*Original Article*

# Standardizing AI Agent Protocols: Definitions, Layered Architectures, and Future Research Directions

Saurabh Atri
(W3C AI Agent Protocol Community Group Member).

*Abstract -* *AI agents are moving from isolated task executors to networked, autonomous participants in software ecosystems. As a current member of the W3C AI Agent Protocol Community Group, I draw on ongoing standards discussions and emerging best practices to frame the interoperability and governance problem. As agents increasingly discover resources, invoke tools, negotiate tasks, and coordinate with humans and other agents, ad hoc interfaces fail to provide interoperability, security, and predictable behavior at scale. This article defines AI Agent Protocols as composable specifications spanning transport, message semantics, capability representation, policy enforcement, and governance. We synthesize recent surveys of agent protocols, standards efforts including the W3C AI Agent Protocol Community Group, and emerging interoperability layers such as the Model Context Protocol to propose a reference architecture and taxonomy. We analyze design tradeoffs, threat models, and conformance criteria, and identify research directions: semantic interoperability, verifiable identity, policy-carrying messages, adaptive coordination, and safety-by-construction for large-scale multi-agent systems.*

*Keywords - AI Agents, Agent Protocols, Multi-Agent Systems, Interoperability, Security, Governance, Agentic Web.*

## 1. Introduction

Large language models (LLMs) and tool-augmented systems have accelerated the adoption of autonomous and semi-autonomous agents across customer support, software engineering, data analysis, robotics, and operations. At the same time, the ecosystem is fragmenting into incompatible stacks: agents differ in how they discover tools, represent intent, authenticate peers, and coordinate tasks. This fragmentation is a scaling limiter: interoperability, safety, and governance become integration problems rather than system properties. Standards organizations have begun to treat an "Agentic Web" as a next-layer evolution of web interactions, requiring protocols for agent discovery, identity, and collaboration [1]. Academic work similarly identifies protocol standardization as a prerequisite for robust agent ecosystems and proposes taxonomies and evaluation dimensions (e.g., security, scalability, latency) [2].

## 2. Definition and Scope of AI Agent Protocols

We define an AI Agent Protocol as a formal specification that governs: (i) how an agent represents and exchanges messages (syntax and semantics), (ii) how it discovers and invokes capabilities (tools, services, and resources), (iii) how it coordinates with other agents and humans, and (iv) how it enforces security, policy, and accountability constraints. Protocols are not synonymous with agent "frameworks" (runtime libraries) or "architectures" (internal agent design). Protocols are implementation-agnostic contracts. Their purpose is to make heterogeneous agents composable, testable, and governable.

## 3. Reference Architecture

Agent protocols are best understood as a layered stack. While implementations may collapse layers, separating them clarifies responsibilities and enables independent standardization.

**Table 1: Reference Layer Model for AI Agent Protocols**

| Layer | Primary Concern | Typical Artifacts | Illustrative Standards/Specs |
|---|---|---|---|
| **Transport** | Connectivity and framing | Streams, request/response, SSE(**Server-Sent Events**), stdio | JSON-RPC 2.0 transport in MCP [3] |
| **Messaging** | Envelope + routing | Message IDs, correlation, retries | RPC schemas; event-driven patterns |
| **Semantics** | Meaning of intent and outcomes | Action/goal schema; tool contracts | JSON-LD for linked semantics [7] |
| **Capability** | Tool/resource discovery & invocation | Tool registry; permissions; I/O contracts | MCP servers/clients/hosts [3][4] |

| **Identity & Trust** | Who is speaking, and under what authority | Identifiers, credentials, attestations | W3C DIDs (Decentralized Identifiers) [5], Verifiable Credentials [6] |
|---|---|---|---|
| **Policy & Governance** | Constraints, auditability, risk controls | Policies, logs, safety gates | NIST AI RMF guidance [8]; domain compliance |

This layered model emphasizes that agent interoperability is not achieved by a single protocol. Instead, interoperable ecosystems compose transport, semantics, identity, and policy mechanisms. Conformance tests should therefore be defined per layer and for cross-layer profiles (e.g., "enterprise agent profile" vs. "open web agent profile").

# 4. Core Components of Agent Protocols

## 4.1. Communication and Message Semantics

Communication protocols define message schemas, sequencing, and error behavior. For agents, message semantics must encode not just data, but intent, commitments, and expected postconditions. A practical design pattern is a split between: (i) machine-actionable intent (tool calls, constraints, required approvals) and (ii) human-readable rationale for audit and oversight. Semantic interoperability is the primary differentiator from conventional distributed systems. JSON-LD enables messages to carry machine-interpretable meaning via linked data contexts, supporting schema evolution and cross-vendor interpretation [7].

## 4.2. Capability Discovery and Tool Invocation

Many modern agents rely on external tools and data sources. The Model Context Protocol (MCP) standardizes how an AI host connects to servers that expose tools and context, using JSON-RPC 2.0 messages and explicit roles (host, client, server) [3]. The design goal is to make integrations composable and reduce bespoke glue code, which improves portability and reduces the attack surface from ad hoc connectors [4].

## 4.3. Coordination, Negotiation, and Multi-Agent Interaction

Multi-agent protocols regulate task allocation, negotiation, consensus, and conflict resolution. Classic multi-agent interaction models (e.g., auction-based allocation, contract-net style delegation) provide useful primitives, but LLM-driven agents require additional constraints: explicit coordination state, bounded autonomy, and resource-aware commitments. Historically, the multi-agent systems community developed explicit speech-act based communication languages; for example, FIPA specified communicative acts and message-level semantics to support interoperable agent interactions [9][10]. While modern LLM-agent protocols emphasize tool use and web-native transports, FIPA-style semantics remain relevant as a conceptual baseline for commitments, delegation, and coordination. A key practical requirement is "coordination determinism": protocols should define what must be deterministic (e.g., ordering rules, commit logs) even if internal reasoning is stochastic. Without this, debugging emergent failures becomes infeasible at scale.

## 4.4. Security, Identity, and Accountability

Open agent ecosystems require robust identity and trust. W3C Decentralized Identifiers (DIDs) define a standards-based mechanism for verifiable decentralized identity [5]. W3C Verifiable Credentials (VCs) provide tamper-evident, machine-verifiable claims (issuer-holder-verifier model) [6]. Together, DIDs and VCs support agent identity, delegation, and provenance without depending on a single centralized identity provider. Security protocols must also address tool misuse (unauthorized actions), prompt/tool injection, data exfiltration, and cross-agent spoofing. Protocol-level mitigations include: capability-scoped tokens, least-privilege tool grants, signed messages, replay protection, and mandatory audit logs.

## 4.5. Governance and Risk Management

Governance protocols operationalize safety and compliance. The NIST AI Risk Management Framework (AI RMF 1.0) provides a widely used structure for incorporating trustworthiness considerations into design, development, and evaluation [8]. Agent protocols can map governance controls to protocol artifacts: for example, attaching risk classifications to actions, requiring approvals for high-impact categories, and logging evidence for post-incident analysis.

# 5. Taxonomy and Classification

Protocol taxonomies help evaluate tradeoffs and guide system design. A recent survey of agent protocols proposes a two-dimensional classification separating (i) context-oriented protocols (agent-tool / agent-data) vs. inter-agent protocols, and (ii) general-purpose vs. domain-specific protocols [2]. This framing is useful because context access and inter-agent coordination have different threat models, latency constraints, and conformance requirements. We extend this taxonomy with an additional axis: governance intensity. Protocols used in regulated environments (healthcare, finance) require stronger auditability and identity guarantees than protocols optimized for experimentation.

# 6. Design Challenges and Failure Modes

Key challenges include: (1) schema and semantics drift across vendors; (2) emergent collective behavior not anticipated by protocol designers; (3) security vulnerabilities from tool exposure; (4) balancing autonomy with human oversight; and (5) performance constraints (latency, bandwidth) in large agent swarms. The protocol surface area grows with capability richness, so minimal, composable cores with profile-based extensions are preferable. From an engineering standpoint, the most common failure modes are not model failures but protocol failures: ambiguous schemas, missing idempotency guarantees, lack of correlation identifiers,

and insufficient authentication at boundaries. Treating protocols as first-class artifacts (with test suites, conformance claims, and version negotiation) is essential.

## 7. Future Directions

Research and standardization priorities include: (i) semantic interoperability via shared ontologies and linked-data contexts; (ii) cryptographic provenance for actions and evidence; (iii) policy-carrying messages where constraints travel with requests; (iv) adaptive coordination protocols that can change strategies under load; and (v) safety-by-construction approaches such as formally specified action permissions and runtime monitors. Being a member of the W3C AI Agent Protocol Community Group, we explicitly target open mechanisms for discovery, identity, and cross-agent collaboration on the Web [1]. Bridging such work with practical tool protocols (e.g., MCP) and governance frameworks (e.g., NIST AI RMF) is a plausible path toward interoperable, auditable agent ecosystems.

## 8. Conclusion

AI Agent Protocols are a prerequisite for scalable autonomy. By separating transport, semantics, capability access, identity, and governance, protocol designers can produce composable standards and verifiable implementations. Near-term progress depends on aligning emerging standards efforts with rigorous conformance testing and threat-informed design. Long-term success will be measured by whether agents can collaborate across vendors while remaining secure, auditable, and aligned with human intent.

## References

[1] W3C, "AI Agent Protocol Community Group," W3C Community Group page, accessed January 5, 2026. https://www.w3.org/community/agentprotocol/

[2] Y. Yang et al., "A Survey of AI Agent Protocols," arXiv:2504.16736, April 2025. https://arxiv.org/abs/2504.16736

[3] Model Context Protocol, "Specification (2025-11-25)," modelcontextprotocol.io, November 25, 2025. https://modelcontextprotocol.io/specification/2025-11-25

[4] Anthropic, "Introducing the Model Context Protocol," November 25, 2024. https://www.anthropic.com/news/model-context-protocol

[5] W3C, "Decentralized Identifiers (DIDs) v1.0," W3C Recommendation, July 19, 2022. https://www.w3.org/TR/did-core/

[6] W3C, "Verifiable Credentials Data Model v2.0," W3C Recommendation, May 15, 2025. https://www.w3.org/TR/vc-data-model-2.0/

[7] W3C, "JSON-LD 1.1," W3C Recommendation, July 16, 2020. https://www.w3.org/TR/json-ld11/

[8] NIST, "Artificial Intelligence Risk Management Framework (AI RMF 1.0)," NIST AI 100-1, January 2023. https://nvlpubs.nist.gov/nistpubs/ai/nist.ai.100-1.pdf

[9] Foundation for Intelligent Physical Agents (FIPA). "FIPA Communicative Act Library Specification."

[10] Document No. XC00037H (Experimental), 2001. Mirror PDF: https://jmvidal.cse.sc.edu/library/XC00037H.pdf

[11] (Accessed January 5, 2026).

[12] FIPA, "FIPA Communicative Act Specifications," FIPA Repository, accessed January 5, 2026. https://www.fipa.org/repository/cas.php3