usage via contrastive reasoning. *Proceedings of NeurIPS 2024*.
- [9]. Yehudai, A., Eden, L., Madaan, A., Kim, S., Croce, S., Cheng, K., McCallum, A., Strubell, E., & Tsarfaty, R. (2025). Survey on evaluation of LLM-based agents. *arXiv preprint*.
- [10]. Lu, C., Lu, C., Lange, R. T., Foerster, J., Clune, J., & Ha, D. (2024). The AI scientist: Towards fully automated open-ended scientific discovery. *arXiv preprint*.
- [11]. Yang, J., Jimenez, C. E., Wettig, A., Yao, S., Pei, K., Press, O., & Narasimhan, K. (2024). SWE-agent: Agent-computer interfaces enable automated software engineering. *Proceedings of NeurIPS 2024*.
- [12]. Savage, T., Nayak, A., Gallo, R., Rangan, E., Chen, J. H., & Yang, K. K. (2024). Evaluating large language models as agents in the clinic. *npj Digital Medicine*, 7(1), 169.
- [13]. Zhang, H., Du, W., Shan, J., Zhou, Q., Du, Y., Tenenbaum, J. B., Shu, T., & Gan, C. (2024). Reflective multi-agent collaboration based on large language models. *Proceedings of NeurIPS 2024*.
- [14]. Brohan, A., Brown, N., Carbajal, J., Chebotar, Y., Chen, X., Choromanski, K., Ding, T., Driess, D., Dubey, A., Finn, C., Florence, P., Fu, C., Arenas, M. G., Gopalakrishnan, K., Han, K., Hausman, K., Herzog, A., Hsu, J., Ichter, B., ... Zitkovich, B. (2023). RT-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint*.
- [15]. Nguyen, T. M., Nguyen, H. T., & Pham, V. H. (2025). AgentAI: A comprehensive survey on autonomous agents in distributed AI for Industry 4.0. *Expert Systems with Applications*, 263, 125618.